

# MANUAL

INL/EXT-15-34123

Revision 4

Printed February 2016

## RAVEN User Manual

Cristian Rabiti, Andrea Alfonsi, Joshua Cogliati, Diego Mandelli, Robert Kinoshita,  
Sonat Sen, Congjian Wang

Prepared by  
Idaho National Laboratory  
Idaho Falls, Idaho 83415

The Idaho National Laboratory is a multiprogram laboratory operated by  
Battelle Energy Alliance for the United States Department of Energy  
under DOE Idaho Operations Office. Contract DE-AC07-05ID14517.

Approved for unlimited release.



Issued by the Idaho National Laboratory, operated for the United States Department of Energy by Battelle Energy Alliance.

**NOTICE:** This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government, nor any agency thereof, nor any of their employees, nor any of their contractors, subcontractors, or their employees, make any warranty, express or implied, or assume any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represent that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government, any agency thereof, or any of their contractors or subcontractors. The views and opinions expressed herein do not necessarily state or reflect those of the United States Government, any agency thereof, or any of their contractors.

Printed in the United States of America. This report has been reproduced directly from the best available copy.



INL/EXT-15-34123  
Revision 4  
Printed February 2016

# **RAVEN User Manual**

***Principal Investigator (PI):***

Cristian Rabiti

***Main Developers:***

Andrea Alfonsi

Joshua Cogliati

Diego Mandelli

Congjian Wang

Sonat Sen

Robert Kinoshita

***Contributors:***

Alessandro Bandini

Daniel P. Maljovec

Ivan Rinaldi

Paul W. Talbot



# Contents

1	Introduction .....	13
2	Manual Formats .....	14
3	Installation Overview .....	15
4	RAVEN Dependencies Installation .....	16
4.1	Linux Package Manager .....	16
4.1.1	Ubuntu .....	16
4.1.2	Fedora .....	16
4.2	OSX Dependencies .....	17
4.3	Miniconda .....	17
4.3.1	With MOOSE .....	17
4.3.2	Direct Miniconda .....	17
4.3.3	With Miniconda package .....	18
4.4	Manual Dependency Install .....	18
5	RAVEN Installation .....	20
5.1	Parallel Directory Git .....	20
5.2	Submodule Git .....	22
5.3	RAVEN Whole Devel Branch .....	22
5.4	RAVEN Whole Stable Branch .....	22
5.5	RAVEN Source Code Package .....	23
5.6	RAVEN Compilation .....	23
5.7	RAVEN Complete .....	24
5.8	Troubleshooting the Installation .....	24
6	Running RAVEN .....	25
7	Raven Input Structure .....	26
7.1	Comments .....	26
7.2	Verbosity .....	27
7.3	External Input Files .....	28
8	RunInfo .....	29
8.1	RunInfo: Input of Calculation Flow .....	29
8.2	RunInfo: Input of Queue Modes .....	31
8.3	RunInfo: Example Cluster Usage .....	33
8.4	RunInfo: Advanced Users .....	34
8.5	RunInfo: Examples .....	35
9	Files .....	37
10	Distributions .....	38
10.1	1-Dimensional Probability Distributions .....	38
10.1.1	1-Dimensional Continuous Distributions .....	38
10.1.1.1	Beta Distribution .....	39
10.1.1.2	Exponential Distribution .....	40
10.1.1.3	Gamma Distribution .....	41

10.1.1.4	Logistic Distribution .....	42
10.1.1.5	LogNormal Distribution .....	42
10.1.1.6	Normal Distribution .....	43
10.1.1.7	Triangular Distribution .....	44
10.1.1.8	Uniform Distribution .....	45
10.1.1.9	Weibull Distribution .....	45
10.1.2	1-Dimensional Discrete Distributions. ....	46
10.1.2.1	Bernoulli Distribution .....	47
10.1.2.2	Binomial Distribution .....	47
10.1.2.3	Poisson Distribution .....	48
10.1.2.4	Categorical Distribution .....	49
10.2	N-Dimensional Probability Distributions .....	50
10.2.1	MultivariateNormal Distribution .....	50
10.2.2	NDInverseWeight Distribution .....	53
10.2.3	NDCartesianSpline Distribution .....	54
11	Samplers .....	56
11.1	Forward Samplers .....	59
11.1.1	Monte Carlo .....	59
11.1.2	Grid .....	62
11.1.3	Sparse Grid Collocation .....	66
11.1.4	Sobol .....	70
11.1.5	Stratified .....	72
11.1.6	Response Surface Design .....	78
11.1.7	Factorial Design .....	83
11.2	Dynamic Event Tree (DET) Samplers .....	87
11.2.1	Dynamic Event Tree .....	88
11.2.2	Hybrid Dynamic Event Tree .....	90
11.3	Adaptive Samplers .....	96
11.3.1	Limit Surface Search .....	96
11.3.2	Adaptive Dynamic Event Tree .....	102
11.3.3	Adaptive Hybrid Dynamic Event Tree .....	107
11.3.4	Adaptive Sparse Grid .....	113
11.3.5	Adaptive Sobol Decomposition .....	118
12	DataObjects .....	122
13	Databases .....	127
14	OutStream system .....	129
14.1	Printing system .....	129
14.1.1	<b>DataObjects</b> Printing .....	130
14.1.2	<b>ROM</b> Printing .....	131
14.2	Plotting system .....	132
14.2.1	Plot input structure .....	132
14.2.1.1	“Actions” input block .....	132

14.2.1.2	“plotSettings” input block	138
14.2.1.3	Predefined Plotting System: 2D/3D	141
14.2.2	2D & 3D Scatter plot	142
14.2.3	2D & 3D Line plot	143
14.2.4	2D & 3D Histogram plot	143
14.2.5	2D & 3D Stem plot	145
14.2.6	2D Step plot	146
14.2.7	2D Pseudocolor plot	146
14.2.8	2D Contour or filledContour plots	147
14.2.9	3D Surface Plot	148
14.2.10	3D Wireframe Plot	149
14.2.11	3D Tri-surface Plot	150
14.2.12	3D Contour or filledContour plots	151
14.2.13	DataMining plots	152
14.2.14	Example XML input	153
15	Models	154
15.1	Code	155
15.2	Dummy	157
15.3	ROM	158
15.3.1	NDspline	159
15.3.2	GaussPolynomialRom	160
15.3.3	HDMRRom	162
15.3.4	MSR	164
15.3.5	NDinvDistWeight	167
15.3.6	SciKitLearn	168
15.3.6.1	Linear Models	168
15.3.6.1.1	Linear Model: Automatic Relevance Determination Regression	168
15.3.6.1.2	Linear Model: Bayesian ridge regression	169
15.3.6.1.3	Linear Model: Elastic Net	171
15.3.6.1.4	Linear Model: Elastic Net CV	172
15.3.6.1.5	Linear Model: Least Angle Regression model	173
15.3.6.1.6	Linear Model: Cross-validated Least Angle Regression model	174
15.3.6.1.7	Linear Model trained with L1 prior as regularizer (aka the Lasso)	175
15.3.6.1.8	Lasso linear model with iterative fitting along a regularization path	176
15.3.6.1.9	Lasso model fit with Least Angle Regression	177
15.3.6.1.10	Cross-validated Lasso, using the LARS algorithm	178
15.3.6.1.11	Lasso model fit with Lars using BIC or AIC for model selection	179

15.3.6.1.12	Ordinary least squares Linear Regression . . . . .	180
15.3.6.1.13	Logistic Regression . . . . .	180
15.3.6.1.14	Multi-task Lasso model trained with L1/L2 mixed-norm as regularizer . . . . .	181
15.3.6.1.15	Multi-task Elastic Net model trained with L1/L2 mixed- norm as regularizer . . . . .	182
15.3.6.1.16	Orthogonal Matching Pursuit model (OMP) . . . . .	183
15.3.6.1.17	Cross-validated Orthogonal Matching Pursuit model (OMP) . . . . .	184
15.3.6.1.18	Passive Aggressive Classifier . . . . .	184
15.3.6.1.19	Passive Aggressive Regressor . . . . .	185
15.3.6.1.20	Perceptron . . . . .	186
15.3.6.1.21	Randomized Lasso . . . . .	187
15.3.6.1.22	Randomized Logistic Regression . . . . .	189
15.3.6.1.23	Linear least squares with l2 regularization . . . . .	190
15.3.6.1.24	Classifier using Ridge regression . . . . .	191
15.3.6.1.25	Ridge classifier with built-in cross-validation . . . . .	192
15.3.6.1.26	Ridge regression with built-in cross-validation . . . . .	193
15.3.6.1.27	Linear classifiers (SVM, logistic regression, a.o.) with SGD training . . . . .	194
15.3.6.1.28	Linear model fitted by minimizing a regularized em- pirical loss with SGD . . . . .	196
15.3.6.1.29	Compute Least Angle Regression or Lasso path using LARS algorithm . . . . .	198
15.3.6.1.30	Compute Lasso path with coordinate descent . . . . .	199
15.3.6.1.31	Stability path based on randomized Lasso estimates . .	200
15.3.6.1.32	Gram Orthogonal Matching Pursuit (OMP) . . . . .	201
15.3.6.2	Support Vector Machines . . . . .	201
15.3.6.2.1	Linear Support Vector Classifier . . . . .	202
15.3.6.2.2	C-Support Vector Classification . . . . .	203
15.3.6.2.3	Nu-Support Vector Classification . . . . .	205
15.3.6.2.4	Support Vector Regression . . . . .	206
15.3.6.3	Multi Class . . . . .	208
15.3.6.3.1	One-vs-the-rest (OvR) multiclass/multilabel strategy . .	208
15.3.6.3.2	One-vs-one multiclass strategy . . . . .	208
15.3.6.3.3	Error-Correcting Output-Code multiclass strategy . . . .	209
15.3.6.4	Naive Bayes . . . . .	209
15.3.6.4.1	Gaussian Naive Bayes . . . . .	210
15.3.6.4.2	Multinomial Naive Bayes . . . . .	211
15.3.6.4.3	Bernoulli Naive Bayes . . . . .	212
15.3.6.5	Neighbors . . . . .	212
15.3.6.5.1	Nearest Neighbors . . . . .	213



15.3.6.5.2	K Neighbors Classifier .....	214
15.3.6.5.3	Radius Neighbors Classifier .....	215
15.3.6.5.4	K Neighbors Regressor .....	216
15.3.6.5.5	Radius Neighbors Regressor .....	218
15.3.6.5.6	Nearest Centroid Classifier .....	219
15.3.6.6	Quadratic Discriminant Analysis .....	219
15.3.6.7	Tree .....	220
15.3.6.7.1	Decision Tree Classifier .....	221
15.3.6.7.2	Decision Tree Regressor .....	222
15.3.6.7.3	Extra Tree Classifier .....	223
15.3.6.7.4	Extra Tree Regressor .....	225
15.3.6.8	Gaussian Process .....	226
15.4	External Model .....	229
15.4.1	Method: def readMoreXML .....	230
15.4.2	Method: def initialize .....	231
15.4.3	Method: def createNewInput .....	232
15.4.4	Method: def run .....	233
15.5	PostProcessor .....	234
15.5.1	BasicStatistics .....	235
15.5.2	ComparisonStatistics .....	237
15.5.3	SafestPoint .....	239
15.5.4	LimitSurface .....	241
15.5.5	LimitSurfaceIntegral .....	243
15.5.6	External .....	245
15.5.7	TopologicalDecomposition .....	246
15.5.8	DataMining ( <i>aka: Knowledge Discovery in DataBases</i> ) .....	248
15.5.8.1	Gaussian mixture models .....	249
15.5.8.1.1	GMM classifier .....	249
15.5.8.1.2	Dirichlet Process GMM Classifier (DPGMM) .....	251
15.5.8.1.3	Variational GMM Classifier (VBGMM) .....	252
15.5.8.2	Clustering .....	252
15.5.8.2.1	K-Means Clustering .....	254
15.5.8.2.2	Mini Batch K-Means .....	255
15.5.8.2.3	Affinity Propagation .....	257
15.5.8.2.4	Mean Shift .....	258
15.5.8.2.5	Spectral clustering .....	259
15.5.8.2.6	DBSCAN Clustering .....	261
15.5.8.2.7	Clustering performance evaluation .....	261
15.5.8.3	Decomposing signals in components (matrix factorization problems) .....	262
15.5.8.3.1	Principal component analysis (PCA) .....	262
15.5.8.3.2	Truncated singular value decomposition .....	267

15.5.8.3.3	Fast ICA .....	268
15.5.8.4	Manifold learning .....	269
15.5.8.4.1	Isomap .....	269
15.5.8.4.2	Locally Linear Embedding .....	271
15.5.8.4.3	Spectral Embedding .....	273
15.5.8.4.4	Multi-dimensional Scaling (MDS) .....	274
15.5.9	Interfaced .....	275
15.5.9.1	Data Format .....	276
16	Functions .....	278
17	Steps .....	280
17.1	SingleRun .....	281
17.2	MultiRun .....	283
17.3	IOStep .....	286
17.4	RomTrainer .....	289
17.5	PostProcess .....	290
18	Existing Interfaces .....	292
18.1	Generic Interface .....	292
18.2	RELAP5 Interface .....	294
18.2.1	Sequence .....	294
18.2.2	batchSize and mode .....	295
18.2.3	RunInfo .....	295
18.2.4	Files .....	295
18.2.5	Models .....	296
18.2.6	Distributions .....	296
18.2.7	Samplers .....	297
18.2.8	Steps .....	298
18.2.9	Databases .....	299
18.3	RELAP7 Interface .....	300
18.3.1	Files .....	300
18.3.2	Models .....	300
18.3.3	Distributions .....	301
18.3.4	Samplers .....	301
18.4	MooseBasedApp Interface .....	302
18.4.1	Files .....	302
18.4.2	Models .....	302
18.4.3	Distributions .....	303
18.4.4	Samplers .....	304
18.4.5	Steps .....	305
18.4.6	Databases .....	307
18.4.7	DataObjects .....	307
18.4.8	OutStreamManager .....	308
18.5	MooseVPP Interface .....	309

18.6	OpenModelica Interface	310
18.6.1	Files	311
18.6.2	Models	311
18.6.3	CSV Output	312
18.7	Dymola Interface	313
18.7.1	Files	314
18.7.2	Models	314
18.8	Mesh Generation Coupled Interfaces	315
18.8.1	MooseBasedApp and Cubit Interface	316
18.8.1.1	Files	316
18.8.1.2	Models	317
18.8.1.3	Distributions	318
18.8.1.4	Samplers	318
18.8.1.5	Steps,OutputStreamManager,DataObjects	319
18.8.1.6	File Cleanup	319
18.8.2	MooseBasedApp and Bison Mesh Script Interface	319
18.8.2.1	Files	319
18.8.2.2	Models	320
18.8.2.3	Distributions	321
18.8.2.4	Samplers	321
18.8.2.5	Steps,OutputStreamManager,DataObjects	322
18.8.2.6	File Cleanup	322
19	Advanced Users: How to couple a new code	323
19.1	Pre-requisites.	323
19.2	Code Interface Creation	326
19.2.1	Method: generateCommand	327
19.2.2	Method: createNewInput	328
19.2.3	Method: getInputExtension	329
19.2.4	Method: finalizeCodeOutput	330
19.2.5	Method: checkForOutputFailure	330
19.3	Tools for Developing Code Interfaces	331
19.3.1	File Objects	331

## Appendix

A	Appendix: Example Primer	333
A.1	Example 1.	333
A.2	Example 2.	336
	References	343



# 1 Introduction

RAVEN is a software framework able to perform parametric and stochastic analysis based on the response of complex system codes. The initial development was aimed at providing dynamic risk analysis capabilities to the thermohydraulic code RELAP-7, currently under development at Idaho National Laboratory (INL). Although the initial goal has been fully accomplished, RAVEN is now a multi-purpose stochastic and uncertainty quantification platform, capable of communicating with any system code.

In fact, the provided Application Programming Interfaces (APIs) allow RAVEN to interact with any code as long as all the parameters that need to be perturbed are accessible by input files or via python interfaces. RAVEN is capable of investigating system response and explore input space using various sampling schemes such as Monte Carlo, grid, or Latin hypercube. However, RAVEN strength lies in its system feature discovery capabilities such as: constructing limit surfaces, separating regions of the input space leading to system failure, and using dynamic supervised learning techniques.

The development of RAVEN started in 2012 when, within the Nuclear Energy Advanced Modeling and Simulation (NEAMS) program, the need to provide a modern risk evaluation framework arose. RAVEN's principal assignment is to provide the necessary software and algorithms in order to employ the concepts developed by the Risk Informed Safety Margin Characterization (RISMC) program. RISMC is one of the pathways defined within the Light Water Reactor Sustainability (LWRS) program.

In the RISMC approach, the goal is not just to identify the frequency of an event potentially leading to a system failure, but the proximity (or lack thereof) to key safety-related events. Hence, the approach is interested in identifying and increasing the safety margins related to those events. A safety margin is a numerical value quantifying the probability that a safety metric (e.g. peak pressure in a pipe) is exceeded under certain conditions.

Most of the capabilities, implemented having RELAP-7 as a principal focus, are easily deployable to other system codes. For this reason, several side activates have been employed (e.g. RELAP5-3D, any MOOSE-based App, etc.) or are currently ongoing for coupling RAVEN with several different software. The aim of this document is to detail the input requirements for RAVEN focusing on the input structure.

## 2 Manual Formats

In order to highlight some parts of the Manual having a particular meaning (e.g. input structure, examples, terminal commands, etc.), specific formats have been used. In this sections all the formats with a specific meaning are reported:

- ***Python Coding:***

```
class AClass():
    def aMethodImplementation(self):
        pass
```

- ***XML input example:***

```
<MainXMLBlock>
...
<anXMLnode name='anObjectName' anAttribute='aValue'>
    <aSubNode>body</aSubNode>
</anXMLnode>
...
</MainXMLBlock>
```

- ***Bash Commands:***

```
cd trunk/raven/
./raven_libs_script.sh
cd ../../
```

### 3 Installation Overview

The installation of the RAVEN code is a straightforward procedure; depending on the usage purpose and machine architecture, the installation process slightly differs.

In the following sections, all the different installation procedures are reported.

There are two main requirements to installing RAVEN, installing the RAVEN dependencies (Section 4) and installing RAVEN itself (Section 5). For any particular installation, only one of the raven dependency procedures and one of the raven installation paths needs to be taken. For some versions of OSX, there is a RAVEN complete package that combines these steps (Section 5.7).

For OSX it is recommended that the dependencies be installed with Miniconda (Section 4.3). For Linux, it is recommended that the distribution package manager be used if possible (Section 4.1).

There are several different ways of getting RAVEN as described in Section 5. They vary depending on how easy they are to use and how easy they are to develop with (these tend to be inversely correlated).

## 4 RAVEN Dependencies Installation

There are several different ways to install the dependencies that RAVEN needs. For OSX, Xcode command line tools and XQuartz are needed, and then Miniconda is recommended (Section 4.2). If the MOOSE environment package is installed, Miniconda is recommended (Section 4.3.1). For Linux, the distribution's package manager is recommended (Section 4.1).

### 4.1 Linux Package Manager

If RAVEN is being installed on Linux, it is recommended that the distribution's package manager be used to install the dependencies. Below are instructions for doing so in Ubuntu and Fedora.

#### 4.1.1 Ubuntu

Install the dependencies using the command:

```
sudo apt-get install libtool git python-dev swig g++ python3-dev \
python-numpy python-sklearn python-h5py
```

Optionally, if you want to be able to edit and rebuild the manual, you can install  $\text{\TeX}$  Live and its related packages:

```
sudo apt-get install texlive-latex-base texlive-extra-utils \
texlive-latex-extra texlive-math-extra
```

Now go on to Section 5 for Raven installation.

#### 4.1.2 Fedora

Install the dependencies:

```
dnf install swig libtool gcc-c++ python-devel python3-devel \
numpy h5py scipy python-scikit-learn python-matplotlib-qt4
```

Optionally, if you want to be able to edit and rebuild the manual, you can install  $\text{\TeX}$  Live and its related packages:

```
dnf install texlive texlive-subfigure texlive-stmaryrd \
texlive-titlesec texlive-preprint texlive-placeins
```



Now go on to Section 5 for Raven installation.

## 4.2 OSX Dependencies

For running OSX, there are two additional dependencies that need to be installed. The first is the Xcode command line tools. The second is XQuartz. How to install these is explained at <http://mooseframework.org/getting-started/osx/>

Next Miniconda can be used to install the rest of the dependencies (Section 4.3)

## 4.3 Miniconda

For OSX, the simplest and most robust way to install the RAVEN dependencies is with Miniconda. The details of this vary depending what else is installed. Only one of the following subsections should be used.

### 4.3.1 With MOOSE

If running other MOOSE packages is needed, then MOOSE's MOOSE Environment package comes with miniconda and that should be used to install the RAVEN dependencies.

Then:

```
sudo bash
#Set up http_proxy and https_proxy variables if needed
conda install numpy hdf5 h5py scipy scikit-learn matplotlib swig
exit #exit from root
```

Don't install the RAVEN Miniconda package with the MOOSE Environment package, they conflict.

Now go on to Section 5 for Raven installation.

### 4.3.2 Direct Miniconda

Miniconda can be directly installed using the instructions at: <http://conda.pydata.org/miniconda.html>

Then the packages RAVEN needs can be installed:

```
conda install numpy hdf5 h5py scipy scikit-learn matplotlib swig
```

Now go on to Section 5 for Raven installation.

### 4.3.3 With Miniconda package

For OSX Yosemite and OSX Mavericks there is a package available called `raven_miniconda.dmg`. The package includes both miniconda and the dependencies RAVEN needs. This package should not be used if the MOOSE Environment package is installed.

The files will be installed into `/opt/raven_libs`.

Your `.bash_profile` will be modified to source the `/opt/raven_libs/environments/raven_libs_profile` file. This needs to be sourced. Reopening the terminal will put the correct Miniconda python in the path.

Once this is installed you can go on to Section 5 and install RAVEN.

## 4.4 Manual Dependency Install

If other options don't work, the dependencies can be installed manually. RAVEN uses the following packages (newer versions usually work):

1. `numpy-1.7.0`
2. `hdf5-1.8.12`
3. `Cython-0.18`
4. `h5py-2.2.1`
5. `scipy-0.12.0`
6. `scikit-learn-0.14.1`
7. `matplotlib-1.4.0`
8. `swig-2.0.12`

You may install these dependencies yourself, or by running the `raven_libs_script.sh` script provided within the RAVEN distribution:

```
#Only use raven_libs_script.sh if other methods don't work  
cd full_path_to_raven_distribution/raven/  
./raven_libs_script.sh  
cd ..
```

If the `raven_libs_script.sh` script is successful it will install the raven dependencies in `~/raven_libs`

If it needs to install `swig` or other tools, you may need to update your path. The command below can be used, and to make it run everytime it should be added to `.bash_profile` or `.bashrc` file.

```
export PATH="$HOME/raven_libs/pylibs/bin:$PATH"
```

Now go on to Section 5 for Raven installation.

## 5 RAVEN Installation

Once the RAVEN dependencies have been installed (See Section 4), the rest of RAVEN can be installed.

There are several different ways to get RAVEN. There are trade offs between how easy it is to setup and how easy it is to develop with the method. If you are using several MOOSE applications and want to share the MOOSE directory go to Section 5.1. If you want to do development on RAVEN and RAVEN is the only MOOSE application you are using or you don't mind having duplicate MOOSE directories, go to Section 5.2. If you don't need to do RAVEN development but want to keep up with the latest version of RAVEN, the RAVEN Whole repository can be used as described in Section 5.3. If you don't need to do RAVEN development, and want to track a stable version of RAVEN, the RAVEN Whole repository with a stable branch can be used as described in Section 5.4. If git is not used, either the Source Code Package in Section 5.5 or if in OSX the complete package (Section 5.7) can be used.

### 5.1 Parallel Directory Git

This install method uses git to obtain the software, and the raven, crow, and moose directories are installed at the same level. This is used to allow other MOOSE applications to share the same MOOSE and libmesh.

First follow the MOOSE getting started instructions: <http://mooseframework.org/getting-started/>

This is generally:

```
git clone https://github.com/idaholab/moose.git
cd moose/scripts
./update_and_rebuild_libmesh.sh
cd ../../
```

Next Crow and RAVEN need to be cloned:

```
git clone git@hpcgitlab.inl.gov:idaholab/crow.git
git clone git@hpcgitlab.inl.gov:idaholab/raven.git
```

If RELAP-7 will be used, only IAPWS should be initialized instead of all the submodules:

```
git clone git@hpcgitlab.inl.gov:idaholab/relap-7.git
cd relap-7
git submodule update --init contrib/iapws
```

```
cd ..
```

Next the RAVEN modules should be compiled:

```
cd raven
make framework_modules
```

If RELAP-7 is used, the rest of raven can be compiled:

```
#make sure you are in the raven directory
make
```

Then the testing should be done. If RELAP-7 is not used the framework tests should be run:

```
./run_framework_tests
```

If RELAP-7 is used, all the tests can be run:

```
./run_tests
```

The output should describe why any tests failed.

At the end, there should be a line that looks similar to the output below:

```
8 passed, 19 skipped, 0 pending, 0 failed
```

Normally there are skipped tests because either some of the codes are not available, or some of the test are not currently working. The output will explain why each is skipped.

If all the tests pass, you are ready to read about Running RAVEN in Section 6.

To update the software, the git pull command can be used:

```
cd moose
git pull
#Possibly libmesh needs to be rebuilt at this point.
cd ../crow
git pull
cd ../raven
git pull
```

## 5.2 Submodule Git

This install method uses git to obtain the software and uses submodules to get Crow and MOOSE. This can be used for RAVEN and Crow development.

First RAVEN needs to be cloned, and then the submodules initialized.

```
git clone git@hpcgitlab.inl.gov:idaholab/raven.git
cd raven
git submodule update --init crow moose
```

Next follow the compilation and testing instructions in Section 5.6.

To update the software, the git pull and submodule update commands can be used:

```
git pull
git submodule update
```

## 5.3 RAVEN Whole Devel Branch

This install method uses git to get a custom repository that includes only the parts of Crow and MOOSE that are needed by RAVEN. This repository is updated every time RAVEN or Crow or MOOSE are changed and the RAVEN framework tests pass.

First clone the RAVEN whole repository:

```
git clone git@hpcgitlab.inl.gov:cogljj/raven.git
cd raven
```

Next follow the compilation and testing instructions in Section 5.6.

To update the software, the git pull command can be used:

```
git pull
```

## 5.4 RAVEN Whole Stable Branch

This is the same as in the previous Section 5.3, but after the clone command a stable branch is checked out:

```
git clone git@hpcgitlab.inl.gov:cogljj/raven.git
```

```
cd raven
git checkout stable/2015_july
```

The rest is the same as the RAVEN Whole devel branch. The stable branches are updated less frequently, and input changes that would invalidate existing files are not allowed.

The compilation and testing instructions in Section 5.6 can be followed.

## 5.5 RAVEN Source Code Package

Untar the source install (if there is more than one version of the source tarball, the full filename will need to be used instead of \*):

```
tar -xvzf raven_framework_*_source.tar.gz
cd raven
```

Next follow the compilation and testing instructions in Section 5.6.

## 5.6 RAVEN Compilation

The RAVEN modules should be compiled:

```
#change into the raven directory if needed.
make framework_modules
```

Then the testing should be done:

```
./run_framework_tests
```

The output should describe why any tests failed.

At the end, there should be a line that looks similar to the output below:

```
8 passed, 19 skipped, 0 pending, 0 failed
```

Normally there are skipped tests because either some of the codes are not available, or some of the test are not currently working. The output will explain why each is skipped.

If all the tests pass, you are ready to read about Running RAVEN in Section 6.

If the tests did not pass, check Section 5.8 on troubleshooting.

## 5.7 RAVEN Complete

Open up the file `raven_framework_complete_version.dmg`.

Next, Open up the `raven_libs.pkg` inside, and install it. If you get an error that the package is not signed, then Control click the package, and choose “Open With” and then Installer.

The files will be installed into `/opt/raven_libs`.

Your `.bash_profile` will be modified to source the `/opt/raven_libs/environments/raven_libs_profile` file. This file sets up the environment variables `PYTHONPATH` and `PATH` so that the `raven_framework` command can be used. The `.bash_profile` needs to be sourced. Reopening the terminal will source it and put the executable in the path.

## 5.8 Troubleshooting the Installation

Often the problems result from one or more of the libraries being incorrect or missing. In the raven directory, the command:

```
./run_tests --library_report
```

can be used to check if all the libraries are available, and which ones are being used. If `amsc`, `distribution1D` or `interpolationND` are missing, then the RAVEN modules need to be compiled or recompiled. Otherwise, the RAVEN dependencies need to be fixed.



## 6 Running RAVEN

The RAVEN code is a blend of C++, C, and Python software. The entry point resides on the Python side and is accessible via a command line interface. After following the instructions in the previous Section, RAVEN is ready to be used. The RAVEN driver is contained in the folder “raven/framework.” To run RAVEN, open a terminal and use the following command (replace `<inputFileName.xml>` with your RAVEN input file):

```
python raven/framework/Driver.py <inputFileName.xml>
```

Alternatively, the `raven_framework` script can be used. In this case, the command is:

```
raven_framework <inputFileName.xml>
```

## 7 Raven Input Structure

The RAVEN code does not have a fixed calculation flow, since all of its basic objects can be combined in order to create a user-defined calculation flow. Thus, its input (XML format) is organized in different XML blocks, each with a different functionality. The main input blocks are as follows:

- **<Simulation>**: The root node containing the entire input, all of the following blocks fit inside the *Simulation* block.
- **<RunInfo>**: Specifies the calculation settings (number of parallel simulations, etc.).
- **<Files>**: Specifies the files to be used in the calculation.
- **<Distributions>**: Defines distributions needed for describing parameters, etc.
- **<Samplers>**: Sets up the strategies used for exploring an uncertain domain.
- **<DataObjects>**: Specifies internal data objects used by RAVEN.
- **<Databases>**: Lists the HDF5 databases used as input/output to a RAVEN run.
- **<OutputStreamManager>**: Visualization and Printing system block.
- **<Models>**: Specifies codes, ROMs, post-processing analysis, etc.
- **<Functions>**: Details interfaces to external user-defined functions and modules. the user will be building and/or running.
- **<Steps>**: Combines other blocks to detail a step in the RAVEN workflow including I/O and computations to be performed.

Each of these blocks are explained in dedicated sections in the following chapters.

### 7.1 Comments

Comments may be included in the RAVEN input using standard XML comments, using `<!--` and `-->` as shown in the example below.

```
<Simulation>
...
<!-- An Example Comment -->
<Samplers>
...
```

Comments may be placed anywhere *except* before the `<Simulation>` node or after the `</Simulation>` node. Comments outside the root node will cause errors in maintaining input file compatibility. Additionally, comments must completely surround any nodes they comment out. Comments are intended to completely remove blocks of code, or to add readability. For instance, the following is INCORRECT usage:

```
<!--<Assembler> -->
<!--</Assembler> -->
```

and the following is compatible usage for a code block:

```
<!--<Samplers>
  <Monte Carlo name='mc'>
    ...
  </Monte Carlo>
  ...
</Samplers> -->
```

## 7.2 Verbosity

Each block within RAVEN also makes use of a **verbosity** system, which allows a user to control the level of output to the user interface. These settings are declared globally as attributes in the `<Simulation>` node, and locally in each block. The verbosity levels are

- **'silent'** - Only simulation-breaking errors are displayed.
- **'quiet'** - Errors as well as warnings are displayed.
- **'all'** (default) - Errors, warnings, and messages are displayed.
- **'debug'** - For developers. All errors, warnings, messages, and debug messages are displayed.

Examples of verbosity usage are included in many examples throughout this manual.

At the `<Simulation>` node, global variables can be set, including **verbosity**. In addition, the attribute **printTimeStamps** can be used to either enable or disable prepending RAVEN output with time stamps by setting it to **'true'** or **'false'**.

### 7.3 External Input Files

The `<ExternalXML>` node defines external input file (XML format) that can be used to replace any XML nodes under `<Simulation>` in the RAVEN input file. This node allows a user to load any external input file that contains the required XML nodes into the RAVEN input file. Each `<ExternalXML>` node has the following attributes:

- `node`, *required string attribute*, user-defined XML node of RAVEN input file.
- `xmlToLoad`, *required string attribute*, file name with its absolute or relative path. Note: if a relative path is specified, it must be relative with respect to the RAVEN input file.

For example, if the file `Models.xml` contain the required RAVEN input XML node `<Models>`, the RAVEN input file might appear as:

```
<Simulation>
...
<Steps>
...
</Steps>
...
<ExternalXML node='Models'
    xmlToLoad='external_input/Models.xml' />
...
</Simulation>
```

Another example, if the file `MultiRun.xml` contain the required RAVEN input XML node `<MultiRun>` under node `<Steps>`, the RAVEN input file might appear as:

```
<Simulation>
...
<Steps>
...
    <ExternalXML node='MultiRun'
        xmlToLoad='external_input/MultiRun.xml' />
...
</Steps>
...
</Simulation>
```

## 8 RunInfo

In the **RunInfo** block, the user specifies how the overall computation should be run. This block accepts several input settings that define how to drive the calculation and set up, when needed, particular settings for the machine the code needs to run on (e.g. queueing system, if not PBS, etc.). In the following subsections, we explain all the keywords and how to use them in detail.

### 8.1 RunInfo: Input of Calculation Flow

This sub-section contains the information regarding the XML nodes used to define the settings of the calculation flow that is being performed through RAVEN:

- **<WorkingDir>**, *string, required field*, specifies the absolute or relative (with respect to the location where the xml file is located) path to a directory that will store all the results of the calculations and where RAVEN looks for the files specified in the block **<Files>**. If `runRelative='True'` is used as an attribute, then it will be relative to where raven is run.

*Default: None*

- **<RemoteRunCommand>**, *string, required field*, specifies the absolute or relative (with respect to the framework directory) path to a command that can be used on a remote machine to execute a command. The command is passed in as the environmental variable `COMMAND`.

*Default: raven\_qsub\_command.sh*

- **<batchSize>**, *integer, optional field*, specifies the number of parallel runs executed simultaneously (e.g., the number of driven code instances, e.g. RELAP5-3D, that RAVEN will spawn at the same time). Each parallel run will use `NumThreads * NumMPI` cores.

*Default: 1*

- **<Sequence>**, *comma separated string, required field*, is an ordered list of the step names that RAVEN will run (see Section 17).

- **<JobName>**, *string, optional field*, specifies the name to use for the job when submitting to a pbs queue. Acceptable characters include alphanumeric as well as “-” and “\_”. If more than 15 characters are provided, RAVEN will truncate it using a hyphen between the first 10 and last 4 character, i.e., “1234567890abcdefg” will be truncated to “1234567890-efgh”.

*Default: raven\_qsub*

- **<NumThreads>**, *integer, optional field*, can be used to specify the number of threads RAVEN should associate when running the driven software. For example, if RAVEN is

driving a code named “FOO,” and this code has multi-threading support, this block is used to specify how many threads each instance of FOO should use (e.g. “FOO --n-threads=N” where N is the number of threads).

*Default: 1 (or None when the driven code does not have multi-threading support)*

- **<NumMPI>**, *integer, optional field*, can be used to specify the number of MPI CPUs RAVEN should associate when running the driven software. For example, if RAVEN is driving a code named “FOO,” and this code has MPI support, this block specifies how many MPI CPUs each instance of FOO should use (e.g. “mpexec FOO -np N” where N is the number of CPUs).

In addition, this flag activates the parallel execution of internal Models (e.g. ROMs, ExternalModels, PostProcessors, etc.). If this node is not present, the internal Models are run using a multi-threading approach (i.e. single processor, multiple threads equal to the batch-size)

*Default: 1 (or None when the driven code does not have MPI support)*

- **<totalNumCoresUsed>**, *integer, optional field*, is the global number of CPUs RAVEN is going to use for performing the calculation. When the driven code has MPI and/or multi-threading support and the user specifies NumThreads > 1 and NumMPI > 1, then totalNumCoresUsed is set according to the following formula:

$$\text{totalNumCoresUsed} = \text{NumThreads} * \text{NumMPI} * \text{batchSize}.$$

*Default: 1*

- **<internalParallel>**, *boolean, optional field*, is a boolean flag that controls the type of parallel implementation needs to be used for Internal Objects (e.g. ROMs, External Models, PostProcessors, etc.). If this flag is set to:
  - **False**, the internal parallelism is employed using multi-threading (i.e. 1 processor, multiple threads);
  - **True**, the internal parallelism is employed using an internally-developed multi-processor approach (i.e. **<totalNumCoresUsed>** processors, 1 single thread). This approach works for both Shared Memory Systems (e.g. PC, laptops, workstations, etc.) and Distributed Memory Machines (e.g. High Performance Computing Systems, etc.).

*Default: False*

- **<precommand>**, *string, optional field*, specifies a command that needs to be inserted before the actual command that is used to run the external model (e.g., `mpexec -n 8 precommand ./externalModel.exe (...)`). Note that the precommand as well as the postcommand are ONLY applied to execution commands flagged as “parallel” within the code interface.

*Default: None*

- **<postcommand>**, *string, optional field*, specifies a command that needs to be appended after the actual command that is used to run the external model (e.g., `mpiexec -n 8 ./externalModel.exe (...) postcommand`). Note that the postcommand as well as the precommand are ONLY applied to execution commands flagged as “parallel” within the code interface.  
*Default: None*
- **<clusterParameters>**, *string, optional field*, specifies extra parameters to be used with the cluster submission command. For example, if `qsub` is used to submit a command, then these parameters will be used as extra parameters with the `qsub` command.  
*Default: None*
- **<MaxLogFileSize>**, *integer, optional field*. specifies the maximum size of the log file in bytes. Every time RAVEN drives a code/software, it creates a logfile of the code’s screen output.  
*Default:  $\infty$*   
*( Note: This flag is not implemented yet.)*
- **<deleteOutExtension>**, *comma separated string, optional field*, specifies, if a run of an external model has not failed, which output files should be deleted by their extension (e.g., `<deleteOutExtension>txt,pdf</deleteOutExtension>` will delete all generated txt and pdf files).  
*Default: None*
- **<delSucLogFiles>**, *boolean, optional field*, when True and the run of an external model has not failed (return code = 0), deletes the associated log files.  
*Default: False*

## 8.2 RunInfo: Input of Queue Modes

In this sub-section, all of the keywords (XML nodes) for setting the queue system are reported.

- **<mode>**, *string, optional field*, can specify which kind of protocol the parallel environment should use. RAVEN currently supports one pre-defined “mode”:
  - **mpi**: this “mode” uses `mpiexec` to distribute the running program; more information regarding this protocol can be found in [1]. Mode “MPI” can either generate a `qsub` command or can execute on selected nodes. In order to make the “mpi” mode generate a `qsub` command, an additional keyword (xml sub-node) needs to be specified:
    - \* If RAVEN is executed in the HEAD node of an HPC system using [2], the user needs to input a sub-node, **<runQSUB>**, right after the specification of the mpi mode (i.e.

`<mode>mpi<runQSUB/></mode>`). If the keyword is provided, RAVEN generates a `qsub` command, instantiates itself, and submits itself to the queue system.

- \* If the user decides to execute RAVEN from an “interactive node” (a certain number of nodes that have been reserved in interactive PBS mode), RAVEN, using the “mpi” system, is going to utilize the reserved resources (CPUs and nodes) to distribute the jobs, but, will not generate a `qsub` command.

When the user decides to run in “mpi” mode without making RAVEN generate a `qsub` command, different options are available:

- \* If the user decides to run on the local machine (either in local desktop/workstation or a remote machine), no additional keywords are needed (i.e. `<mode>mpi</mode>`).
- \* If the user is running on multiple nodes, the node ids have to be specified:
  - the node ids can be specified in an external text file (node ids separated by blank space). This file needs to be provided in the XML node `<mode>`, introducing a sub-node named `<nodefile>` (e.g. `<mode>mpi<nodefile>/tmp/nodes</nodefile></mode>`).
  - the node ids can be contained in an environmental variable (node ids separated by blank space). This variable needs to be provided in the `<mode>` XML node, introducing a sub-node named `<nodefileenv>` (e.g. `<mode>mpi<nodefileenv>NODEFILE</nodefileenv></mode>>`).
  - If none of the above options are used, RAVEN will attempt to find the nodes’ information in the environment variable `PBS_NODEFILE`.

In addition, this flag activates the remote (PBS) execution of internal Models (e.g. ROMs, ExternalModels, PostProcessors, etc.). If this node is not present, the internal Models are run using a multi-threading approach (i.e. master processor, multiple parallel threads)

- `<CustomMode>`, *xml node, optional field*, is an xml node where “advanced” users can implement newer “modes.” Please refer to sub-section 8.4 for advanced users.
- `<queueingSoftware>`, *string, optional field*. RAVEN has support for the PBS queueing system. If the platform provides a different queueing system, the user can specify its name here (e.g., PBS PROFESSIONAL, etc.).  
*Default: PBS PROFESSIONAL*
- `<expectedTime>`, *column separated string, optional field (mpi or custom mode)*, specifies the time the whole calculation is expected to last. The syntax of this node is *hours:minutes:seconds* (e.g. 40:10:30 equals 40 hours, 10 minutes, 30 seconds). After this period of time, the HPC system will automatically stop the simulation (even if the simulation is not completed). It is preferable to rationally overestimate the needed time.  
*Default: 10:00:00 (10 hours.)*



### 8.3 RunInfo: Example Cluster Usage

For this example, we have a PBSPro cluster, and there are thousands of node, and each node has 4 processors that share memory. There are a couple different ways this can be used. One way is to use interactive mode and have a RunInfo block:

```
<RunInfo>
  <WorkingDir>./</WorkingDir>
  <Sequence>FirstMRun</Sequence>
  <batchSize>3</batchSize>
  <NumThreads>4</NumThreads>
  <mode>mpi</mode>
  <NumMPI>2</NumMPI>
</RunInfo>
```

Then the commands can be used:

```
#Note: select=NumMPI*batchSize, ncpus=NumThreads
qsub -l select=6:ncpus=4:mpiprocs=1 -l walltime=10:00:00 -I
#wait for processes to be allocated and interactive shell to start

#Switch to the correct directory
cd $PBS_O_WORKDIR

#Load the module with the raven libraries
module load raven-devel-gcc

#Start Raven
python ../../framework/Driver.py test_mpi.xml
```

Alternatively, RAVEN can be asked to submit the qsub directory. With this, the RunInfo is:

```
<RunInfo>
  <WorkingDir>./</WorkingDir>
  <Sequence>FirstMQRun</Sequence>
  <batchSize>3</batchSize>
  <NumThreads>4</NumThreads>
  <mode>
    mpi
    <runQSUB/>
  </mode>
  <NumMPI>2</NumMPI>
  <expectedTime>10:00:00</expectedTime>
```

```
</RunInfo>
```

In this case, the command run from the cluster submit node:

```
python ../../framework/Driver.py test_mpiqsub_local.xml
```

## 8.4 RunInfo: Advanced Users

This sub-section addresses some customizations of the running environment that are possible in RAVEN. Firstly, all the keywords reported in the previous sections can be pre-defined by the user in an auxiliary XML input file. Every time RAVEN gets instantiated (i.e. the code is run), it looks for an optional file, named “default\_runinfo.XML” contained in the “\home\username\.raven\” directory (i.e. “\home\username\.raven\default\_runinfo.XML”). This file (same syntax as the RunInfo block defined in the general input file) will be used for defining default values for the data in the RunInfo block. In addition to the keywords defined in the previous sections, in the **<RunInfo>** node, an additional keyword can be defined:

- **<DefaultInputFile>**, *string, optional field*. In this block, the user can change the default xml input file RAVEN is going to look for if none have been provided as a command-line argument.  
*Default: “test.xml”.*

As already mentioned, this file is read to define default data for the RunInfo block. This means that all the keywords defined here will be overridden by any values specified in the actual RAVEN input file.

In section 8.2, it is explained how RAVEN can handle the queue and parallel systems. If the currently available “modes” are not suitable for the user’s system (workstation, HPC system, etc.), it is possible to define a custom “mode” modifying the **<RunInfo>** block as follows:

```
<RunInfo>
...
  <CustomMode file="newMode.py" class="NewMode">
    aNewMode
  </CustomMode>
  <mode>aNewMode</mode>
...
</RunInfo>
```

The file field can use `%BASE_WORKING_DIR%` and `%FRAMEWORK_DIR%` to specify the location of the file with respect to the base working directory or the framework directory.

The python file should define a class that inherits from `Simulation.SimulationMode` of the RAVEN framework and overrides the necessary functions. Generally, `modifySimulation` will be overridden to change the precommand or postcommand parts which will be added before and after the executable command. An example Python class is given below with the functions that can and should be overridden:

```
import Simulation
class NewMode(Simulation.SimulationMode):
    def doOverrideRun(self):
        # If doOverrideRun is true, then use runOverride instead of
        # running the simulation normally.
        # This method should call simulation.run somehow
        return True

    def runOverride(self):
        # this can completely override the Simulation's run method
        pass

    def modifySimulation(self):
        # modifySimulation is called after the runInfoDict has been
        # setup and allows the mode to change any parameters that
        # need changing. This typically modifies the precommand and
        # the postcommand that are put before/after the command.
        pass

    def XMLread(self, XMLNode):
        # XMLread is called with the mode node, and can be used to
        # get extra parameters needed for the simulation mode.
        pass
```

## 8.5 RunInfo: Examples

Here we present a few examples using different components of the RunInfo node:

```
<RunInfo>
  <WorkingDir>externalModel</WorkingDir>
  <Sequence>MonteCarlo</Sequence>
  <batchSize>100</batchSize>
  <NumThreads>4</NumThreads>
  <mode>mpi</mode>
  <NumMPI>2</NumMPI>
</RunInfo>
```

```
<Files>
  <Input name='lorenzAttractor.py'
    type=''>lorenzAttractor.py</Input>
</Files>
```

This examples specifies the working directory (`WorkingDir`) where the necessary file (`Files`) is located and to run a series of 100 (`batchSize`) Monte-Carlo calculations (`Sequence`). MPI mode (`mode`) is used along with 4 threads (`NumThreads`) and 2 MPI processes per run (`NumMPI`).

## 9 Files

The `<Files>` block defines any files that might be needed within the RAVEN run. This could include inputs to the Model, pickled ROM files, or CSV files for postprocessors, to name a few. Each entry in the `<Files>` block is a tag with the file type. Files given through the input XML at this point are all `<Input>` type. Each `<Input>` node has the following attributes:

- **name**, *required string attribute*, user-defined name of the file. This does not need to be the actual filename; this is the name by which RAVEN will identify the file. **Note:** As with other objects, this name can be used to refer to this specific entity from other input blocks in the XML.
- **type**, *optional string attribute*, a type label for this file. While RAVEN does not directly make use of file types, they are available in the CodeInterface as identifiers. If not provided, the type will be stored as python `None` type.
- **perturbable**, *optional boolean attribute*, flag to indicate whether a file can be perturbed or not. RAVEN does not directly use this attribute, but it is available in the CodeInterface. If not provided, defaults to `True`.

For example, if the files `templateInput.i`, `materials.i`, `history.i`, `mesh.e` are required to run a Model, the `<Files>` block might appear as:

```
...
<Files>
  <Input name='main' type='maininput'>templateInput.i</Input>
  <Input name='mat' type='mtlinput' >materials.i</Input>
  <Input name='hist' type='histinput'>history.i</Input>
  <Input name='mesh' type='mesh'
    perturbable='false'>mesh.e</Input>
</Files>
...
</Simulation>
```

## 10 Distributions

RAVEN provides support for several probability distributions. Currently, the user can choose among several 1-dimensional distributions and  $N$ -dimensional ones, either custom or multidimensional normal.

The user will specify the probability distributions, that need to be used during the simulation, within the `<Distributions>` XML block:

```
<Simulation>
  ...
  <Distributions>
    <!-- All the necessary distributions will be listed here -->
  </Distributions>
  ...
</Simulation>
```

In the next two sub-sections, the input requirements for all of the distributions are reported.

### 10.1 1-Dimensional Probability Distributions

This sub-section is organized in two different parts: 1) continuous 1-D distributions and 2) discrete 1-D distributions. These two paragraphs cover all the requirements for using the different distribution entities.

#### 10.1.1 1-Dimensional Continuous Distributions

In this paragraph all the 1-D distributions currently available in RAVEN are reported.

Firstly, all the probability distributions functions in the code can be truncated by using the following keywords:

```
<Distributions>
  ...
  <aDistributionType>
    ...
    <lowerBound>aFloatValue</lowerBound>
    <upperBound>aFloatValue</upperBound>
    ...
  </aDistributionType>
```

Each distribution has a pre-defined, default support (domain) based on its definition, however these domains can be shifted/stretched using the appropriate **<low>** and **<high>** parameters where applicable, and/or truncated using the nodes in the example above, namely **<lowerBound>** and **<upperBound>**. For example, the Normal distribution domain is  $[-\infty, +\infty]$ , and thus cannot be shifted or stretched, as it is already unbounded, but can be truncated. RAVEN currently provides support for 13 1-Dimensional distributions. In the following paragraphs, all the input requirements are reported and commented.

### 10.1.1.1 Beta Distribution

The **Beta** distribution is parameterized by two positive shape parameters, denoted by  $\alpha$  and  $\beta$ , that appear as exponents of the random variable. Its default support (domain) is  $x \in (0, 1)$ . The distribution domain can be changed, specifying new boundaries, to fit the user's needs. The user can specify a **Beta** distribution in two ways. The standard is to provide the parameters **<low>**, **<high>**, **<alpha>**, and **<beta>**. Alternatively, to approximate a normal distribution that falls to 0 at the endpoints, the user may provide the parameters **<low>**, **<high>**, and **<peakFactor>**. The peak factor is a value between 0 and 1 that determines the peakedness of the distribution. At 0 it is dome-like ( $\alpha = \beta = 4$ ) and at 1 it is very strongly peaked around the mean ( $\alpha = \beta = 100$ ). A reasonable approximation to a Gaussian normal is a peak factor of 0.5.

The specifications of this distribution must be defined within a **<Beta>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following children nodes:

- Standard initialization:
  - **<alpha>**, *float, conditional required parameter*, first shape parameter. If specified, **<beta>** must also be inputted and **<peakFactor>** can not be specified.
  - **<beta>**, *float, conditional required parameter*, second shape parameter. If specified, **<alpha>** must also be inputted and **<peakFactor>** can not be specified.
  - **<low>**, *float, optional parameter*, lower domain boundary.  
*Default: 0.0*
  - **<high>**, *float, optional parameter*, upper domain, boundary.  
*Default: 1.0*

- Alternative initialization:
  - **<peakFactor>**, *float, optional parameter*, alternative to specifying **<alpha>** and **<beta>**. Acceptable values range from 0 to 1.
  - **<low>**, *float, optional parameter*, lower domain boundary.  
*Default: 0.0*
  - **<high>**, *float, optional parameter*, upper domain, boundary.  
*Default: 1.0*

Example:

```

<Distributions>
...
<Beta name='aUserDefinedName'>
  <low>aFloatValue</low>
  <high>aFloatValue</high>
  <alpha>aFloatValue</alpha>
  <beta>aFloatValue</beta>
</Beta>
<Beta name='aUserDefinedName2'>
  <low>aFloatValue</low>
  <high>aFloatValue</high>
  <peakFactor>aFloatValue</peakFactor>
</Beta>
...
</Distributions>

```

### 10.1.1.2 Exponential Distribution

The **Exponential** distribution has a default support of  $x \in [0, +\infty)$ .

The specifications of this distribution must be defined within an **<Exponential>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following child node:

- **<lambda>**, *float, required parameter*, rate parameter.



- **<low>**, *float, optional parameter*, lower domain boundary.  
Default: 0.0

Example:

```
<Distributions>
...
<Exponential name='aUserDefinedName'>
  <lambda>aFloatValue</lambda>
  <low>aFloatValue</low>
</Exponential>
...
</Distributions>
```

### 10.1.1.3 Gamma Distribution

The **Gamma** distribution is a two-parameter family of continuous probability distributions. The common exponential distribution and  $\chi$ -squared distribution are special cases of the gamma distribution. Its default support is  $x \in (0, +\infty)$ .

The specifications of this distribution must be defined within a **<Gamma>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following children nodes:

- **<alpha>**, *float, required parameter*, shape parameter.
- **<beta>**, *float, optional parameter*, 1/scale or the inverse scale parameter.  
Default: 1.0
- **<low>**, *float, optional parameter*, lower domain boundary.  
Default: 0.0

Example:

```
<Distributions>
...
<Gamma name='aUserDefinedName'>
  <alpha>aFloatValue</alpha>
```

```

    <beta>aFloatValue</beta>
    <low>aFloatValue</low>
  </Gamma>
  ...
</Distributions>

```

#### 10.1.1.4 Logistic Distribution

The **Logistic** distribution is similar to the normal distribution with a CDF that is an instance of a logistic function ( $Cdf(x) = \frac{1}{1+e^{-\frac{(x-location)}{scale}}}$ ). It resembles the normal distribution in shape but has heavier tails (higher kurtosis). Its default support is  $x \in (-\infty, +\infty)$ .

The specifications of this distribution must be defined within a **<Logistic>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following children nodes:

- **<location>**, *float, required parameter*, the distribution mean.
- **<scale>**, *float, required parameter*, scale parameter that is proportional to the standard deviation ( $\sigma^2 = \frac{1}{3}\pi^2 scale^2$ ).

Example:

```

<Distributions>
  ...
  <Logistic name='aUserDefinedName'>
    <location>aFloatValue</location>
    <scale>aFloatValue</scale>
  </Logistic>
  ...
</Distributions>

```

#### 10.1.1.5 LogNormal Distribution

The **LogNormal** distribution is a distribution with the logarithm of the random variable being normally distributed. Its default support is  $x \in (0, +\infty)$ .

The specifications of this distribution must be defined within a `<LogNormal>` XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following children nodes:

- `<mean>`, *float, required parameter*, the distribution mean or expected value (in log-scale).
- `<sigma>`, *float, required parameter*, standard deviation.
- `<low>`, *float, optional parameter*, lower domain boundary.  
*Default: 0.0*

Example:

```
<Distributions>
...
<LogNormal name='aUserDefinedName'>
  <mean>aFloatValue</mean>
  <sigma>aFloatValue</sigma>
  <low>aFloatValue</low>
</LogNormal>
...
</Distributions>
```

### 10.1.1.6 Normal Distribution

The **Normal** distribution is an extremely useful continuous distribution. Its utility is due to the central limit theorem, which states that, under mild conditions, the mean of many random variables independently drawn from the same distribution is distributed approximately normally, irrespective of the form of the original distribution. Its default support is  $x \in (-\infty, +\infty)$ .

The specifications of this distribution must be defined within a `<Normal>` XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following children nodes:

- **<mean>**, *float, required parameter*, the distribution mean or expected value.
- **<sigma>**, *float, required parameter*, the standard deviation.

Example:

```
<Distributions>
...
<Normal name='aUserDefinedName' >
  <mean>aFloatValue</mean>
  <sigma>aFloatValue</sigma>
</Normal>
...
</Distributions>
```

### 10.1.1.7 Triangular Distribution

The **Triangular** distribution is a continuous distribution that has a triangular shape for its PDF. Like the uniform distribution, upper and lower limits are “known,” but a “best guess,” of the mode or center point is also added. It has been recommended as a “proxy” for the beta distribution. Its default support is  $min \leq x \leq max$ .

The specifications of this distribution must be defined within a **<Triangular>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following children nodes:

- **<apex>**, *float, required parameter*, peak location
- **<min>**, *float, required parameter*, domain lower boundary.
- **<max>**, *float, required parameter*, domain upper boundary.

Example:

```
<Distributions>
...
<Triangular name='aUserDefinedName' >
  <apex>aFloatValue</apex>
```

```
<min>aFloatValue</min>
<max>aFloatValue</max>
</Triangular>
...
</Distributions>
```

### 10.1.1.8 Uniform Distribution

The **Uniform** distribution is a continuous distribution with a rectangular-shaped PDF. It is often used where the distribution is only vaguely known, but upper and lower limits are known. Its default support is  $lower \leq x \leq upper$ .

The specifications of this distribution must be defined within a **<Uniform>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following children nodes:

- **<lowerBound>**, *float, required parameter*, domain lower boundary.
- **<upperBound>**, *float, required parameter*, domain upper boundary.

**Note:** Since the Uniform distribution is a rectangular-shaped PDF, the truncation does not have any effect; this is the reason why the children nodes are the ones generally used for truncated distributions. Example:

```
<Distributions>
...
<Uniform name='aUserDefinedName'>
  <lowerBound>aFloatValue</lowerBound>
  <upperBound>aFloatValue</upperBound>
</Uniform>
...
</Distributions>
```

### 10.1.1.9 Weibull Distribution

The **Weibull** distribution is a continuous distribution that is often used in the field of failure analysis; in particular, it can mimic distributions where the failure rate varies over time. If the failure

rate is:

- constant over time, then  $k = 1$ , suggests that items are failing from random events;
- decreases over time, then  $k < 1$ , suggesting “infant mortality”;
- increases over time, then  $k > 1$ , suggesting “wear out” - more likely to fail as time goes by.

Its default support is  $x \in [0, +\infty)$ .

The specifications of this distribution must be defined within a `<Weibull>` XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following children nodes:

- `<k>`, *float, required parameter*, shape parameter.
- `<lambda>`, *float, required parameter*, scale parameter.
- `<low>`, *float, optional parameter*, lower domain boundary.  
*Default: 0.0*

Example:

```
<Distributions>
...
<Weibull name='aUserDefinedName' >
  <lambda>aFloatValue</lambda>
  <k>aFloatValue</k>
  <low>aFloatValue</low>
</Weibull>
...
</Distributions>
```

### 10.1.2 1-Dimensional Discrete Distributions.

RAVEN currently supports 3 discrete distributions. In the following paragraphs, the input requirements are reported.

### 10.1.2.1 Bernoulli Distribution

The **Bernoulli** distribution is a discrete distribution of the outcome of a single trial with only two results, 0 (failure) or 1 (success), with a probability of success  $p$ . It is the simplest building block on which other discrete distributions of sequences of independent Bernoulli trials can be based. Basically, it is the binomial distribution ( $k = 1, p$ ) with only one trial. Its default support is  $k \in 0, 1$ .

The specifications of this distribution must be defined within a `<Bernoulli>` XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following child node:

- `<p>`, *float, required parameter*, probability of success.

Example:

```
<Distributions>
...
<Bernoulli name='aUserDefinedName'>
  <p>aFloatValue</p>
</Bernoulli>
...
</Distributions>
```

### 10.1.2.2 Binomial Distribution

The **Binomial** distribution is the discrete probability distribution of the number of successes in a sequence of  $n$  independent yes/no experiments, each of which yields success with probability  $p$ . Its default support is  $k \in 0, 1, 2, \dots, n$ .

The specifications of this distribution must be defined within a `<Binomial>` XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following children nodes:

- `<p>`, *float, required parameter*, probability of success.
- `<n>`, *integer, required parameter*, number of experiments.

Example:

```
<Distributions>
...
<Binomial name='aUserDefinedName'>
  <n>aIntegerValue</n>
  <p>aFloatValue</p>
</Binomial>
...
</Distributions>
```

### 10.1.2.3 Poisson Distribution

The **Poisson** distribution is a discrete probability distribution that expresses the probability of a given number of events occurring in a fixed interval of time and/or space if these events occur with a known average rate and independently of the time since the last event. Its default support is  $k \in 1, 2, 3, 4, \dots$

The specifications of this distribution must be defined within a `<Poisson>` XML block. This XML node accepts one attribute:

- `name`, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following child node:

- `<mu>`, *float, required parameter*, mean rate of events/time.

Example:

```
<Distributions>
...
<Poisson name='aUserDefinedName'>
  <mu>aFloatValue</mu>
</Poisson>
...
</Distributions>
```



#### 10.1.2.4 Categorical Distribution

The **Categorical** distribution is a discrete distribution that describes the result of a random variable that can have K possible outcomes. The probability of each outcome is separately specified. The possible outcomes must be numerical values (float) There is not necessarily an underlying ordering of these outcomes, but labels are assigned in describing the distribution (in the range 1 to K). The specifications of this distribution must be defined within a **<Categorical>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

This distribution can be initialized with the following child node:

- **<state>**, *float, required parameter*, probability for outcome 1
  - **outcome**, *float, required parameter*, outcome value.
- **<state>**, *float, required parameter*, probability for outcome 2
  - **outcome**, *float, required parameter*, outcome value.
- ...
- **<state>**, *float, required parameter*, probability for outcome K
  - **outcome**, *float, required parameter*, outcome value.

Example:

```
<Distributions>
...
  <Categorical name='testCategorical'>
    <state outcome="10">0.1</state>
    <state outcome="20">0.2</state>
    <state outcome="50">0.15</state>
    <state outcome="60">0.4</state>
    <state outcome="90">0.15</state>
  </Categorical>
...
</Distributions>
```

## 10.2 N-Dimensional Probability Distributions

The group of  $N$ -Dimensional distributions allow the user to model stochastic dependences between parameters. Thus instead of using  $N$  distributions for  $N$  parameters, the user can define a single distribution lying in a  $N$ -Dimensional space. The following  $N$ -Dimensional Probability Distributions are available within RAVEN:

- `MultivariateNormal`: Multivariate normal distribution (see Section 10.2.1)
- `NDInverseWeight`: ND Inverse Weight interpolation distribution (see Section 10.2.2)
- `NDCartesianSpline`: ND spline interpolation distribution (see Section 10.2.3)

For `NDInverseWeight` and `NDCartesianSpline` distributions, the user provides the sampled values of either CDF or PDF of the distribution. The sampled values can be scattered distributed (for `NDInverseWeight`) or over a cartesian grid (for `NDCartesianSpline`).

The user could specify, for each  $N$ -Dimensional distribution, the parameters of the random number generator function:

- `<initialGridDisc>`, *positive integer, optional field*, user-defined initial grid discretization. This parameter specifies the number of discretizations that need to be performed, initially, for each Dimension in order to find  $N$ -Dimensional coordinate that corresponds to the CDF represented by a random number (0-1);
- `<tolerance>`, *float, optional field*, user-defined tolerance in order to find the  $N$ -D coordinates corresponding to a random number. This tolerance is expressed in terms of CDF.

in the `<samplerInit>` block defined in sampler block `<samplerInit>` (see Section 11).

### 10.2.1 MultivariateNormal Distribution

the multivariate normal distribution or multivariate Gaussian distribution, is a generalization of the one-dimensional (univariate) normal distribution to higher dimensions. The multivariate normal distribution is often used to describe, at least approximately, any set of (possibly) correlated real-valued random variables each of which clusters around a mean value. The multivariate normal distribution of a  $k$ -dimensional random vector  $\mathbf{x} = [x_1, x_2, \dots, x_k]$  can be written in the following notation:  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with with  $k$ -dimensional mean vector

$$\boldsymbol{\mu} = [E[x_1], E[x_2], \dots, E[x_k]]$$

and  $k \times k$  covariance matrix

$$\Sigma = [Cov[x_i, x_j]], i = 1, 2, \dots, k; j = 1, 2, \dots, k$$

The probability distribution function for this distribution is the following:

$$f_{\mathbf{x}}(x_1, \dots, x_k) = \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right),$$

The specifications of this distribution must be defined within the xml block **<MultivariateNormal>**. This XML node needs to contain the attributes:

- **name**, *required string attribute*, user-defined identifier of this multivariate normal distribution. **Note:** As with other objects, this is the name that can be used to refer to this specific entity from other input XML blocks.
- **method**, *required string attribute*, defines which method is used to generate the multivariate normal distribution. The only allowable methods are 'spline' and 'pca'.

In RAVEN the MultivariateNormal distribution can be initialized through the following keywords:

- **<mu>**, list of mean values of each dimension
- **<covariance>**, list of element values in the covariance matrix. There are two types of **<covariance>**, based on the **type**:
  - **type**, *string, optional field*, specifies the type of covariance, the default **type** is 'abs'. Possible values for **type** are 'abs' and 'rel'. **Note:** 'abs' indicates the covariance is a normal covariance matrix, while 'rel' indicates the covariance is a relative covariance matrix. In addition, method 'pca' can be combined with both types, and method 'spline' only accept the type 'abs'
- **<transformation>**, *XML node, optional field*, option to enable input parameter transformation using principal component analysis (PCA) approach. If this node is provided, PCA will be used to compute the principal components of input covariance matrix. The subnode **<rank>** is used to indicate the number of principal components that will be used for the input transformation. The content will specify one attribute:
  - **<rank>**, *positive integer, required field*, user-defined dimensionality reduction.

Example:

```
<Distributions>
...
  <MultivariateNormal name='MultivariateNormal_test'
    method='spline'>
    <mu>0.0 60.0</mu>
    <covariance>
```

```

    1.0 0.7
    0.7 1.0
  </covariance>
</MultivariateNormal>
<MultivariateNormal name='MultivariateNormal_abs'
  method='pca'>
  <mu>0.0 60.0</mu>
  <covariance type='abs'>
    1.0 0.7
    0.7 1.0
  </covariance>
</MultivariateNormal>
<MultivariateNormal name='MultivariateNormal_rel'
  method='pca'>
  <mu>0.0 60.0</mu>
  <covariance type='rel'>
    1.0 0.7
    0.7 1.0
  </covariance>
</MultivariateNormal>
...
</Distributions>

```

In the following, we defined a distribution with a transformation node using PCA method. The number of principal components is defined in `<rank>`. In this distribution, PCA is employed to restruct the multivariate normal distribution. In addition, the size of uncorrelated variables is also determined by `<rank>`.

```

<Distributions>
...
  <MultivariateNormal name='MultivariateNormal_test'
    method='pca'>
      <mu>0.0 10.0 20.0</mu>
      <covariance type="abs">
        1.0    0.7    -0.2
        0.7    1.0    0.4
        -0.2   0.4    1.0
      </covariance>
      <transformation>
        <rank>2</rank>
      </transformation>
    </MultivariateNormal>
...

```

```
</Distributions>
```

## 10.2.2 NDInverseWeight Distribution

The NDInverseWeight distribution creates a  $N$ -Dimensional distribution given a set of points scattered distributed. These points sample the PDF of the original distribution. Distribution values (PDF or CDF) are calculated using the inverse weight interpolation scheme.

The specifications of this distribution must be defined within a `<NDInverseWeight>` XML block. This XML node accepts the following attributes:

- `name`, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

In RAVEN the NDInverseWeight distribution can be initialized through the following nodes:

- `<p>`, *float, required parameter*, power parameter. Greater values of  $p$  assign greater influence to values closest to the interpolated point.
- `<data_filename>`, *string, required parameter*, name of the data file containing scattered values (file type '.txt').
  - `type`, *required string attribute*, indicates if the data in indicated file is PDF or CDF.
- `<working_dir>`, *string, required parameter*, folder location of the data file

Example:

```
<Distributions>
...
<NDInverseWeight name='...'>
  <p>...</p>
  <dataFilename type='...'>...</dataFilename>
  <workingDir>...</workingDir>
</NDInverseWeight>
...
</Distributions>
```

Each data entry contained in `data_filename` is listed row by row and must be listed as follows:

- number of dimensions

- number of sampled points
- ND coordinate of each sampled point
- value of each sampled point

As an example, the following shows the data entries contained in `data_filename` for a 3-dimensional data set that contained two sampled CDF values:  $([0.0,0.0,0.0], 0.1)$  and  $([1.0, 1.0,0.0], 0.8)$

Example scattered data file:

```
3
2
0.0
0.0
0.0
1.0
1.0
0.0
0.1
0.8
```

### 10.2.3 NDCartesianSpline Distribution

The NDCartesianSpline distribution creates a  $N$ -Dimensional distribution given a set of points regularly distributed on a cartesian grid. These points sample the PDF of the original distribution. Distribution values (PDF or CDF) are calculated using the ND spline interpolation scheme.

The specifications of this distribution must be defined within a `<NDCartesianSpline>` XML block. This XML node accepts the following attributes:

- `name`, *required string attribute*, user-defined name of this distribution. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

In RAVEN the NDCartesianSpline distribution can be initialized through the following nodes:

- `<data_filename>`, *string, required parameter*, name of the data file containing scattered values (file type '.txt').
  - `type`, *required string attribute*, indicates if the data in indicated file is PDF or CDF.
- `<working_dir>`, *string, required parameter*, folder location of the data file

Example:

```
<Distributions>
...
<NDCartesianSpline name='...'>
  <dataFilename type='...'>...</dataFilename>
  <workingDir></workingDir>
</NDCartesianSpline>
...
</Distributions>
```

Each data entry contained in data \_filename is listed row by row and must be listed as follows:

- number of dimensions
- number of discretization for each dimension
- discretization values for each dimension
- value of each sampled point

As an example, the following shows the data entries contained in data \_filename for a 2-dimensional CDF data set on the following grid  $(x, y)$ :

- first dimension (x): -0.5, 0.5
- first dimension (y): 1.0 2.0 3.0

Example scattered data file:

```
2
2
3
-0.5
0.5
1.0
2.0
3.0
CDF value of (-0.5, 1.0)
CDF value of (+0.5, 1.0)
CDF value of (-0.5, 2.0)
CDF value of (+0.5, 2.0)
CDF value of (-0.5, 3.0)
CDF value of (+0.5, 3.0)
```

## 11 Samplers

The sampler is probably the most important entity in the RAVEN framework. It performs the driving of the specific sampling strategy and, hence, determines the effectiveness of the analysis, from both an accuracy and computational point of view. The samplers, that are available in RAVEN, can be categorized into three main classes:

- **Forward** (see Section 11.1)
- **Dynamic Event Tree (DET)** (see Section 11.2)
- **Adaptive** (see Section 11.3)

Before analyzing each sampler in detail, it is important to mention that each type has a similar syntax to input the variables to be “sampled”. In the example below, the variable ‘**variableName**’ is going to be sampled by the Sampler ‘**whatever**’ using the distribution named ‘**aDistribution**’.

```
<Simulation>
...
<Samplers>
...
<WhateverSampler name='whatever'>
...
  <variable name='variableName'>
    ...
    <distribution>aDistribution</distribution>
    ...
  </variable>
  ...
</WhateverSampler>
...
</Samplers>
...
</Simulation>
```

As reported in section 18, the variable naming syntax, for external driven codes, depends on the way the “code interface” has been implemented. For example, if the code has an input structure like the one reported below (YAML), the variable name might be ‘**I-Level | II-Level | variable**’ . In this way, the relative code interface (and input parser) will know which variable needs to be perturbed and the “recipe” to access it. As reported in 18, its syntax is chosen by the developer of



the “code interface” and is implemented in the interface only (no modifications are needed in the RAVEN code).

Example YAML based Input:

```
[I-Level]
  [./II-Level]
    variable = xxx
  [../]
[]
```

Example XML block to define the variables and associated distributions:

```
<variable name='I-Level|II-Level|variable'>
  <distribution>exampleDistribution</distribution>
</variable>
```

If the variable is associated to a multi-dimensional ND distribution, it is needed to specify which dimension of the ND distribution is associated to such variable. An example is shown below: the variable “variableX” is associated to the third dimension of the ND distribution “ND-distribution”.

```
<variable name='variableX'>
  <distribution dim='3'>NDdistribution</distribution>
</variable>
```

For most codes, it is prudent that there are no redundant inputs; however there are cases where this is not reality. For example, if there is a variable ‘**inner\_radius**’ and a variable ‘**outer\_radius**’, there may be a third variable ‘**thickness**’ that is actually derived from the previous two, as ‘**thickness**’ = ‘**outer\_radius**’ - ‘**inner\_radius**’. RAVEN supports this type of redundant input through a Function entity. In this case, instead of a **<distribution>** node in the **<variable>** block, there is a **<Function>** node, specifying the name of the function (defined in the **<Functions>** block). In order to work properly, this function must have a method that is identical to the variable name that returns a single python float object. For example,

```
...
<Functions>
  <External name='torus_calcs' file='torus_calcs.py'>
    <variable>outer_radius</variable>
    <variable>inner_radius</variable>
  </External>
</Functions>
...
<Samplers>
```

```

<WhateverSampler name='myExampleSampler'>
  <variable name='inner_radius'>
    <distribution>inner_dist</distribution>
  </variable>
  <variable name='outer_radius'>
    <distribution>outer_dist</distribution>
  </variable>
  <variable name='thickness'>
    <Function>torus_calcs</Function>
  </variable>
</WhateverSampler>
</Samplers>

```

The corresponding function file '`torus_calcs.py`' needs the following method:

```

def thickness(self):
    return self.outer_radius - self.inner_radius

```

The '`thickness`' parameter will still be treated as an input for the sake of csv printing and DataObjects storage.

In the sampler class a special node exists: the `<sampler_init>` node. This node contains specific parameters that characterize each particular sampler. In addition, `<sampler_init>` might contain the information regarding the random generator function for each  $N$ -Dimensional distribution (specified in the `<dist_init>` node):

- `initial_grid_disc`
- `tolerance`

An example of `<dist_init>` node is provided below:

```

<distInit>
  <distribution name= 'ND_dist_name'>
    <initialGridDisc>5</initialGridDisc>
    <tolerance>0.2</tolerance>
  </distribution>
</distInit>

```

In the `<sampler_init>` node it is possible to add also the subnode `<globalGrid>`. The `<globalGrid>` can be used in two cases:

- 1D distributions: an identical grid that is associated to several distributions

- ND distribution: a grid associated to a single ND distribution. This is the case when a stratified sampling is performed on the CDF of an ND distribution: the `<globalGrid>` is shared among the variables associated to the Nd distribution

## 11.1 Forward Samplers

The Forward sampler category collects all the strategies that perform the sampling of the input space without exploiting, through dynamic learning approaches, the information made available from the outcomes of calculations previously performed (adaptive sampling) and the common system evolution (patterns) that different sampled calculations can generate in the phase space (dynamic event tree). In the RAVEN framework, five different “Forward” samplers are available:

- **Monte Carlo (MC)**
- **Stratified**
- **Grid Based**
- **Sparse Grid Collocation**
- **Sobol Decomposition**
- **Response Surface Design of Experiment**
- **Factorial Design of Experiment**

From a practical point of view, these sampling strategies represent different ways to explore the input space. In the following paragraphs, the input requirements and a small explanation of the different sampling methodologies are reported.

### 11.1.1 Monte Carlo

The **Monte-Carlo** sampling approach is one of the most well-known and widely used approaches to perform exploration of the input space. The main idea behind MonteCarlo sampling is to randomly perturbate the input space according to uniform or parameter-based probability density functions.

The specifications of this sampler must be defined within a `<MonteCarlo>` XML block. This XML node accepts the following attributes:

- **name**, *required string attribute*, user-defined name of this Sampler. N.B. As for the other objects, this is the name that can be used to refer to this specific entity from other input blocks (xml);

In the **MonteCarlo** input block, the user needs to specify the variables need to be sampled. As already mentioned, these variables are inputted within consecutive xml blocks called **<variable>**. In addition, the settings for this sampler need to be specified in the **<samplerInit>** XML block:

- **<samplerInit>**, *XML node, required parameter*. In this xml-node, the following xml sub-nodes need to be specified:
  - **<limit>**, *integer, required field*, number of MonteCarlo samples needs to be generated;
  - **<initialSeed>**, *integer, optional field*, initial seeding of random number generator
  - **<reseedEachIteration>**, *boolean/string(case insensitive), optional field*, perform a re-seeding for each sample generated (True values = True, yes, y, t).  
*Default: False;*
  - **<distInit>**, *integer, optional field*, in this node the user specifies the initialization of the random number generator function for each N-Dimensional Probability Distributions (see Section 10.2).
- **<variable>**, *XML node, required parameter* will specify one attribute:
  - **name**, *required string attribute*, user-defined name of this variable.

This **<variable>** recognizes the following child nodes:

- **<distribution>**, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the **<Distributions>** block explained in Section 10. **Note:** Alternatively, this node must be omitted if the **<function>** node is supplied.
- **<function>**, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the **<Functions>** block explained in Section 16. **Note:** Alternatively, this node must be omitted if the **<distribution>** node is supplied.

If the input parameters are correlated, the **MonteCarlo** sampling approach can be also used if the user specified a multivariate distributions inside the **<Distributions>** (see Section 10.2). Furthermore, if the covariance matrix is provided and the input parameters is assumed to have the multivariate normal distribution, one can also use **MonteCarlo** approach to sample the input parameters in the transformed space (aka subspace, reduced space). If this is the case, the user needs to provide additional information, i.e. the **<transformation>** under **<MultivariateNormal>** of **<Distributions>** (more information can be found in Section 10.2). In addition, the node **<variablesTransformation>** is also required for **MonteCarlo** sampling. This node is used to transform the variables specified by **<latentVariables>** in

the transformed space of input into variables specified by `<manifestVariables>` in the input space. The variables listed in `<latentVariables>` should be predefined in `<variable>`, and the variables listed in `<manifestVariables>` are used by the `<Models>`.

- `<variablesTransformation>`, *optional field*. this XML node accepts one attribute:
  - `model`, *required string attribute*, the name for the model defined in the XML node `<Models>`. This attribute indicates the `<manifestVariables>` defined in the subnode will be used by the `model`.

In addition, this XML node also accepts three children nodes:

- `<latentVariables>`, *comma separated string, required field*, user-defined latent variables that are used for the variables transformation. All the variables listed under this node should be also mentioned in `<variable>`.
- `<manifestVariables>`, *comma separated string, required field*, user-defined manifest variables that are used by the `model`.
- `<method>`, *string, required field*, the method that is used for the variables transformation. The currently available method is 'pca'.

**Assembler Objects** These objects are either required or optional depending on the functionality of the MonteCarlo Sampler. The objects must be listed with a rigorous syntax that, except for the XML node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to identify them within the simulation framework:

- `class`, *required string attribute*, the main “class” of the listed object. For example, it can be 'Models', 'Functions', etc.
- `type`, *required string attribute*, the object identifier or sub-type. For example, it can be 'ROM', 'External', etc.

The **MonteCarlo** approach requires or optionally accepts the following object types:

- `<Restart>`, *string, optional field*, the body of this XML node must contain the name of an appropriate **PointSet** defined in the `<DataObjects>` block (see Section 12). It is used as a “restart” tool, where it accepts pre-existing solutions in the PointSet instead of recalculating solutions.

Example:

```

<Samplers>
...
<MonteCarlo name='MCname'>
  <samplerInit>
    <limit>10</limit>
    <initialSeed>200286</initialSeed>
    <reseedEachIteration>>false</reseedEachIteration>
    <distInit>
      <distribution name= 'ND_InverseWeight_P'>
        <initialGridDisc>10</initialGridDisc>
        <tolerance>0.2</tolerance>
      </distribution>
    </distInit>
  </samplerInit>
  <variable name='var1'>
    <distribution>aDistributionNameDefinedInDistributionBlock
  </distribution>
  </variable>
</MonteCarlo>
...
</Samplers>

```

### 11.1.2 Grid

The **Grid** sampling approach is probably the simplest exploration approach that can be employed to explore an uncertain domain. The idea is to construct an  $N$ -dimensional grid where each dimension is represented by one uncertain variable. This approach performs the sampling at each node of the grid. The sampling of the grid consists in evaluating the answer of the system under all possible combinations among the different variables' values with respect to a predefined discretization metric. In RAVEN two discretization metrics are available: 1) cumulative distribution function, and 2) value. Thus, the grid meshing can be input via probability or variable values. Regarding the  $N$ -dimensional distributions, the user can specify for each dimension the type of grid to be used (i.e., value or CDF). Note the discretization of the CDF, only for the grid sampler, is performed on the marginal distribution for the specific variable considered.

The specifications of this sampler must be defined within a **<Grid>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in

the XML.

In the `<Grid>` input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive `<variable>` XML blocks:

- `<variable>`, *XML node, required parameter* will specify one attribute:
  - `name`, *required string attribute*, user-defined name of this variable.

This `<variable>` recognizes the following child nodes:

- `<distribution>`, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the `<Distributions>` block explained in Section 10. **Note:** Alternatively, this node must be omitted if the `<function>` node is supplied.
- `<function>`, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the `<Functions>` block explained in Section 16. **Note:** Alternatively, this node must be omitted if the `<distribution>` node is supplied.
- `<grid>`, *space separated floats, required field*, the content of this XML node depends on the definition of the associated attributes:
  - \* `type`, *required string attribute*, user-defined discretization metric type: 1) `'CDF'`, the grid will be specified based on cumulative distribution function probability thresholds, and 2) `'value'`, the grid will be provided using variable values.
  - \* `construction`, *required string attribute*, how the grid needs to be constructed, independent of its type (i.e. `'CDF'` or `'value'`).

Based on the `construction` type, the content of the `<grid>` XML node and the requirements for other attributes change:

- \* `construction='equal'`. The grid is going to be constructed equally-spaced (`type='value'`) or equally probable (`type='CDF'`). This construction type requires the definition of additional attributes:
  - `steps`, *required integer attribute*, number of equally spaced/probable discretization steps.

This construction type requires that the content of the `<grid>` node represents the lower and upper bounds (either in probability or value). Two values need to be specified; the lowest one will be considered as the *lowerBound*, the largest, the *upperBound*. The lower and upper bounds are checked against the associated `<distribution>` bounds. If one or both of them falls outside the distribution's bounds, the code will raise an error. The *stepSize* is determined as follows:

$$stepSize = (upperBound - lowerBound) / steps$$



- \* **construction='custom'**. The grid will be directly specified by the user. No additional attributes are needed. This construction type requires that the `<grid>` node contains the actual mesh bins. For example, if the grid **type** is 'CDF', in the body of `<grid>`, the user will specify the CDF probability thresholds (nodalization in probability). All the bins are checked against the associated `<distribution>` bounds. If one or more of them falls outside the distribution's bounds, the code will raise an error.

**Note:** The `<grid>` node is only required if a `<distribution>` node is supplied. In the case of a `<function>` node, no grid information is requested.

If the input parameters are correlated, the **Grid** sampling approach can be also used if the user specified a multivariate distributions inside the `<Distributions>` (see Section 10.2). Furthermore, if the covariance matrix is provided and the input parameters is assumed to have the multivariate normal distribution, one can also use **Grid** approach to sample the input parameters in the transformed space (aka subspace, reduced space). This means one creates the grids of variables listed by `<latentVariables>` in the transformed space. If this is the case, the user needs to provide additional information, i.e. the `<transformation>` under `<MultivariateNormal>` of `<Distributions>` (more information can be found in Section 10.2). In addition, the node `<variablesTransformation>` is also required for **Grid** sampling. This node is used to transform the variables specified by `<latentVariables>` in the transformed space of input into variables specified by `<manifestVariables>` in the input space. The variables listed in `<latentVariables>` should be predefined in `<variable>`, and the variables listed in `<manifestVariables>` are used by the `<Models>`.

- `<variablesTransformation>`, *optional field*. this XML node accepts one attribute:
  - **model**, *required string attribute*, the name for the model defined in the XML node `<Models>`. This attribute indicates the `<manifestVariables>` defined in the subnode will be used by the **model**.

In addition, this XML node also accepts three children nodes:

- `<latentVariables>`, *comma separated string, required field*, user-defined latent variables that are used for the variables transformation. All the variables listed under this node should be also mentioned in `<variable>`.
- `<manifestVariables>`, *comma separated string, required field*, user-defined manifest variables that are used by the **model**.
- `<method>`, *string, required field*, the method that is used for the variables transformation. The currently available method is 'pca'.

**Assembler Objects** These objects are either required or optional depending on the functionality of the Grid Sampler. The objects must be listed with a rigorous syntax that, except for the XML



node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to identify them within the simulation framework:

- **class**, *required string attribute*, the main “class” of the listed object. For example, it can be 'Models', 'Functions', etc.
- **type**, *required string attribute*, the object identifier or sub-type. For example, it can be 'ROM', 'External', etc.

The **Grid** approach requires or optionally accepts the following object types:

- **<Restart>**, *string, optional field*, the body of this XML node must contain the name of an appropriate **PointSet** defined in the **<DataObjects>** block (see Section 12). It is used as a “restart” tool, where it accepts pre-existing solutions in the PointSet instead of recalculating solutions.

Example:

```
<Samplers>
...
<Grid name='Gridname'>
  <variable name='var1'>
    <distribution>aDistributionNameDefinedInDistributionBlock1
    </distribution>
    <grid type='value' construction='equal' steps='100' >0.2
      10</grid>
  </variable>
  <variable name='var2'>
    <distribution>aDistributionNameDefinedInDistributionBlock2
    </distribution>
    <grid type='CDF' construction='equal' steps='5' >0.2
      0.8</grid>
  </variable>
  <variable name='var3'>
    <distribution>aDistributionNameDefinedInDistributionBlock3
    </distribution>
    <grid type='value' construction='equal' steps='100' >0.2
      21.0</grid>
  </variable>
  <variable name='var4'>
    <distribution>aDistributionNameDefinedInDistributionBlock4
    </distribution>
```

```

    <grid type='CDF' construction='equal' steps='5' >0.2
      1.0</grid>
  </variable>
  <variable name='var5'>
    <distribution>aDistributionNameDefinedInDistributionBlock5
    </distribution>
    <grid type='value' construction='custom'>0.2 0.5
      10.0</grid>
  </variable>
  <variable name='var6'>
    <distribution>aDistributionNameDefinedInDistributionBlock6
    </distribution>
    <grid type='CDF' construction='custom'>0.2 0.5 1.0</grid>
  </variable>
</Grid>
...
</Samplers>

```

### 11.1.3 Sparse Grid Collocation

**Sparse Grid Collocation** builds on generic **Grid** sampling by selecting evaluation points based on characteristic quadratures as part of stochastic collocation for generalized polynomial chaos uncertainty quantification. In collocation you construct an N-dimensional grid, with each uncertain variable providing an axis. Along each axis, the points of evaluation correspond to quadrature points necessary to integrate polynomials (see 15.3.2). In the simplest (and most naive) case, a N-Dimensional tensor product of all possible combinations of points from each dimension's quadrature is constructed as sampling points. The number of necessary samples can be reduced by employing Smolyak-like sparse grid algorithms, which use reduced combinations of polynomial orders to reduce the necessary sampling space. The specifications of this sampler must be defined within a **<SparseGridCollocation>** XML block. .

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.
- **parallel**, *optional string attribute*, option to disable parallel construction of the sparse grid. Because of increasing computational expense with increasing input space dimension, RAVEN will default to parallel construction of the sparse grid.
- **outfile**, *optional string attribute*, option to allow the generated sparse grid points and weights to be printed to a file with the given name.  
*Default: True*

In the `<SparseGridCollocation>` input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive `<variable>` XML blocks:

- `<variable>`, *XML node, required parameter* will specify one attribute:
  - `name`, *required string attribute*, user-defined name of this variable.

In the variable node, the following xml-node needs to be specified:

- `<distribution>`, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the `<Distributions>` block explained in Section 10. **Note:** Alternatively, this node must be omitted if the `<function>` node is supplied.
- `<function>`, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the `<Functions>` block explained in Section 16. **Note:** Alternatively, this node must be omitted if the `<distribution>` node is supplied.

Because of the tight coupling between the Sampler and the ROM in stochastic collocation for generalized polynomial chaos, the Sampler needs access to the ROM via the assembler to determine the polynomials, quadratures, and importance weights to use in each dimension (see 15.3.2).

**Assembler Objects** These objects are either required or optional depending on the functionality of the SparseGridCollocation Sampler. The objects must be listed with a rigorous syntax that, except for the XML node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to identify them within the simulation framework:

- `class`, *required string attribute*, the main “class” of the listed object. For example, it can be `'Models'`, `'Functions'`, etc.
- `type`, *required string attribute*, the object identifier or sub-type. For example, it can be `'ROM'`, `'External'`, etc.

The `SparseGridCollocation` approach requires or optionally accepts the following object types:

- `<ROM>`, *string, required field*, the body of this XML node must contain the name of an appropriate ROM defined in the `<Models>` block (see Section 15.3).
- `<Restart>`, *string, optional field*, the body of this XML node must contain the name of an appropriate `PointSet` defined in the `<DataObjects>` block (see Section 12). It is used as a “restart” tool, where it accepts pre-existing solutions in the PointSet instead of recalculating solutions.

Example:

```
<Samplers>
...
<SparseGridCollocation name="mySG" parallel="0">
  <variable name="x1">
    <distribution>myDist1</distribution>
  </variable>
  <variable name="x2">
    <distribution>myDist2</distribution>
  </variable>
  <ROM class = 'Models' type = 'ROM' >SCROM</ROM>
  <Restart class = 'DataObjects' type = 'PointSet' >solns</Restart>
</SparseGridCollocation>
...
</Samplers>
```

In general, **SparseGridCollocation** requires uncorrelated input parameters. If the input parameters are correlated, one can transform the correlated parameters into uncorrelated parameters; the **SparseGridCollocation** can also be used with the uncorrelated parameters in the transformed space. Like in the **Grid** sampler, if the covariance matrix is provided and the input parameters are assumed to have the multivariate normal distribution, the **SparseGridCollocation** can be used. This means one creates the sparse grids of variables listed by **<latentVariables>** in the transformed space. If this is the case, the user needs to provide additional information, i.e. the **<transformation>** under **<MultivariateNormal>** of **<Distributions>** (more information can be found in Section 10.2). In addition, the node **<variablesTransformation>** is also required for **SparseGridCollocation** sampler. This node is used to transform the variables specified by **<latentVariables>** in the transformed space of input into variables specified by **<manifestVariables>** in the input space. The variables listed in **<latentVariables>** should be predefined in **<variable>**, and the variables listed in **<manifestVariables>** are used by the **<Models>**.

- **<variablesTransformation>**, *optional field*. this XML node accepts one attribute:
  - **model**, *required string attribute*, the name for the model defined in the XML node **<Models>**. This attribute indicates the **<manifestVariables>** defined in the subnode will be used by the **model**.

In addition, this XML node also accepts three children nodes:

- **<latentVariables>**, *comma separated string, required field*, user-defined latent variables that are used for the variables transformation. All the variables listed under this node should be also mentioned in **<variable>**.
- **<manifestVariables>**, *comma separated string, required field*, user-defined manifest variables that are used by the **model**.

- `<method>`, *string, required field*, the method that is used for the variables transformation. The currently available method is 'pca'.

```

...
<Models>
  ...
  <ExternalModel ModuleToLoad="lorenzAttractor_noK"
    name="PythonModule" subType="">
    <variable>sigma</variable>
    <variable>rho</variable>
    <variable>beta</variable>
    <variable>x</variable>
    <variable>y</variable>
    <variable>z</variable>
    <variable>time</variable>
    <variable>x0</variable>
    <variable>y0</variable>
    <variable>z0</variable>
  </ExternalModel>
  <ROM name="SCROM" subType="GaussPolynomialRom">
    <Target>and</Target>
    <Features>x1,y1,z1</Features>
    <IndexSet>TensorProduct</IndexSet>
    <PolynomialOrder>1</PolynomialOrder>
  </ROM>
  ...
</Models>

<Distributions>
  ...
  <MultivariateNormal name='MVNDist' method='pca'>
    <transformation>
      <rank>3</rank>
    </transformation>
    <mu>0.0 1.0 2.0</mu>
    <covariance type="abs">
      1.0      0.6      -0.4
      0.6      1.0      0.2
      -0.4     0.2      0.8
    </covariance>
  </MultivariateNormal>
  ...
</Distributions>

```

```

<Samplers>
...
<SparseGridCollocation name='SC'>
  <variable name='x0'>
    <distribution dim='1'>MVNDist</distribution>
  </variable>
  <variable name='y0'>
    <distribution dim='2'>MVNDist</distribution>
  </variable>
  <variable name='z0'>
    <distribution dim='3'>MVNDist</distribution>
  </variable>
  <variablesTransformation model="PythonModule">
    <latentVariables>x1, y1, z1</latentVariables>
    <manifestVariables>x0, y0, z0</manifestVariables>
    <method>pca</method>
  </variablesTransformation>
  <ROM class = 'Models' type = 'ROM' >SCROM</ROM>
</SparseGridCollocation>
...
</Samplers>
...

```

#### 11.1.4 Sobol

The **Sobol** sampler uses high-density model reduction (HDMR) a.k.a. Sobol decomposition to approximate a function as the sum of increasing-complexity interactions. At its lowest level (order 1), it treats the function as a sum of the reference case plus a functional of each input dimension separately. At order 2, it adds functionals to consider the pairing of each dimension with each other dimension. The benefit to this approach is considering several functions of small input cardinality instead of a single function with large input cardinality. This allows reduced order models like generalized polynomial chaos (see 15.3.2) to approximate the functionals accurately with few computations runs. This Sobol sampler uses the associated HDMRRom (see 15.3.3) to determine at what points the input space need be evaluated.

**Note:** Sobol sampler relies on SparseGridCollocation, which requires uncorrelated input dimensions; thus, for the present, it is not compatible with NDDistribution objects.

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in

the XML.

- **parallel**, *optional string attribute*, option to disable parallel construction of the sparse grid. Because of increasing computational expense with increasing input space dimension, RAVEN will default to parallel construction of the sparse grid.

*Default: True*

In the `<Sobol>` input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive `<variable>` XML blocks:

- `<variable>`, *XML node, required parameter* will specify one attribute:
  - **name**, *required string attribute*, user-defined name of this variable.

In the variable node, the following xml-node needs to be specified:

- `<distribution>`, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the `<Distributions>` block explained in Section 10. **Note:** Alternatively, this node must be omitted if the `<function>` node is supplied.
- `<function>`, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the `<Functions>` block explained in Section 16. **Note:** Alternatively, this node must be omitted if the `<distribution>` node is supplied.

Because of the tight coupling between the Sobol sampler and the HDMRROM, the Sampler needs access to the ROM via the assembler to determine the polynomials, quadratures, Sobol order, and importance weights to use in each dimension (see 15.3.3).

**Assembler Objects** These objects are either required or optional depending on the functionality of the Sobol Sampler. The objects must be listed with a rigorous syntax that, except for the XML node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to identify them within the simulation framework:

- **class**, *required string attribute*, the main “class” of the listed object. For example, it can be `'Models'`, `'Functions'`, etc.
- **type**, *required string attribute*, the object identifier or sub-type. For example, it can be `'ROM'`, `'External'`, etc.

The **Sobol** approach requires or optionally accepts the following object types:

- **<ROM>**, *string, required field*, the body of this XML node must contain the name of an appropriate ROM defined in the **<Models>** block (see Section 15.3).
- **<Restart>**, *string, optional field*, the body of this XML node must contain the name of an appropriate **PointSet** defined in the **<DataObjects>** block (see Section 12). It is used as a “restart” tool, where it accepts pre-existing solutions in the PointSet instead of recalculating solutions.

Example:

```

<Samplers>
...
<Sobol name="mySobol" parallel="0">
  <variable name="x1">
    <distribution>myDist1</distribution>
  </variable>
  <variable name="x2">
    <distribution>myDist2</distribution>
  </variable>
  <ROM class = 'Models' type = 'ROM' >myHDMR</ROM>
</Sobol>
...
</Samplers>

```

### 11.1.5 Stratified

The **Stratified** sampling approach is a method for the exploration of the input space that consists of dividing the uncertain domain into subgroups before sampling. In the “stratified” sampling, these subgroups must be:

- mutually exclusive: every element in the population must be assigned to only one stratum (subgroup);
- collectively exhaustive: no population element can be excluded.

Then simple random sampling or systematic sampling is applied within each stratum. It is worthwhile to note that the well-known Latin hypercube sampling represents a specialized version of the stratified approach, when the domain strata are constructed in equally-probable CDF bins.

The specifications of this sampler must be defined within a **<Stratified>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.



In the `<Stratified>` input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive `<variable>` XML blocks:

- `<variable>`, *XML node, required parameter* will specify one attribute:
  - `name`, *required string attribute*, user-defined name of this variable.

This `<variable>` recognizes the following child nodes:

- `<distribution>`, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the `<Distributions>` block explained in Section 10. **Note:** Alternatively, this node must be omitted if the `<function>` node is supplied.
- `<function>`, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the `<Functions>` block explained in Section 16. **Note:** Alternatively, this node must be omitted if the `<distribution>` node is supplied.
- `<grid>`, *space separated floats, required field*, the content of this XML node depends on the definition of the associated attributes:
  - \* `type`, *required string attribute*, user-defined discretization metric type: 1) `'CDF'`, the grid will be specified based on cumulative distribution function probability thresholds, and 2) `'value'`, the grid will be provided using variable values.
  - \* `construction`, *required string attribute*, how the grid needs to be constructed, independent of its type (i.e. `'CDF'` or `'value'`).

Based on the `construction` type, the content of the `<grid>` XML node and the requirements for other attributes change:

- \* `construction='equal'`. The grid is going to be constructed equally-spaced (`type='value'`) or equally probable (`type='CDF'`). This construction type requires the definition of additional attributes:
  - `steps`, *required integer attribute*, number of equally spaced/probable discretization steps.

This construction type requires that the content of the `<grid>` node represents the lower and upper bounds (either in probability or value). Two values need to be specified; the lowest one will be considered as the *lowerBound*, the largest, the *upperBound*. The lower and upper bounds are checked against the associated `<distribution>` bounds. If one or both of them falls outside the distribution's bounds, the code will raise an error. The *stepSize* is determined as follows:

$$stepSize = (upperBound - lowerBound) / steps$$

- \* `construction='custom'`. The grid will be directly specified by the user. No additional attributes are needed. This construction type requires that the `<grid>` node contains the actual mesh bins. For example, if the grid `type` is `'CDF'`, in the

body of `<grid>`, the user will specify the CDF probability thresholds (nodalization in probability). All the bins are checked against the associated `<distribution>` bounds. If one or more of them falls outside the distribution's bounds, the code will raise an error.

**Note:** The `<grid>` node is only required if a `<distribution>` node is supplied. In the case of a `<function>` node, no grid information is requested.

As one can see, the input specifications for the **Stratified** sampler are similar to that of the **Grid** sampler. It is important to mention again that for each zone (grid mesh) only a point, randomly selected, is picked and not all the nodal combinations (like in the **Grid** sampling).

**Assembler Objects** These objects are either required or optional depending on the functionality of the Stratified Sampler. The objects must be listed with a rigorous syntax that, except for the XML node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to identify them within the simulation framework:

- **class**, *required string attribute*, the main “class” of the listed object. For example, it can be 'Models', 'Functions', etc.
- **type**, *required string attribute*, the object identifier or sub-type. For example, it can be 'ROM', 'External', etc.

The **Stratified** approach requires or optionally accepts the following object types:

- `<Restart>`, *string, optional field*, the body of this XML node must contain the name of an appropriate **PointSet** defined in the `<DataObjects>` block (see Section 12). It is used as a “restart” tool, where it accepts pre-existing solutions in the PointSet instead of recalculating solutions.

Example:

```
<Samplers>
...
<Stratified name='StratifiedName'>
  <variable name='var1'>
    <distribution>aDistributionNameDefinedInDistributionBlock1
  </distribution>
  <grid type='CDF' construction='equal' steps='5' >0.2
    0.8</grid>
  </variable>
  <variable name='var2'>
```

```

    <distribution>aDistributionNameDefinedInDistributionBlock2
  </distribution>
  <grid type='value' construction='equal' steps='100' >0.2
    21.0</grid>
</variable>
<variable name='var3'>
  <distribution>aDistributionNameDefinedInDistributionBlock3
  </distribution>
  <grid type='CDF' construction='custom'>0.2 0.5 1.0</grid>
</variable>
</Stratified>
...
</Samplers>

```

For N-dimensional (ND) distributions, there are two different approaches to perform the stratified sampling. In the first approach, the subgroups is determined by the joint CDF of given multivariate distributions. If this approach is used, the sampling is performed on a grid on a CDF, while the user is required to specify the same CDF grid for all the dimensions of the ND distribution. This is possible by defining a `<globalGrid>` node and associate such `<globalGrid>` to each variable belonging to the ND distribution as follows.

```

<Samplers>
...
<Stratified name='StratifiedName'>
  <variable name='x0'>
    <distribution
      dim='1'>ND_InverseWeight_P</distribution>
    <grid type='globalGrid'>name_grid1</grid>
  </variable>
  <variable name='y0, z0'>
    <distribution
      dim='2'>ND_InverseWeight_P</distribution>
    <grid type='globalGrid'>name_grid1</grid>
  </variable>
  <globalGrid>
    <grid name='name_grid1' type='CDF'
      construction='custom'>0.1 1.0 0.2</grid>
  </globalGrid>
</Stratified>
...
</Samplers>
...

```

The second approach is different than the first approach. Like in the **Grid** sampling, if the covariance matrix is provided and the input parameters is assumed to have the multivariate normal distribution, one can also use **Stratified** approach to sample the input parameters in the transformed space (aka subspace, reduced space). This means one creates the grids of variables listed by `<latentVariables>` in the transformed space. If this is the case, the user needs to provide additional information, i.e. the `<transformation>` under `<MultivariateNormal>` of `<Distributions>` (more information can be found in Section 10.2). In addition, the node `<variablesTransformation>` is also required for **Stratified** sampler. This node is used to transform the variables specified by `<latentVariables>` in the transformed space of input into variables specified by `<manifestVariables>` in the input space. The variables listed in `<latentVariables>` should be predefined in `<variable>`, and the variables listed in `<manifestVariables>` are used by the `<Models>`. In addition, `<globalGrid>` will be not used for approach.

- `<variablesTransformation>`, *optional field*. this XML node accepts one attribute:
  - `model`, *required string attribute*, the name for the model defined in the XML node `<Models>`. This attribute indicates the `<manifestVariables>` defined in the subnode will be used by the `model`.

In addition, this XML node also accepts three children nodes:

- `<latentVariables>`, *comma separated string, required field*, user-defined latent variables that are used for the variables transformation. All the variables listed under this node should be also mentioned in `<variable>`.
- `<manifestVariables>`, *comma separated string, required field*, user-defined manifest variables that are used by the `model`.
- `<method>`, *string, required field*, the method that is used for the variables transformation. The currently available method is 'pca'.

```

...
<Models>
  ...
  <ExternalModel ModuleToLoad="lorenzAttractor_noK"
    name="PythonModule" subType="">
    <variable>sigma</variable>
    <variable>rho</variable>
    <variable>beta</variable>
    <variable>x</variable>
    <variable>y</variable>
    <variable>z</variable>
    <variable>time</variable>
  </ExternalModel>

```

```

    <variable>x0</variable>
    <variable>y0</variable>
    <variable>z0</variable>
  </ExternalModel>
  ...
</Models>

<Distributions>
  ...
  <MultivariateNormal name='MVNDist' method='pca'>
    <transformation>
      <rank>3</rank>
    </transformation>
    <mu>0.0 1.0 2.0</mu>
    <covariance type="abs">
      1.0      0.6      -0.4
      0.6      1.0      0.2
      -0.4     0.2      0.8
    </covariance>
  </MultivariateNormal>
  ...
</Distributions>

<Samplers>
  ...
  <Stratified name='StratifiedName'>
    <variable name='x0'>
      <distribution dim='1'>MVNDist</distribution>
      <grid type='CDF' construction='equal' steps='3'>0.1
        0.9</grid>
    </variable>
    <variable name='y0'>
      <distribution dim='2'>MVNDist</distribution>
      <grid type='value' construction='equal'
        steps='3'>0.1 0.9</grid>
    </variable>
    <variable name='z0'>
      <distribution dim='3'>MVNDist</distribution>
      <grid type='CDF' construction='equal' steps='3'>0.2
        0.8</grid>
    </variable>
    <variablesTransformation model="PythonModule">

```

```

        <latentVariables>x1, y1, z1</latentVariables>
        <manifestVariables>x0, y0, z0</manifestVariables>
        <method>pca</method>
    </variablesTransformation>
</Stratified>
...
</Samplers>
...

```

### 11.1.6 Response Surface Design

The **Response Surface Design**, or Response Surface Modeling (RSM), approach is one of the most common Design of Experiment (DOE) methodologies currently in use. It explores the relationships between several explanatory variables and one or more response variables. The main idea of RSM is to use a sequence of designed experiments to obtain an optimal response. RAVEN currently employs two different algorithms that can be classified within this family of methods:

- **Box-Behnken:** This methodology aims to achieve the following goals:
  - Each factor, or independent variable, is placed at one of three equally spaced values, usually coded as -1, 0, +1. (At least three levels are needed for the following goal);
  - The design should be sufficient to fit a quadratic model, that is, one squared term per factor and the products of any two factors;
  - The ratio of the number of experimental points to the number of coefficients in the quadratic model should be reasonable (in fact, their designs keep it in the range of 1.5 to 2.6);
  - The estimation variance should more or less depend only on the distance from the center (this is achieved exactly for the designs with 4 and 7 factors), and should not vary too much inside the smallest (hyper)cube containing the experimental points.

Each design can be thought of as a combination of a two-level (full or fractional) factorial design with an incomplete block design. In each block, a certain number of factors are put through all combinations for the factorial design, while the other factors are kept at the central values.

- **Central Composite:** This design consists of three distinct sets of experimental runs:
  - A factorial (perhaps fractional) design in the factors are studied, each having two levels;
  - A set of center points, experimental runs whose values of each factor are the medians of the values used in the factorial portion. This point is often replicated in order to improve the precision of the experiment;

- A set of axial points, experimental runs identical to the centre points except for one factor, which will take on values both below and above the median of the two factorial levels, and typically both outside their range. All factors are varied in this way.

This methodology is useful for building a second order (quadratic) model for the response variable without needing to use a complete three-level factorial experiment.

All the parameters, needed for setting up the algorithms reported above, must be defined within a **<ResponseSurfaceDesign>** block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

In the **<ResponseSurfaceDesign>** input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive **<variable>** XML blocks:

- **<variable>**, *XML node, required parameter* will specify one attribute:
  - **name**, *required string attribute*, user-defined name of this variable.

This **<variable>** recognizes the following child nodes:

- **<distribution>**, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the **<Distributions>** block explained in Section 10. **Note:** Alternatively, this node must be omitted if the **<function>** node is supplied.
- **<function>**, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the **<Functions>** block explained in Section 16. **Note:** Alternatively, this node must be omitted if the **<distribution>** node is supplied.
- **<grid>**, *space separated floats, required field*, the content of this XML node depends on the definition of the associated attributes:
  - \* **type**, *required string attribute*, user-defined discretization metric type: 1) 'CDF', the grid will be specified based on cumulative distribution function probability thresholds, and 2) 'value', the grid will be provided using variable values.
  - \* **construction**, *required string attribute*, how the grid needs to be constructed, independent of its type (i.e. 'CDF' or 'value').

Based on the **construction** type, the content of the **<grid>** XML node and the requirements for other attributes change. In this case, only the following is available:



- \* **construction='custom'**. The grid will be directly specified by the user. This construction type requires that the **<grid>** node contains the actual mesh bins. For example, if the grid **type** is 'CDF', in the body of **<grid>**, the user will specify the CDF probability thresholds (nodalization in probability). All the bins are checked against the associated **<distribution>** bounds. If one or more of them falls outside the distribution's bounds, the code will raise an error. No additional attributes are needed.

**Note:** The **<grid>** node is only required if a **<distribution>** node is supplied. In the case of a **<function>** node, no grid information is requested. **Note:** Only the construction "custom" is available. In the **<grid>** body only the lower and upper bounds can be inputted (2 numbers only).

- **<ResponseSurfaceDesignSettings>**, *required*, In this sub-node, the user needs to specify different settings depending on the algorithm being used:

- **<algorithmType>**, *string, required field*, this XML node will contain the name of the algorithm to be used. Based on the chosen algorithm, other nodes need to be defined:

- \* **<algorithmType>**BoxBehnken**<algorithmType/>**. If Box-Behnken is specified, the following additional node is recognized:

- **<ncenters>**, *integer, optional field*, the number of center points to include in the box. If this parameter is not specified, then a pre-determined number of points are automatically included.

*Default: Automatic Generation.*

**Note:** In order to employ the "Box-Behnken" design, at least 3 variables must be used.

- \* **<algorithmType>**CentralComposite**<algorithmType/>**. If Central Composite is specified, the following additional nodes will be recognized:

- **<centers>**, *comma separated integers, optional field*, the number of center points to be included. This block needs to contain 2 integers values separated by a comma. The first entry represents the number of centers to be added for the factorial block; the second one is the one for the star block.

*Default: 4,4.*

- **<alpha>**, *string, optional field*, in this node, the user decides how an  $\alpha$  factor needs to be determined. Two options are available:

**orthogonal** for orthogonal design.

**rotatable** for rotatable design.

*Default: orthogonal.*

- **<face>**, *string, optional field*, in this node, the user defines how faces should be constructed. Three options are available:



**circumscribed** for circumscribed facing  
**inscribed** for inscribed facing  
**faced** for faced facing.

*Default: circumscribed.*

**Note:** In order to employ the “Central Composite” design, at least 2 variables must be used.

Furthermore, if the covariance matrix is provided and the input parameters are assumed to have a multivariate normal distribution, one can use **ResponseSurfaceDesign** approach to sample the input parameters in the transformed space (aka subspace, reduced space). In this case, the user needs to provide additional information, i.e. the **<transformation>** under **<MultivariateNormal>** of **<Distributions>** (more information can be found in Section 10.2). In addition, the node **<variablesTransformation>** is also required for **ResponseSurfaceDesign** sampling. This node is used to transform the variables specified by **<latentVariables>** in the transformed space of input into variables specified by **<manifestVariables>** in the input space. The variables listed in **<latentVariables>** should be predefined in **<variable>**, and the variables listed in **<manifestVariables>** are used by the **<Models>**.

- **<variablesTransformation>**, *optional field*. this XML node accepts one attribute:
  - **model**, *required string attribute*, the name for the model defined in the XML node **<Models>**. This attribute indicates the **<manifestVariables>** defined in the subnode will be used by the **model**.

In addition, this XML node also accepts three children nodes:

- **<latentVariables>**, *comma separated string, required field*, user-defined latent variables that are used for the variables transformation. All the variables listed under this node should be also mentioned in **<variable>**.
- **<manifestVariables>**, *comma separated string, required field*, user-defined manifest variables that are used by the **model**.
- **<method>**, *string, required field*, the method that is used for the variables transformation. The currently available method is 'pca'.

Example:

```
<Samplers>
...
  <ResponseSurfaceDesign name='BoxBehnkenRespDesign' >
    <ResponseSurfaceDesignSettings>
```

```

        <algorithmType>BoxBehnken</algorithmType>
        <ncenters>1</ncenters>
    </ResponseSurfaceDesignSettings>
    <variable name='var1' >
        <distribution >Gauss1</distribution>
        <grid type='CDF' construction='custom' >0.2
            0.8</grid>
    </variable>
    <!-- N.B. at least 3 variables need to inputted
        in order to employ this algorithm
    -->
</ResponseSurfaceDesign>
<ResponseSurfaceDesign name='CentralCompositeRespDesign'>
    <ResponseSurfaceDesignSettings>
        <algorithmType>CentralComposite</algorithmType>
        <centers>1, 2</centers>
        <alpha>orthogonal</alpha>
        <face>circumscribed</face>
    </ResponseSurfaceDesignSettings>
    <variable name='var4' >
        <distribution >Gauss1</distribution>
        <grid type='CDF' construction='custom' >0.2
            0.8</grid>
    </variable>
    <!-- N.B. at least 2 variables need to inputted
        in order to employ this algorithm
    -->
</ResponseSurfaceDesign>
<ResponseSurfaceDesign name='transformedSpaceSampling'>
    <ResponseSurfaceDesignSettings>
        <algorithmType>BoxBehnken</algorithmType>
        <ncenters>1</ncenters>
    </ResponseSurfaceDesignSettings>
    <variable name='var1' >
        <distribution >Gauss1</distribution>
        <grid type='CDF' construction='custom' >0.2
            0.8</grid>
    </variable>
    ...
    <variablesTransformation model="givenModel">
        <latentVariables>var1, ...</latentVariables>
        <manifestVariables>...</manifestVariables>

```

```
        <method>pca</method>
      </variablesTransformation>
</ResponseSurfaceDesign>
...
</Samplers>
```

### 11.1.7 Factorial Design

The **Factorial Design** method is an important method to determine the effects of multiple variables on a response. A factorial design can reduce the number of samples one has to perform by studying multiple factors simultaneously. Additionally, it can be used to find both main effects (from each independent factor) and interaction effects (when both factors must be used to explain the outcome). A factorial design tests all possible conditions. Because factorial designs can lead to a large number of trials, which can become expensive and time-consuming, they are best used for small numbers of variables with only a few domain discretizations (1 to 3). Factorial designs work well when interactions between variables are strong and important and where every variable contributes significantly. RAVEN currently employs three different algorithms that can be classified within this family of techniques:

- **General Full Factorial** explores the input space by investigating all possible combinations of a set of factors (variables).
- **2-Level Fractional-Factorial** consists of a carefully chosen subset (fraction) of the experimental runs of a full factorial design. The subset is chosen so as to exploit the sparsity-of-effects principle exposing information about the most important features of the problem studied, while using a fraction of the effort of a full factorial design in terms of experimental runs and resources.
- **Plackett-Burman** identifies the most important factors early in the experimentation phase when complete knowledge about the system is usually unavailable. It is an efficient screening method for identifying the active factors (variables) using as few samples as possible. In Plackett-Burman designs, main effects have a complicated confounding relationship with two-factor interactions. Therefore, these designs should be used to study main effects when it can be assumed that two-way interactions are negligible.

All the parameters needed for setting up the algorithms reported above must be defined within a **<FactorialDesign>** block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

In the `<FactorialDesign>` input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive `<variable>` XML blocks:

- `<variable>`, *XML node, required parameter* will specify one attribute:
  - `name`, *required string attribute*, user-defined name of this variable.

This `<variable>` recognizes the following child nodes:

- `<distribution>`, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the `<Distributions>` block explained in Section 10. **Note:** Alternatively, this node must be omitted if the `<function>` node is supplied.
- `<function>`, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the `<Functions>` block explained in Section 16. **Note:** Alternatively, this node must be omitted if the `<distribution>` node is supplied.
- `<grid>`, *space separated floats, required field*, the content of this XML node depends on the definition of the associated attributes:
  - \* `type`, *required string attribute*, user-defined discretization metric type: 1) `'CDF'`, the grid will be specified based on cumulative distribution function probability thresholds, and 2) `'value'`, the grid will be provided using variable values.
  - \* `construction`, *required string attribute*, how the grid needs to be constructed, independent of its type (i.e. `'CDF'` or `'value'`).

Based on the `construction` type, the content of the `<grid>` XML node and the requirements for other attributes change:

- \* `construction='equal'`. The grid is going to be constructed equally-spaced (`type='value'`) or equally probable (`type='CDF'`). This construction type requires the definition of additional attributes:
  - `steps`, *required integer attribute*, number of equally spaced/probable discretization steps.

This construction type requires that the content of the `<grid>` node represents the lower and upper bounds (either in probability or value). Two values need to be specified; the lowest one will be considered as the *lowerBound*, the largest, the *upperBound*. The lower and upper bounds are checked against the associated `<distribution>` bounds. If one or both of them falls outside the distribution's bounds, the code will raise an error. The *stepSize* is determined as follows:

$$stepSize = (upperBound - lowerBound) / steps$$

- \* `construction='custom'`. The grid will be directly specified by the user. No additional attributes are needed. This construction type requires that the `<grid>` node contains the actual mesh bins. For example, if the grid `type` is `'CDF'`, in the

body of `<grid>`, the user will specify the CDF probability thresholds (nodalization in probability). All the bins are checked against the associated `<distribution>` bounds. If one or more of them falls outside the distribution's bounds, the code will raise an error.

**Note:** The `<grid>` node is only required if a `<distribution>` node is supplied. In the case of a `<function>` node, no grid information is requested.

The main `<FactorialDesign>` block needs to contain an additional sub-node called `<FactorialSettings>`. In this sub-node, the user needs to specify different settings depending on the algorithm being used:

- `<algorithmType>`, *string, required field*, specifies the algorithm to be used. Based on the chosen algorithm, other nodes may be defined:
  - `<algorithmType>full<algorithmType/>`. Full factorial design. If `full` is specified, no additional nodes are necessary.
    - Note:** The full factorial design does not have any limitations on the number of discretization bins that can be used in the `<grid>` XML node for each `<variable>` specified.
  - `<algorithmType>2levelFract<algorithmType/>`. Two-level Fractional-Factorial design. If `2levelFract` is specified, the following additional nodes must be specified:
    - \* `<gen>`, *space separated strings, required field*, specifies the confounding mapping. For instance, in this block the user defines the decisions on a fraction of the full-factorial by allowing some of the factor main effects to be compounded with other factor interaction effects. This is done by defining an alias structure that defines, symbolically, these interactions. These alias structures are written like “C = AB” or “I = ABC”, or “AB = CD”, etc. These define how a column is related to the others.
    - \* `<genMap>`, *space separated strings, required field*, defines the mapping between the `<gen>` symbolic aliases and the variables that have been inputted in the `<FactorialDesign>` main block.
      - Note:** The Two-levels Fractional-Factorial design is limited to 2 discretization bins in the `<grid>` node for each `<variable>`.
  - `<algorithmType>pb<algorithmType/>`. Plackett-Burman design. If `pb` is specified, no additional nodes are necessary.
    - Note:** The Plackett-Burman design does not have any limitations on the number of discretization bins allowed in the `<grid>` node for each `<variable>`.

Example:

```

<Samplers>
...
<FactorialDesign name='fullFactorial'>
  <FactorialSettings>
    <algorithmType>full</algorithmType>
  </FactorialSettings>
  <variable name='var1' >
    <distribution>aDistributionNameDefinedInDistributionBlock1
    </distribution>
    <grid type='value' construction='custom' >0.02 0.03
      0.5</grid>
  </variable>
  <variable name='var2' >
    <distribution>aDistributionNameDefinedInDistributionBlock2
    </distribution>
    <grid type='CDF' construction='custom'>0.5 0.7 1.0</grid>
  </variable>
</FactorialDesign>
<FactorialDesign name='2levelFractFactorial'>
  <FactorialSettings>
    <algorithmType>2levelFract</algorithmType>
    <gen>a,b,ab</gen>
    <genMap>var1,var2,var3</genMap>
  </FactorialSettings>
  <variable name='var1' >
    <distribution>aDistributionNameDefinedInDistributionBlock3
    </distribution>
    <grid type='value' construction='custom' >0.02 0.5</grid>
  </variable>
  <variable name='var2' >
    <distribution>aDistributionNameDefinedInDistributionBlock
    </distribution>
    <grid type='CDF' construction='custom'>0.5 1.0</grid>
  </variable>
  <variable name='var3'>
    <distribution>aDistributionNameDefinedInDistributionBlock5
    </distribution>
    <grid type='value' upperBound='4' construction='equal'
      steps='1'>0.5</grid>
  </variable>
</FactorialDesign>
<FactorialDesign name='pbFactorial'>

```

```

<FactorialSettings>
  <algorithmType>pb</algorithmType>
</FactorialSettings>
<variable name='var1' >
  <distribution>aDistributionNameDefinedInDistributionBlock6
  </distribution>
  <grid type='value' construction='custom' >0.02 0.5</grid>
</variable>
<variable name='VarGauss2' >
  <distribution>aDistributionNameDefinedInDistributionBlock7
  </distribution>
  <grid type='CDF' construction='custom'>0.5 1.0</grid>
</variable>
</FactorialDesign>
...
</Samplers>

```

## 11.2 Dynamic Event Tree (DET) Samplers

The **Dynamic Event Tree** methodologies are designed to take the timing of events explicitly into account, which can become very important especially when uncertainties in complex phenomena are considered. Hence, the main idea of this methodology is to let a system code determine the pathway of an accident scenario within a probabilistic environment. In this family of methods, a continuous monitoring of the system evolution in the phase space is needed. In order to use the DET-based methods, the generic driven code needs to have, at least, an internal trigger system and, consequently, a “restart” capability. In the RAVEN framework, 4 different DET samplers are available:

- **Dynamic Event Tree (DET)**
- **Hybrid Dynamic Event Tree (HDET)**
- **Adaptive Dynamic Event Tree (ADET)**
- **Adaptive Hybrid Dynamic Event Tree (AHDET)**

The ADET and the AHDET methodologies represent a hybrid between the DET/HDET and adaptive sampling approaches. For this reason, its input requirements are reported in the Adaptive Samplers’ section (11.3).



### 11.2.1 Dynamic Event Tree

The **Dynamic Event Tree** sampling approach is a sampling strategy that is designed to take the timing of events, in transient/accident scenarios, explicitly into account. From an application point of view, an  $N$ -Dimensional grid is built on the CDF space. A single simulation is spawned and a set of triggers is added to the system code control logic. Every time a trigger is activated (one of the CDF thresholds in the grid is exceeded), a new set of simulations (branches) is spawned. Each branch carries its conditional probability. In the RAVEN code, the triggers are defined by specifying a grid using a predefined discretization metric in the mode input space. RAVEN provides two discretization metrics: 1) CDF, and 2) value. Thus, the trigger thresholds can be entered either in probability or value space.

The specifications of this sampler must be defined within a `<DynamicEventTree>` XML block. This XML node accepts the following attributes:

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.
- **printEndXML**, *optional string/boolean attribute*, controls the dumping of a “summary” of the DET performed into an external XML.  
*Default: False.*
- **maxSimulationTime**, *optional float attribute*, this attribute controls the maximum “mission” time of the simulation underneath.  
*Default: None.*

In the `<DynamicEventTree>` input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive `<variable>` XML blocks:

- `<variable>`, *XML node, required parameter* will specify one attribute:
  - **name**, *required string attribute*, user-defined name of this variable.

This `<variable>` recognizes the following child nodes:

- `<distribution>`, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the `<Distributions>` block explained in Section 10. **Note:** Alternatively, this node must be omitted if the `<function>` node is supplied.
- `<function>`, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the `<Functions>` block explained in Section 16. **Note:** Alternatively, this node must be omitted if the `<distribution>` node is supplied.



- **<grid>**, *space separated floats, required field*, the content of this XML node depends on the definition of the associated attributes:
  - \* **type**, *required string attribute*, user-defined discretization metric type: 1) 'CDF', the grid will be specified based on cumulative distribution function probability thresholds, and 2) 'value', the grid will be provided using variable values.
  - \* **construction**, *required string attribute*, how the grid needs to be constructed, independent of its type (i.e. 'CDF' or 'value').

Based on the **construction** type, the content of the **<grid>** XML node and the requirements for other attributes change:

- \* **construction='equal'**. The grid is going to be constructed equally-spaced (**type='value'**) or equally probable (**type='CDF'**). This construction type requires the definition of additional attributes:
  - **steps**, *required integer attribute*, number of equally spaced/probable discretization steps.

This construction type requires that the content of the **<grid>** node represents the lower and upper bounds (either in probability or value). Two values need to be specified; the lowest one will be considered as the *lowerBound*, the largest, the *upperBound*. The lower and upper bounds are checked against the associated **<distribution>** bounds. If one or both of them falls outside the distribution's bounds, the code will raise an error. The *stepSize* is determined as follows:

$$stepSize = (upperBound - lowerBound) / steps$$

- \* **construction='custom'**. The grid will be directly specified by the user. No additional attributes are needed. This construction type requires that the **<grid>** node contains the actual mesh bins. For example, if the grid **type** is 'CDF', in the body of **<grid>**, the user will specify the CDF probability thresholds (nodalization in probability). All the bins are checked against the associated **<distribution>** bounds. If one or more of them falls outside the distribution's bounds, the code will raise an error.

**Note:** The **<grid>** node is only required if a **<distribution>** node is supplied. In the case of a **<function>** node, no grid information is requested.

Example:

```
<Samplers>
...
<DynamicEventTree name='DETname' >
  <variable name='var1' >
    <distribution>aDistributionNameDefinedInDistributionBlock1
    </distribution>
    <grid type='value' construction='equal' steps='100'
      lowerBound='1.0' >0.2</grid>
```

```

</variable>
<variable name='var2'>
  <distribution>aDistributionNameDefinedInDistributionBlock2
  </distribution>
  <grid type='CDF' construction='equal' steps='5'
    lowerBound='0.0'>0.2</grid>
</variable>
<variable name='var3'>
  <distribution>aDistributionNameDefinedInDistributionBlock3
  </distribution>
  <grid type='value' construction='equal' steps='100'
    upperBound='21.0'>0.2</grid>
</variable>
<variable name='var4'>
  <distribution>aDistributionNameDefinedInDistributionBlock4
  </distribution>
  <grid type='CDF' construction='equal' steps='5'
    upperBound='1.0'>0.2</grid>
</variable>
<variable name='var5'>
  <distribution>aDistributionNameDefinedInDistributionBlock5
  </distribution>
  <grid type='value' construction='custom'>0.2 0.5
    10.0</grid>
</variable>
<variable name='var6'>
  <distribution>aDistributionNameDefinedInDistributionBlock6
  </distribution>
  <grid type='CDF' construction='custom'>0.2 0.5 1.0</grid>
</variable>
</DynamicEventTree>
...
</Samplers>

```

### 11.2.2 Hybrid Dynamic Event Tree

The **Hybrid Dynamic Event Tree** sampling approach is a sampling strategy that represents an evolution of the Dynamic Event Tree method for the simultaneous exploration of the epistemic and aleatory uncertain space. In similar approaches, the uncertainties are generally treated by employing a Monte-Carlo sampling approach (epistemic) and DET methodology (aleatory). The HDET methodology, developed within the RAVEN code, can reproduce the capabilities employed by this

approach, but provides additional sampling strategies to the user. The epistemic or epistemic-like uncertainties can be sampled through the following strategies:

- Monte-Carlo;
- Grid sampling;
- Stratified (e.g., Latin Hyper Cube).

From a practical point of view, the user defines the parameters that need to be sampled by one or more different approaches. The HDET module samples those parameters creating an  $N$ -dimensional grid characterized by all the possible combinations of the input space coordinates coming from the different sampling strategies. Each coordinate in the input space represents a separate and parallel standard DET exploration of the uncertain domain. The HDET methodology allows the user to explore the uncertain domain employing the best approach for each variable kind. The addition of a grid sampling strategy among the usable approaches allows the user to perform a discrete parametric study under aleatory and epistemic uncertainties.

Regarding the input requirements, the HDET sampler is a “sub-type” of the `<DynamicEventTree>` sampler. For this reason, its specifications must be defined within a `<DynamicEventTree>` block. This XML node accepts the following attributes:

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.
- **printEndXML**, *optional string/boolean attribute*, controls the dumping of a “summary” of the DET performed into an external XML.  
*Default: False.*
- **maxSimulationTime**, *optional float attribute*, this attribute controls the maximum “mission” time of the simulation underneath.  
*Default: None.*

In the `<DynamicEventTree>` input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive `<variable>` XML blocks:

- `<variable>`, *XML node, required parameter* will specify one attribute:
  - **name**, *required string attribute*, user-defined name of this variable.

This `<variable>` recognizes the following child nodes:

- **<distribution>**, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the **<Distributions>** block explained in Section 10. **Note:** Alternatively, this node must be omitted if the **<function>** node is supplied.
- **<function>**, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the **<Functions>** block explained in Section 16. **Note:** Alternatively, this node must be omitted if the **<distribution>** node is supplied.
- **<grid>**, *space separated floats, required field*, the content of this XML node depends on the definition of the associated attributes:
  - \* **type**, *required string attribute*, user-defined discretization metric type: 1) 'CDF', the grid will be specified based on cumulative distribution function probability thresholds, and 2) 'value', the grid will be provided using variable values.
  - \* **construction**, *required string attribute*, how the grid needs to be constructed, independent of its type (i.e. 'CDF' or 'value').

Based on the **construction** type, the content of the **<grid>** XML node and the requirements for other attributes change:

- \* **construction='equal'**. The grid is going to be constructed equally-spaced (**type='value'**) or equally probable (**type='CDF'**). This construction type requires the definition of additional attributes:
  - **steps**, *required integer attribute*, number of equally spaced/probable discretization steps.

This construction type requires that the content of the **<grid>** node represents the lower and upper bounds (either in probability or value). Two values need to be specified; the lowest one will be considered as the *lowerBound*, the largest, the *upperBound*. The lower and upper bounds are checked against the associated **<distribution>** bounds. If one or both of them falls outside the distribution's bounds, the code will raise an error. The *stepSize* is determined as follows:

$$stepSize = (upperBound - lowerBound) / steps$$
- \* **construction='custom'**. The grid will be directly specified by the user. No additional attributes are needed. This construction type requires that the **<grid>** node contains the actual mesh bins. For example, if the grid **type** is 'CDF', in the body of **<grid>**, the user will specify the CDF probability thresholds (nodalization in probability). All the bins are checked against the associated **<distribution>** bounds. If one or more of them falls outside the distribution's bounds, the code will raise an error.

**Note:** The **<grid>** node is only required if a **<distribution>** node is supplied. In the case of a **<function>** node, no grid information is requested.

In order to activate the **Hybrid Dynamic Event Tree** sampler, the main **<DynamicEventTree>** block needs to contain, at least, an additional sub-node called **<HybridSampler>**. As already

mentioned, the user can combine the Monte-Carlo, Stratified, and Grid approaches in order to create a “pre-sampling”  $N$ -dimensional grid, from whose nodes a standard DET method is employed. For this reason, the user can specify a maximum of three `<HybridSampler>` sub-nodes (i.e. one for each of the available Forward samplers). This sub-node needs to contain the following attribute:

- **type**, *required string attribute*, type of pre-sampling strategy to be used. Available options are 'MonteCarlo', 'Grid', and 'Stratified'.

Independent of the type of “pre-sampler” that has been specified, the `<HybridSampler>` must contain the variables that need to be sampled. As already mentioned, these variables are specified within consecutive `<variable>` XML blocks:

- `<variable>`, *XML node, required parameter* will specify one attribute:
  - **name**, *required string attribute*, user-defined name of this variable.

This `<variable>` recognizes the following child nodes:

- `<distribution>`, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the `<Distributions>` block explained in Section 10. **Note:** Alternatively, this node must be omitted if the `<function>` node is supplied.
- `<function>`, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the `<Functions>` block explained in Section 16. **Note:** Alternatively, this node must be omitted if the `<distribution>` node is supplied.

If a pre-sampling strategy **type** is either 'Grid' or 'Stratified', within the `<variable>` blocks, the user needs to specify the sub-node `<grid>`. As with the standard DET, the content of this XML node depends on the definition of the associated attributes:

- **type**, *required string attribute*, user-defined discretization metric type:
  - 'CDF', the grid is going to be specified based on the cumulative distribution function probability thresholds
  - 'value', the grid is going to be provided using variable values.
- **construction**, *required string attribute*, how the grid needs to be constructed, independent of its type (i.e. 'CDF' or 'value').

Based on the `construction` type, the content of the `<grid>` XML node and the requirements for other attributes change:

- `construction='equal'`. The grid is going to be constructed equally-spaced (`type='value'`) or equally probable (`type='CDF'`). This construction type requires the definition of additional attributes:
  - `steps`, *required integer attribute*, number of equally spaced/probable discretization steps.

This construction type requires that the content of the `<grid>` node represents the lower and upper bounds (either in probability or value). Two values need to be specified; the lowest one will be considered as the *lowerBound*, the largest, the *upperBound*. The lower and upper bounds are checked against the associated `<distribution>` bounds. If one or both of them falls outside the distribution's bounds, the code will raise an error. The *stepSize* is determined as follows:

$$stepSize = (upperBound - lowerBound) / steps$$

- `construction='custom'`. The grid will be directly specified by the user. No additional attributes are needed. This construction type requires that the `<grid>` node contains the actual mesh bins. For example, if the grid `type` is `'CDF'`, in the body of `<grid>`, the user will specify the CDF probability thresholds (nodalization in probability). All the bins are checked against the associated `<distribution>` bounds. If one or more of them falls outside the distribution's bounds, the code will raise an error.

Example:

```
<Samplers>
...
<DynamicEventTree name='HybridDETname' print_end_XML="True">
  <HybridSampler type='MonteCarlo' limit='2'>
    <variable name='var1' >
      <distribution>aDistributionNameDefinedInDistributionBlock1
    </distribution>
    </variable>
    <variable name='var2' >
      <distribution>aDistributionNameDefinedInDistributionBlock2
    </distribution>
    <grid type='CDF' construction='equal' steps='1'
      lowerBound='0.1'>0.1</grid>
    </variable>
  </HybridSampler>
  <HybridSampler type='Grid'>
    <!-- Point sampler way (directly sampling the variable) -->
```

```

<variable name='var3' >
  <distribution>aDistributionNameDefinedInDistributionBlock3
  </distribution>
  <grid type='CDF' construction='equal' steps='1'
    lowerBound='0.1'>0.1</grid>
</variable>
<variable name='var4' >
  <distribution>aDistributionNameDefinedInDistributionBlock4
  </distribution>
  <grid type='CDF' construction='equal' steps='1'
    lowerBound='0.1'>0.1</grid>
</variable>
</HybridSampler>
<HybridSampler type='Stratified'>
  <!-- Point sampler way (directly sampling the variable )
  -->
  <variable name='var5' >
    <distribution>aDistributionNameDefinedInDistributionBlock5
    </distribution>
    <grid type='CDF' construction='equal' steps='1'
      lowerBound='0.1'>0.1</grid>
  </variable>
  <variable name='var6' >
    <distribution>aDistributionNameDefinedInDistributionBlock6
    </distribution>
    <grid type='CDF' construction='equal' steps='1'
      lowerBound='0.1'>0.1</grid>
  </variable>
</HybridSampler>
<!-- DYNAMIC EVENT TREE INPUT (it goes outside an inner
  block like HybridSamplerSettings) -->
<Distribution name='dist7'>
  <distribution>aDistributionNameDefinedInDistributionBlock7
  </distribution>
  <grid type='CDF' construction='custom'>0.1 0.8</grid>
</Distribution>
</DynamicEventTree>
...
</Samplers>

```



## 11.3 Adaptive Samplers

The Adaptive Samplers' family provides the possibility to perform smart sampling (also known as adaptive sampling) as an alternative to classical "Forward" techniques. The motivation is that system simulations are often computationally expensive, time-consuming, and high dimensional with respect to the number of input parameters. Thus, exploring the space of all possible simulation outcomes is infeasible using finite computing resources. During simulation-based probabilistic risk analysis, it is important to discover the relationship between a potentially large number of input parameters and the output of a simulation using as few simulation trials as possible.

The description above characterizes a typical context for performing adaptive sampling where a few observations are obtained from the simulation, a reduced order model (ROM) is built to represent the simulation space, and new samples are selected based on the model constructed. The reduced order model (see section 15.3) is then updated based on the simulation results of the sampled points. In this way, an attempt is made to gain the most information possible with a small number of carefully selected sample points, limiting the number of expensive trials needed to understand features of the system space.

Currently, RAVEN provides support for the following adaptive algorithms:

- Limit Surface Search
- Adaptive Dynamic Event Tree
- Adaptive Hybrid Dynamic Event Tree
- Adaptive Sparse Grid
- Adaptive Sobol Decomposition

In the following paragraphs, the input requirements and a small explanation of the different sampling methods are reported.

### 11.3.1 Limit Surface Search

The **Limit Surface Search** approach is an advanced methodology that employs a smart sampling around transition zones that determine a change in the status of the system (limit surface). To perform such sampling, RAVEN uses ROMs for predicting, in the input space, the location(s) of these transitions, in order to accelerate the exploration of the input space in proximity of the limit surface.

The specifications of this sampler must be defined within an `<LimitSurfaceSearch>` XML block. This XML node accepts one attribute:



- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

In the `<LimitSurfaceSearch>` input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive `<variable>` XML blocks:

- This `<variable>` recognizes the following child nodes:
  - `<distribution>`, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the `<Distributions>` block explained in Section 10. **Note:** Alternatively, this node must be omitted if the `<function>` node is supplied.
  - `<function>`, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the `<Functions>` block explained in Section 16. **Note:** Alternatively, this node must be omitted if the `<distribution>` node is supplied.

In addition to the `<variable>` nodes, the main XML node `<Adaptive>` needs to contain two supplementary sub-nodes:

- `<Convergence>`, *float, required field*, Convergence tolerance. The meaning of this tolerance depends on the definition of other attributes that might be defined in this XML node:
  - **limit**, *optional integer attribute*, the maximum number of adaptive samples (iterations).  
*Default: infinite.*
  - **forceIteration**, *optional boolean attribute*, this attribute controls if at least a number of iterations equal to **limit** must be performed.  
*Default: False.*
  - **weight**, *optional string attribute (case insensitive)*, defines on what the convergence check needs to be performed.
    - \* **'CDF'**, the convergence is checked in terms of probability (Cumulative Distribution Function). From a practical point of view, this means that full uncertain domain is discretized in a way that the probability volume of each cell is going to be equal to the tolerance specified in the body of the node `<Convergence>`
    - \* **'value'**, the convergence is checked on the hyper-volume in terms of variable values. From a practical point of view, this means that full uncertain domain is

discretized in a way that the “volume” fraction of each cell is going to be equal to the tolerance specified in the body of the node **<Convergence>**. In other words, each cell volume is going to be equal to the total volume times the tolerance.

*Default: CDF.*

- **persistence**, *optional integer attribute*, offers an additional convergence check. It represents the number of times the computed error needs to be below the inputted tolerance before convergence is reported.

*Default: 5.*

- **subGridTol**, *optional float attribute*, this attribute is used to activate the multi-grid approach (adaptive meshing) of the constructed evaluation grid (see attribute **weight**). In case this attribute is specified, the final grid discretization (cell’s “volume content” aka convergence confidence) is represented by the value here specified. The sampler converges on the initial coarse grid, defined by the tolerance specified in the body of the node **<Convergence>**. When the Limit Surface has been identified on the coarse grid, the sampler starts refining the grid until the “volume content” of each cell is equal to the value specified in this attribute (Multi-grid approach).

*Default: None.*

In summary, this XML node contains the information that is needed in order to control this sampler’s convergence criterion.

- **<batchStrategy>**, *string, optional field*, defines how points should be selected within a batch of size  $n$  where  $n$  is given by the **<maxBatchSize>** parameter below. Four options are available:
  - ‘**none**’ If this is specified then the **<maxBatchSize>** parameter below will be ignored and the functionality will replicate the LimitSurfaceSearch, in that the limit surface will be rebuilt and the points will be re-scored after each trial is completed.
  - ‘**naive**’ The top  $n$  candidates will be queued for adaptive sampling before retraining the limit surface and re-scoring the new candidate set.
  - ‘**maxP**’ The topology of the limit surface given the scoring function values will be decomposed and the top  $n$  highest topologically persistent features (local maxima) will be queued for adaptive sampling before retraining and re-scoring the new candidate set.
  - ‘**maxV**’ The topology of the limit surface given the scoring function values will be decomposed and the top  $n$  highest topological features (local maxima) will be queued for adaptive sampling before retraining and re-scoring the new candidate set.

*Default: none.*

- **<maxBatchSize>**, *integer, optional field*, specifies the number of points to select for adaptive sampling before retraining the limit surface and re-scoring the candidates. This is

the equivalent of the  $n$  parameter used in the `<batchStrategy>` description.

*Default: 1.*

- `<scoring>`, *string, optional field*, defines the scoring function to use on the candidate limit surface points in order to select the next adaptive point. Two options are available:
  - `'distance'` will scoring the candidate points by their distance to the closest realized point, in this way preference is given to unexplored regions of the limit surface.
  - `'distancePersistence'` augments the distance above by multiplying it with the inverse persistence of a candidate point which measures how many times the label of the candidate point has changed throughout the lifespan of the algorithm.

*Default: distancePersistence.*

- `<simplification>`, *float in the range [0,1], optional field*, specifies the percent of the scoring function range (on the candidate set) as the amount of topological simplification to do before extracting the topological features from the candidate set (local maxima). This only applies when the `<batchStrategy>` is set to `'maxP'` or `'maxV'`. Thus, one may end up with a batch size less than that specified by `<maxBatchSize>`.

*Default: 0.*

- `<thickness>`, *positive integer, optional field*, specifies how much the limit surface should be expanded (in terms of grid distance) when constructing a candidate set. A value of 1 implies only the points bounding the limit surface.

*Default: 1.*

- `<threshold>`, *float in the range [0,1], optional field*, once the candidates have been ranked and selected, before queueing them for adaptive sampling, this value is used to threshold any points whose score is less than this percentage of the scoring function range (on the candidate set). Thus, one may end up with a batch size less than that specified by `<maxBatchSize>`.

*Default: 0*

- **Assembler Objects** These objects are either required or optional depending on the functionality of the LimitSurfaceSearch Sampler. The objects must be listed with a rigorous syntax that, except for the XML node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to identify them within the simulation framework:

- `class`, *required string attribute*, the main “class” of the listed object. For example, it can be `'Models'`, `'Functions'`, etc.
- `type`, *required string attribute*, the object identifier or sub-type. For example, it can be `'ROM'`, `'External'`, etc.

The **LimitSurfaceSearch** approach requires or optionally accepts the following object types:

- **<Function>**, *string, required field*, the body of this XML block needs to contain the name of an external function object defined within the **<Functions>** main block (see Section 16). This object represents the boolean function that defines the transition boundaries. This function must implement a method called `__residuumSign(self)`, that returns either -1 or 1, depending on the system conditions (see Section 16).
- **<ROM>**, *string, optional field*, if used, the body of this XML node must contain the name of a ROM defined in the **<Models>** block (see Section 15.3). The ROM here specified is going to be used as “acceleration model” to speed up the convergence of the sampling strategy. The **<Target>** XML node in the ROM input block (within the **<Models>** section) needs to match the name of the goal **<Function>** (e.g. if the goal function is named “transitionIdentifier”, the **<Target>** of the ROM needs to report the same name: **<Target>transitionIdentifier<Target>**).
- **<TargetEvaluation>**, *string, required field*, represents the container where the system evaluations are stored. From a practical point of view, this XML node must contain the name of a data object defined in the **<DataObjects>** block (see Section 12). The Limit Surface Search sampling accepts “DataObjects” of type “Point” and “PointSet” only.

Example:

```

<Samplers>
...
<LimitSurfaceSearch name='LSSName'>
  <ROM class='Models' type='ROM'>ROMname</ROM>
  <Function class='Functions' type='External'
    >FunctionName</Function>
  <TargetEvaluation class='DataObjects'
    type='PointSet'>DataName</TargetEvaluation>
  <Convergence limit='3000' forceIteration='False'
    weight='CDF' subGridTol='1e-4' persistence='5'>
    1e-2
  </Convergence>
  <variable name='var1'>
    <distribution>aDistributionNameDefinedInDistributionBlock1
    </distribution>
  </variable>
  <variable name='var2'>
    <distribution>aDistributionNameDefinedInDistributionBlock2
    </distribution>
  </variable>
  <variable name='var3'>
    <distribution>aDistributionNameDefinedInDistributionBlock3
    </distribution>

```

```

    </variable>
</LimitSurfaceSearch>
...
</Samplers>

```

Batch sampling Example:

```

<Samplers>
...
<LimitSurfaceSearch name='LSBSName'>
  <ROM class='Models' type='ROM'>ROMname</ROM>
  <Function class='Functions' type='External'
    >FunctionName</Function>
  <TargetEvaluation class='DataObjects'
    type='PointSet'>DataName</TargetEvaluation>
  <Convergence limit='3000' forceIteration='False'
    weight='CDF' subGridTol='1e-4' persistence='5'>
    1e-2
  </Convergence>
  <scoring>distancePersistence</scoring>
  <batchStrategy>maxP</batchStrategy>
  <thickness>1</thickness>
  <maxBatchSize>4</maxBatchSize>
  <variable name='var1'>
    <distribution>aDistributionNameDefinedInDistributionBlock1
    </distribution>
  </variable>
  <variable name='var2'>
    <distribution>aDistributionNameDefinedInDistributionBlock2
    </distribution>
  </variable>
  <variable name='var3'>
    <distribution>aDistributionNameDefinedInDistributionBlock3
    </distribution>
  </variable>
</LimitSurfaceSearch>
...
</Samplers>

```

Associated External Python Module:

```

def __residuumSign(self) :
  if self.whatEverValue < self.OtherValue :

```

```
    return 1
else:
    return -1
```

### 11.3.2 Adaptive Dynamic Event Tree

The **Adaptive Dynamic Event Tree** approach is an advanced methodology employing a smart sampling around transition zones that determine a change in the status of the system (limit surface), using the support of a Dynamic Event Tree methodology. The main idea of the application of the previously explained adaptive sampling approach to the DET comes from the observation that the DET, when evaluated from a limit surface perspective, is intrinsically adaptive. For this reason, it appears natural to use the DET approach to perform a goal-function oriented pre-sampling of the input space.

RAVEN uses ROMs for predicting, in the input space, the location(s) of these transitions, in order to accelerate the exploration of the input space in proximity of the limit surface.

The specifications of this sampler must be defined within an **<AdaptiveDynamicEventTree>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.
- **printEndXML**, *optional string/boolean attribute*, this attribute controls the dumping of a “summary” of the DET performed in to an external XML.  
*Default: False.*
- **maxSimulationTime**, *optional float attribute*, this attribute controls the maximum “mission” time of the simulation underneath.  
*Default: None.*
- **mode**, *optional string attribute*, controls when the adaptive search needs to begin. Two options are available:
  - **'post'**, if this option is activated, the sampler first performs a standard Dynamic Event Tree analysis. At end of it, it uses the outcomes to start the adaptive search in conjunction with the DET support.
  - **'online'**, if this option is activated, the adaptive search starts at the beginning, during the initial standard Dynamic Event Tree analysis. Whenever a transition is detected, the **Adaptive Dynamic Event Tree** starts its goal-oriented search using the DET as support;

*Default: post.*

- **updateGrid**, *optional boolean attribute*, if true, each adaptive request is going to update the meshing of the initial DET grid.

*Default: True.*

In the **<AdaptiveDynamicEventTree>** input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive **<variable>** XML blocks:

- **<variable>**, *XML node, required parameter* will specify one attribute:
  - **name**, *required string attribute*, user-defined name of this variable.

This **<variable>** recognizes the following child nodes:

- **<distribution>**, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the **<Distributions>** block explained in Section 10. **Note:** Alternatively, this node must be omitted if the **<function>** node is supplied.
- **<function>**, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the **<Functions>** block explained in Section 16. **Note:** Alternatively, this node must be omitted if the **<distribution>** node is supplied.
- **<grid>**, *space separated floats, required field*, the content of this XML node depends on the definition of the associated attributes:
  - \* **type**, *required string attribute*, user-defined discretization metric type: 1) 'CDF', the grid will be specified based on cumulative distribution function probability thresholds, and 2) 'value', the grid will be provided using variable values.
  - \* **construction**, *required string attribute*, how the grid needs to be constructed, independent of its type (i.e. 'CDF' or 'value').

Based on the **construction** type, the content of the **<grid>** XML node and the requirements for other attributes change:

- \* **construction='equal'**. The grid is going to be constructed equally-spaced (**type='value'**) or equally probable (**type='CDF'**). This construction type requires the definition of additional attributes:
  - **steps**, *required integer attribute*, number of equally spaced/probable discretization steps.

This construction type requires that the content of the **<grid>** node represents the lower and upper bounds (either in probability or value). Two values need to



be specified; the lowest one will be considered as the *lowerBound*, the largest, the *upperBound*. The lower and upper bounds are checked against the associated **<distribution>** bounds. If one or both of them falls outside the distribution's bounds, the code will raise an error. The *stepSize* is determined as follows:  

$$stepSize = (upperBound - lowerBound) / steps$$

- \* **construction='custom'**. The grid will be directly specified by the user. No additional attributes are needed. This construction type requires that the **<grid>** node contains the actual mesh bins. For example, if the grid **type** is 'CDF', in the body of **<grid>**, the user will specify the CDF probability thresholds (nodalization in probability). All the bins are checked against the associated **<distribution>** bounds. If one or more of them falls outside the distribution's bounds, the code will raise an error.

**Note:** The **<grid>** node is only required if a **<distribution>** node is supplied. In the case of a **<function>** node, no grid information is requested.

In addition to the **<variable>** nodes, the main **<AdaptiveDynamicEventTree>** node needs to contain two supplementary sub-nodes:

- **<Convergence>**, *float, required field*, Convergence tolerance. The meaning of this tolerance depends on the definition of other attributes that might be defined in this XML node:
  - **limit**, *optional integer attribute*, the maximum number of adaptive samples (iterations).  
*Default: infinite.*
  - **forceIteration**, *optional boolean attribute*, this attribute controls if at least a number of iterations equal to **limit** must be performed.  
*Default: False.*
  - **weight**, *optional string attribute (case insensitive)*, defines on what the convergence check needs to be performed.
    - \* 'CDF', the convergence is checked in terms of probability (Cumulative Distribution Function). From a practical point of view, this means that full uncertain domain is discretized in a way that the probability volume of each cell is going to be equal to the tolerance specified in the body of the node **<Convergence>**
    - \* 'value', the convergence is checked on the hyper-volume in terms of variable values. From a practical point of view, this means that full uncertain domain is discretized in a way that the "volume" fraction of each cell is going to be equal to the tolerance specified in the body of the node **<Convergence>**. In other words, each cell volume is going to be equal to the total volume times the tolerance.

*Default: CDF.*



- **persistence**, *optional integer attribute*, offers an additional convergence check. It represents the number of times the computed error needs to be below the inputted tolerance before convergence is reported.  
*Default: 5.*
- **subGridTol**, *optional float attribute*, this attribute is used to activate the multi-grid approach (adaptive meshing) of the constructed evaluation grid (see attribute **weight**). In case this attribute is specified, the final grid discretization (cell’s “volume content” aka convergence confidence) is represented by the value here specified. The sampler converges on the initial coarse grid, defined by the tolerance specified in the body of the node **<Convergence>**. When the Limit Surface has been identified on the coarse grid, the sampler starts refining the grid until the “volume content” of each cell is equal to the value specified in this attribute (Multi-grid approach).  
*Default: None.*

In summary, this XML node contains the information that is needed in order to control this sampler’s convergence criterion.

- **Assembler Objects** These objects are either required or optional depending on the functionality of the AdaptiveDynamicEventTree Sampler. The objects must be listed with a rigorous syntax that, except for the XML node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to identify them within the simulation framework:
  - **class**, *required string attribute*, the main “class” of the listed object. For example, it can be **'Models'**, **'Functions'**, etc.
  - **type**, *required string attribute*, the object identifier or sub-type. For example, it can be **'ROM'**, **'External'**, etc.

The **AdaptiveDynamicEventTree** approach requires or optionally accepts the following object types:

- **<Function>**, *string, required field*, the body of this XML block needs to contain the name of an external function object defined within the **<Functions>** main block (see Section 16). This object represents the boolean function that defines the transition boundaries. This function must implement a method called `__residuumSign(self)`, that returns either -1 or 1, depending on the system conditions (see Section 16).
- **<ROM>**, *string, optional field*, if used, the body of this XML node must contain the name of a ROM defined in the **<Models>** block (see Section 15.3). The ROM here specified is going to be used as “acceleration model” to speed up the convergence of the sampling strategy. The **<Target>** XML node in the ROM input block (within the **<Models>** section) needs to match the name of the goal **<Function>** (e.g. if the goal function is named “transitionIdentifier”, the **<Target>** of the ROM needs to report the same name: **<Target>transitionIdentifier<Target>**).

- `<TargetEvaluation>`, *string, required field*, represents the container where the system evaluations are stored. From a practical point of view, this XML node must contain the name of a data object defined in the `<DataObjects>` block (see Section 12). The adaptive sampling accepts “DataObjects” of type “Point” and “PointSet” only.

Example:

```

<Samplers>
...
<AdaptiveDynamicEventTree name = 'AdaptiveName'>
  <ROM class = 'Models' type = 'ROM' ROMname</ROM>
  <Function class = 'Functions' type =
    'External'>FunctionName</Function>
  <TargetEvaluation class = 'DataObjects' type =
    'PointSet'>DataName</TargetEvaluation>
  <Convergence limit = '3000' subGridTol= '0.001'
    forceIteration = 'False' weight = 'CDF'
    subGriTol='''1e-5' persistence = '5'>
    1e-2
  </Convergence>
  <variable name = 'var1'>
    <distribution>
      aDistributionNameDefinedInDistributionBlock1
    </distribution>
    <grid type='CDF' construction='custom'>0.1 0.8</grid>
  </variable>
  <variable name = 'var2'>
    <distribution>
      aDistributionNameDefinedInDistributionBlock2
    </distribution>
    <grid type='CDF' construction='custom'>0.1 0.8</grid>
  </variable>
  <variable name = 'var3'>
    <distribution>
      aDistributionNameDefinedInDistributionBlock3
    </distribution>
    <grid type='CDF' construction='custom'>0.1 0.8</grid>
  </variable>
</AdaptiveDynamicEventTree>
...
</Samplers>

```

Associated External Python Module:

```
def __residuumSign(self):
    if self.whatEverValue < self.OtherValue:
        return 1
    else:
        return -1
```

### 11.3.3 Adaptive Hybrid Dynamic Event Tree

The **Adaptive Hybrid Dynamic Event Tree** approach is an advanced methodology employing a smart sampling around transition zones that determine a change in the status of the system (limit surface), using the support of the Hybrid Dynamic Event Tree methodology. Practically, this methodology represents a conjunction between the previously described Adaptive DET and the Hybrid DET method for the treatment of the epistemic variables.

Regarding the input requirements, the AHDET sampler is a “sub-type” of the `<AdaptiveDynamicEventTree>` sampler. For this reason, its specifications must be defined within a `<AdaptiveDynamicEventTree>` block.

The specifications of this sampler must be defined within an `<AdaptiveDynamicEventTree>` XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.
- **printEndXML**, *optional string/boolean attribute*, this attribute controls the dumping of a “summary” of the DET performed in to an external XML.  
*Default: False.*
- **maxSimulationTime**, *optional float attribute*, this attribute controls the maximum “mission” time of the simulation underneath.  
*Default: None.*
- **mode**, *optional string attribute*, controls when the adaptive search needs to begin. Two options are available:
  - **'post'**, if this option is activated, the sampler first performs a standard Dynamic Event Tree analysis. At end of it, it uses the outcomes to start the adaptive search in conjunction with the DET support.
  - **'online'**, if this option is activated, the adaptive search starts at the beginning, during the initial standard Dynamic Event Tree analysis. Whenever a transition is detected,

the **Adaptive Dynamic Event Tree** starts its goal-oriented search using the DET as support;

*Default: post.*

- **updateGrid**, *optional boolean attribute*, if true, each adaptive request is going to update the meshing of the initial DET grid.

*Default: True.*

In the **<AdaptiveDynamicEventTree>** input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive **<variable>** XML blocks:

- **<variable>**, *XML node, required parameter* will specify one attribute:
  - **name**, *required string attribute*, user-defined name of this variable.

This **<variable>** recognizes the following child nodes:

- **<distribution>**, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the **<Distributions>** block explained in Section 10. **Note:** Alternatively, this node must be omitted if the **<function>** node is supplied.
- **<function>**, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the **<Functions>** block explained in Section 16. **Note:** Alternatively, this node must be omitted if the **<distribution>** node is supplied.
- **<grid>**, *space separated floats, required field*, the content of this XML node depends on the definition of the associated attributes:
  - \* **type**, *required string attribute*, user-defined discretization metric type: 1) 'CDF', the grid will be specified based on cumulative distribution function probability thresholds, and 2) 'value', the grid will be provided using variable values.
  - \* **construction**, *required string attribute*, how the grid needs to be constructed, independent of its type (i.e. 'CDF' or 'value').

Based on the **construction** type, the content of the **<grid>** XML node and the requirements for other attributes change:

- \* **construction='equal'**. The grid is going to be constructed equally-spaced (**type='value'**) or equally probable (**type='CDF'**). This construction type requires the definition of additional attributes:
  - **steps**, *required integer attribute*, number of equally spaced/probable discretization steps.

This construction type requires that the content of the `<grid>` node represents the lower and upper bounds (either in probability or value). Two values need to be specified; the lowest one will be considered as the *lowerBound*, the largest, the *upperBound*. The lower and upper bounds are checked against the associated `<distribution>` bounds. If one or both of them falls outside the distribution's bounds, the code will raise an error. The *stepSize* is determined as follows:

$$stepSize = (upperBound - lowerBound) / steps$$

- \* `construction='custom'`. The grid will be directly specified by the user. No additional attributes are needed. This construction type requires that the `<grid>` node contains the actual mesh bins. For example, if the grid `type` is 'CDF', in the body of `<grid>`, the user will specify the CDF probability thresholds (nodalization in probability). All the bins are checked against the associated `<distribution>` bounds. If one or more of them falls outside the distribution's bounds, the code will raise an error.

**Note:** The `<grid>` node is only required if a `<distribution>` node is supplied. In the case of a `<function>` node, no grid information is requested.

In addition to the `<variable>` nodes, the main `<AdaptiveDynamicEventTree>` node needs to contain two supplementary sub-nodes:

- `<Convergence>`, *float, required field*, Convergence tolerance. The meaning of this tolerance depends on the definition of other attributes that might be defined in this XML node:
  - `limit`, *optional integer attribute*, the maximum number of adaptive samples (iterations).  
*Default: infinite.*
  - `forceIteration`, *optional boolean attribute*, this attribute controls if at least a number of iterations equal to `limit` must be performed.  
*Default: False.*
  - `weight`, *optional string attribute (case insensitive)*, defines on what the convergence check needs to be performed.
    - \* 'CDF', the convergence is checked in terms of probability (Cumulative Distribution Function). From a practical point of view, this means that full uncertain domain is discretized in a way that the probability volume of each cell is going to be equal to the tolerance specified in the body of the node `<Convergence>`
    - \* 'value', the convergence is checked on the hyper-volume in terms of variable values. From a practical point of view, this means that full uncertain domain is discretized in a way that the "volume" fraction of each cell is going to be equal to the tolerance specified in the body of the node `<Convergence>`. In other words, each cell volume is going to be equal to the total volume times the tolerance.

*Default: CDF.*

- **persistence**, *optional integer attribute*, offers an additional convergence check. It represents the number of times the computed error needs to be below the inputted tolerance before convergence is reported.

*Default: 5.*

- **subGridTol**, *optional float attribute*, this attribute is used to activate the multi-grid approach (adaptive meshing) of the constructed evaluation grid (see attribute **weight**). In case this attribute is specified, the final grid discretization (cell’s “volume content” aka convergence confidence) is represented by the value here specified. The sampler converges on the initial coarse grid, defined by the tolerance specified in the body of the node **<Convergence>**. When the Limit Surface has been identified on the coarse grid, the sampler starts refining the grid until the “volume content” of each cell is equal to the value specified in this attribute (Multi-grid approach).

*Default: None.*

In summary, this XML node contains the information that is needed in order to control this sampler’s convergence criterion.

- **Assembler Objects** These objects are either required or optional depending on the functionality of the AdaptiveDynamicEventTree Sampler. The objects must be listed with a rigorous syntax that, except for the XML node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to identify them within the simulation framework:

- **class**, *required string attribute*, the main “class” of the listed object. For example, it can be **'Models'**, **'Functions'**, etc.
- **type**, *required string attribute*, the object identifier or sub-type. For example, it can be **'ROM'**, **'External'**, etc.

The **AdaptiveDynamicEventTree** approach requires or optionally accepts the following object types:

- **<Function>**, *string, required field*, the body of this XML block needs to contain the name of an external function object defined within the **<Functions>** main block (see Section 16). This object represents the boolean function that defines the transition boundaries. This function must implement a method called `__residuumSign(self)`, that returns either -1 or 1, depending on the system conditions (see Section 16).
- **<ROM>**, *string, optional field*, if used, the body of this XML node must contain the name of a ROM defined in the **<Models>** block (see Section 15.3). The ROM here specified is going to be used as “acceleration model” to speed up the convergence of the sampling strategy. The **<Target>** XML node in the ROM input block (within the **<Models>** section) needs to match the name of the goal **<Function>** (e.g. if

the goal function is named “transitionIdentifier”, the **<Target>** of the ROM needs to report the same name: **<Target>transitionIdentifier<Target>**).

- **<TargetEvaluation>**, *string, required field*, represents the container where the system evaluations are stored. From a practical point of view, this XML node must contain the name of a data object defined in the **<DataObjects>** block (see Section 12). The adaptive sampling accepts “DataObjects” of type “Point” and “PointSet” only.

As it can be noticed, the basic specifications of the Adaptive Hybrid Dynamic Event Tree method are consistent with the ones for the ADET methodology. In order to activate the **Adaptive Hybrid Dynamic Event Tree** sampler, the main **<AdaptiveDynamicEventTree>** block needs to contain an additional sub-node called **<HybridSampler>**. This sub-node needs to contain the following attribute:

- **type**, *required string attribute*, type of pre-sampling strategy to be used. Up to now only one option is available:
  - **'LimitSurface'**. With this option, the epistemic variables here listed are going to be part of the LS search. This means that the discretization of the domain of these variables is determined by the **<Convergece>** node.

Independent of the type of HybridSampler that has been specified, the **<HybridSampler>** must contain the variables that need to be sampled. As already mentioned, these variables are specified within consecutive **<variable>** XML blocks:

- **<variable>**, *XML node, required parameter* will specify one attribute:
  - **name**, *required string attribute*, user-defined name of this variable.

This **<variable>** recognizes the following child nodes:

- **<distribution>**, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the **<Distributions>** block explained in Section 10. **Note:** Alternatively, this node must be omitted if the **<function>** node is supplied.
- **<function>**, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the **<Functions>** block explained in Section 16. **Note:** Alternatively, this node must be omitted if the **<distribution>** node is supplied.

Example:



```

<Samplers>
...
<AdaptiveDynamicEventTree name = 'AdaptiveName'>
  <ROM class = 'Models' type = 'ROM'ROMname</ROM>
  <Function class = 'Functions' type =
    'External'>FunctionName</Function>
  <TargetEvaluation class = 'DataObjects' type =
    'PointSet'>DataName</TargetEvaluation>
  <Convergence limit = '3000' subGridTol= '0.001'
    forceIteration = 'False' weight = 'CDF'
    subGriTol='''1e-5' persistence = '5'>
    1e-2
  </Convergence>
  <HybridSampler type='LimitSurface'>
    <variable name = 'epistemicVar1'>
      <distribution>
        aDistributionNameDefinedInDistributionBlock1
      </distribution>
    </variable>
    <variable name = 'epistemicVar2'>
      <distribution>
        aDistributionNameDefinedInDistributionBlock2
      </distribution>
    </variable>
  </HybridSampler>
  <variable name = 'var1'>
    <distribution>
      aDistributionNameDefinedInDistributionBlock3
    </distribution>
    <grid type='CDF' construction='custom'>0.1 0.8</grid>
  </variable>
  <variable name = 'var2'>
    <distribution>
      aDistributionNameDefinedInDistributionBlock4
    </distribution>
    <grid type='CDF' construction='custom'>0.1 0.8</grid>
  </variable>
  <variable name = 'var3'>
    <distribution>
      aDistributionNameDefinedInDistributionBlock5
    </distribution>
    <grid type='CDF' construction='custom'>0.1 0.8</grid>

```



```
</variable>

</AdaptiveDynamicEventTree>
...
</Samplers>
```

Associated External Python Module:

```
def __residuuumSign(self) :
    if self.whatEverValue < self.OtherValue:
        return 1
    else:
        return -1
```

### 11.3.4 Adaptive Sparse Grid

The **Adaptive Sparse Grid** approach is an advanced methodology that employs an intelligent search for the most suitable sparse grid quadrature to characterize a model. To perform such sampling, RAVEN adaptively builds an index set and generates sparse grids in a similar manner to Sparse Grid Collocation samplers. In each iterative step, the adaptive index set determines the next possible quadrature orders to add in each dimension, and determines the index set point that would offer the largest impact to one of the convergence metrics. This process continues until the total impact of all the potential index set points is less than tolerance. For many models, this function converges after fewer runs than a traditional Sparse Grid Collocation sampling. However, it should be noted that this algorithm fails in the event that the partial derivative of the response surface with respect to any single input dimension is zero at the origin of the input domain. For example, the adaptive algorithm fails for the model  $f(x) = x \cdot y$ .

The specifications of this sampler must be defined within an **<Adaptive Sparse Grid>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

In the **<Adaptive Sparse Grid>** input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive **<variable>** XML blocks:

- **<variable>**, *XML node, required parameter* will specify one attribute:

- **name**, *required string attribute*, user-defined name of this variable.

This **<variable>** recognizes the following child nodes:

- **<distribution>**, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the **<Distributions>** block explained in Section 10. **Note:** Alternatively, this node must be omitted if the **<function>** node is supplied.
- **<function>**, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the **<Functions>** block explained in Section 16. **Note:** Alternatively, this node must be omitted if the **<distribution>** node is supplied.

In addition to the **<variable>** nodes, the main XML node **<AdaptiveSparseGrid>** needs to contain the following supplementary sub-nodes:

- **<Convergence>**, *float, required field*, Convergence tolerance. The meaning of this tolerance depends on the **target** attribute of this node.
  - **target**, *required string attribute*, the metric for convergence. The following metrics are available: 'variance', which converges the sparse quadrature integration of the second moment of the model.
  - **maxPolyOrder**, *optional integer attribute*, limits the maximum size equivalent polynomial for any one dimension.  
*Default: 10.*
  - **persistence**, *optional integer attribute*, defines the number of index set points that are required to be found before calculation can exit. Setting this to a higher value can help if the adaptive process is not finding significant indices on its own.  
*Default: 2.*

In summary, this XML node contains the information that is needed in order to control this sampler's convergence criterion.

- **<convergenceStudy>**, *optional node*, if included, triggers writing state points at particular numbers of model solves for the purpose of a convergence study. The study is performed by writing XML output files as described in the OutStreamManager for ROMs at the state points requested, using 'all' as the requested **<what>** values. The state points are identified when a certain number of model runs is passed, as specified by the **<runStatePoints>** node. This node has the following sub-nodes to define its parameters:

- **<runStatePoints>**, *list of integers, required node*, lists the number of model runs at which state points should be written. Note that these will be written when the requested number of runs is met or passed, so the actual value is often somewhat more than the requested value, and the exact value will be listed in the XML output.
- **<baseFilename>**, *string, optional node*, if specified determines the base file name for the state point outputs. If not specified, defaults to 'out\_'.
- **<pickle>**, *no text, optional node*, if this node is included, serialized (pickled) versions of the ROM at each of the run states is also created in the working directory, with the format <baseFilename><numRuns>.pk, such as out\_100.pk.
- **<logfile>**, *optional node*, if included, the log file onto which the adaptive step progress can be printed. The log includes the values of included polynomial coefficients as well as the expected impacts of polynomial coefficients not yet included. This is different from the convergenceStudy print, which will give statistical moments at certain steps.
- **<maxRuns>**, *optional node*, if included, the adaptive sampler will track the number of computational solves necessary to construct the associated GaussPolynomialROM. If at any point the number of solves exceeds the value given, it will not initiate any additional solves, and will exit when existing solves finish.

**Assembler Objects** These objects are either required or optional depending on the functionality of the Adaptive Sparse Grid Sampler. The objects must be listed with a rigorous syntax that, except for the XML node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to identify them within the simulation framework:

- **class**, *required string attribute*, the main “class” of the listed object. For example, it can be 'Models', 'Functions', etc.
- **type**, *required string attribute*, the object identifier or sub-type. For example, it can be 'ROM', 'External', etc.

The **Adaptive Sparse Grid** approach requires or optionally accepts the following object types:

- **<ROM>**, *string, required field*, the body of this XML node must contain the name of an appropriate ROM defined in the **<Models>** block (see Section 15.3).
- **<TargetEvaluation>**, *string, required field*, represents the container where the system evaluations are stored. From a practical point of view, this XML node must contain the name of a data object defined in the **<DataObjects>** block (see Section 12). The Adaptive Sparse Grid sampling accepts “DataObjects” of type “Point” and “PointSet” only.

Example:

```
<Samplers>
...
<AdaptiveSparseGrid name="ASG" verbosity='debug'>
  <Convergence target='coeffs'>1e-2</Convergence>
  <variable name="x1">
    <distribution>UniDist</distribution>
  </variable>
  <variable name="x2">
    <distribution>UniDist</distribution>
  </variable>
  <ROM class = 'Models' type = 'ROM'>gausspolyrom</ROM>
  <TargetEvaluation class = 'DataObjects' type =
    'PointSet'>solns</TargetEvaluation>
</AdaptiveSparseGrid>
...
</Samplers>
```

Like in the **SparseGridCollocation** sampler, if the covariance matrix is provided and the input parameters are assumed to have the multivariate normal distribution, the **AdaptiveSparseGrid** can be also used. This means one creates the sparse grids of variables listed by **<latentVariables>** in the transformed space. If this is the case, the user needs to provide additional information, i.e. the **<transformation>** under **<MultivariateNormal>** of **<Distributions>** (more information can be found in Section 10.2). In addition, the node **<variablesTransformation>** is also required for **AdaptiveSparseGrid** sampler. This node is used to transform the variables specified by **<latentVariables>** in the transformed space of input into variables specified by **<manifestVariables>** in the input space. The variables listed in **<latentVariables>** should be predefined in **<variable>**, and the variables listed in **<manifestVariables>** are used by the **<Models>**.

- **<variablesTransformation>**, *optional field*. this XML node accepts one attribute:
  - **model**, *required string attribute*, the name for the model defined in the XML node **<Models>**. This attribute indicates the **<manifestVariables>** defined in the subnode will be used by the **model**.

In addition, this XML node also accepts three children nodes:

- **<latentVariables>**, *comma separated string, required field*, user-defined latent variables that are used for the variables transformation. All the variables listed under this node should be also mentioned in **<variable>**.
- **<manifestVariables>**, *comma separated string, required field*, user-defined manifest variables that are used by the **model**.

- `<method>`, *string, required field*, the method that is used for the variables transformation. The currently available method is 'pca'.

```

...
<Models>
  ...
  <ExternalModel ModuleToLoad="lorenzAttractor_noK"
    name="PythonModule" subType="">
    <variable>sigma</variable>
    <variable>rho</variable>
    <variable>beta</variable>
    <variable>x</variable>
    <variable>y</variable>
    <variable>z</variable>
    <variable>time</variable>
    <variable>x0</variable>
    <variable>y0</variable>
    <variable>z0</variable>
  </ExternalModel>
  <ROM name="gausspolyrom" subType="GaussPolynomialRom">
    <Target>ans</Target>
    <Features>x1,y1,z1</Features>
    <IndexSet>TensorProduct</IndexSet>
    <PolynomialOrder>1</PolynomialOrder>
  </ROM>
  ...
</Models>

<Distributions>
  ...
  <MultivariateNormal name='MVNDist' method='pca'>
    <transformation>
      <rank>3</rank>
    </transformation>
    <mu>0.0 1.0 2.0</mu>
    <covariance type="abs">
      1.0      0.6      -0.4
      0.6      1.0      0.2
      -0.4     0.2      0.8
    </covariance>
  </MultivariateNormal>
  ...
</Distributions>

```

```

<Samplers>
  ...
  <AdaptiveSparseGrid name='ASC'>
    <variable name='x0'>
      <distribution dim='1'>MVNDist</distribution>
    </variable>
    <variable name='y0'>
      <distribution dim='2'>MVNDist</distribution>
    </variable>
    <variable name='z0'>
      <distribution dim='3'>MVNDist</distribution>
    </variable>
    <variablesTransformation model="PythonModule">
      <latentVariables>x1, y1, z1</latentVariables>
      <manifestVariables>x0, y0, z0</manifestVariables>
      <method>pca</method>
    </variablesTransformation>
    <ROM class = 'Models' type = 'ROM'>gausspolyrom</ROM>
    <TargetEvaluation class = 'DataObjects' type =
      'PointSet'>sols</TargetEvaluation>
  </AdaptiveSparseGrid>
  ...
</Samplers>
  ...

```

### 11.3.5 Adaptive Sobol Decomposition

The **Adaptive Sobol Decomposition** approach is an advanced methodology that decomposes an uncertainty space into subsets and adaptively includes the most influential ones. For example, for a response function  $f(a, b, c)$ , the full list of subsets include  $(a)$ ,  $(b)$ ,  $(c)$ ,  $(a, b)$ ,  $(a, c)$ ,  $(b, c)$ ,  $(a, b, c)$ . A Gauss Polynomial ROM is constructed for each included subset using the Adaptive Sparse Grid sampler. The importance of each subset is estimated based on the importance of preceding subsets; that is, the impact of  $(a, b)$  on the representation of  $f$  is estimated using the impact of  $(a)$  and  $(b)$ . Because of the excellent performance of Gauss Polynomial ROMs for small-dimension spaces, this sampler used to construct an HDMR ROM can be very efficient. Note that the ROM specified for this sampler *must* be an HDMRRom specified in the Models block.

The specifications of this sampler must be defined within an **<Adaptive Sobol>** XML block. This XML node accepts one attribute:

- **name**, *required string attribute*, user-defined name of this sampler. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.

In the `<Adaptive Sobol>` input block, the user needs to specify the variables to sample. As already mentioned, these variables are specified within consecutive `<variable>` XML blocks:

- `<variable>`, *XML node, required parameter* will specify one attribute:
  - **name**, *required string attribute*, user-defined name of this variable.

This `<variable>` recognizes the following child nodes:

- `<distribution>`, *string, required field*, name of the distribution that is associated to this variable. Its name needs to be contained in the `<Distributions>` block explained in Section 10. **Note:** Alternatively, this node must be omitted if the `<function>` node is supplied.
- `<function>`, *string, required field*, name of the function that defines the calculation of this variable from other distributed variables. Its name needs to be contained in the `<Functions>` block explained in Section 16. **Note:** Alternatively, this node must be omitted if the `<distribution>` node is supplied.

In addition to the `<variable>` nodes, the main XML node `<AdaptiveSobol>` needs to contain the following supplementary sub-nodes:

- `<Convergence>`, *required node*, Convergence properties. This node contains the following properties that can be set by sub-nodes:
  - `<relTolerance>`, *required float*, the relative tolerance to converge. This will compare to the estimate of subset polynomial errors and additional subset polynomials over the variance of the expansion so far to determine convergence.
  - `<maxRuns>`, *optional integer field*, a limit for the number of model calls. Once this limit is reached, no additional subsets will be generated or considered; however, existing subsets will continue to be trained. If not specified, no limit on solves is imposed.
  - `<maxSobolOrder>`, *optional integer field*, the largest polynomial orders to use in subset GaussPolynomialRom objects. If specified, polynomial indices with a value larger than the value given will be rejected during adaptive construction.
  - `<progressParam>`, *optional float field*, a favoritism parameter ranging between 0 and 2. At 0, the algorithm will always prefer adding polynomials to adding new subsets in the HDMR expansion. At 2, the opposite is true. Default is 1.



- **<logFile>**, *optional string field*, a file to which adaptive progress is recorded. If specified, each adaptive step will trigger printing progress to the file given, including the estimated error at the step, the next adaptive step to take, the coefficient of each polynomial within each gPC expansion, and the actual and expected Sobol sensitivities of each HDMR subset. Default is no printing.
- **<subsetVerbosity>**, *optional string field*, the verbosity for components constructed during the adaptive HDMR process. Options are *silent*, *quiet*, *all*, or *debug*, in order of verbosity. If an invalid entry is provided, will resort to default. Default is *quiet*.

In summary, this XML node contains the information that is needed in order to control this sampler’s convergence criterion.

- **<convergenceStudy>**, *optional node*, if included, triggers writing state points at particular numbers of model solves for the purpose of a convergence study. The study is performed by writing XML output files as described in the OutStreamManager for ROMs at the state points requested, using *'all'* as the requested **<what>** values. The state points are identified when a certain number of model runs is passed, as specified by the **<runStatePoints>** node. This node has the following sub-nodes to define its parameters:
  - **<runStatePoints>**, *list of integers, required node*, lists the number of model runs at which state points should be written. Note that these will be written when the requested number of runs is met or passed, so the actual value is often somewhat more than the requested value, and the exact value will be listed in the XML output.
  - **<baseFilename>**, *string, optional node*, if specified determines the base file name for the state point outputs. If not specified, defaults to *'out\_.'*
  - **<pickle>**, *no text, optional node*, if this node is included, serialized (pickled) versions of the ROM at each of the run states is also created in the working directory, with the format `<baseFilename><numRuns>.pk`, such as `out_100.pk`.

**Assembler Objects** These objects are either required or optional depending on the functionality of the AdaptiveSobol Sampler. The objects must be listed with a rigorous syntax that, except for the XML node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to identify them within the simulation framework:

- **class**, *required string attribute*, the main “class” of the listed object. For example, it can be *'Models'*, *'Functions'*, etc.
- **type**, *required string attribute*, the object identifier or sub-type. For example, it can be *'ROM'*, *'External'*, etc.



The **AdaptiveSobol** approach requires or optionally accepts the following object types:

- **<ROM>**, *string, required field*, the body of this XML node must contain the name of an appropriate ROM defined in the **<Models>** block (see Section 15.3).
- **<TargetEvaluation>**, *string, required field*, represents the container where the system evaluations are stored. From a practical point of view, this XML node must contain the name of a data object defined in the **<DataObjects>** block (see Section 12). The Adaptive Sobol sampling accepts “DataObjects” of type “Point” and “PointSet” only.

Example:

```
<Samplers>
...
<AdaptiveSobol name="AS" verbosity='debug'>
  <Convergence>
    <relTolerance>1e-5</relTolerance>
    <maxRuns>150</maxRuns>
    <maxSobolOrder>3</maxSobolOrder>
    <progressParam>1</progressParam>
    <logFile>progress.txt</logFile>
    <subsetVerbosity>silent</subsetVerbosity>
  </Convergence>
  <variable name="x1">
    <distribution>UniDist</distribution>
  </variable>
  <variable name="x2">
    <distribution>UniDist</distribution>
  </variable>
  <ROM class = 'Models' type = 'ROM'>hdmrrom</ROM>
  <TargetEvaluation class = 'DataObjects' type =
    'PointSet'>solns</TargetEvaluation>
</AdaptiveSobol>
...
</Samplers>
```

## 12 DataObjects

As seen in the previous chapters, different entities in the RAVEN code interact with each other in order to create, ideally, an infinite number of different calculation flows. These interactions are made possible through a data handling system that each entity understands. This system, neglecting the grammar imprecision, is called the “DataObjects” system.

The `<DataObjects>` tag is a container of data objects of various types that can be constructed during the execution of a particular calculation flow. These data objects can be used as input or output for a particular **Model** (see Roles’ meaning in section 15), etc. Currently, RAVEN supports 4 different data types, each with a particular conceptual meaning. These data types are instantiated as sub-nodes in the `<DataObjects>` block of an input file:

- `<Point>`, as the name suggests, describes the state of the system at a certain point (e.g. in time). In other words, it can be considered a mapping between a set of parameters in the input space and the resulting outcomes in the output space at a particular point in the phase space (e.g. in time).
- `<PointSet>` is a collection of individual *Point* objects. It can be considered a mapping between multiple sets of parameters in the input space and the resulting sets of outcomes in the output space at a particular point (e.g. in time).
- `<History>` describes the temporal evolution of the state of the system within a certain input domain.
- `<HistorySet>` is a collection of individual *History* objects. It can be considered a mapping between multiple sets of parameters in the input space and the resulting sets of temporal evolutions in the output space.

As noted above, each data object represents a mapping between a set of parameters and the resulting outcomes. The data objects are defined within the main XML block called `<DataObjects>`:

```
<Simulation>
...
<DataObjects>
  <Point name='***'>...</Point>
  <PointSet name='***'>...</PointSet>
  <History name='***'>...</History>
  <HistorySet name='***'>...</HistorySet>
</DataObjects>
...
</Simulation>
```

Independent of the type of data, the respective XML node has the following available attributes:

- **name**, *required string attribute*, is a user-defined identifier for this data object. **Note:** As with other objects, this name can be used to refer to this specific entity from other input blocks in the XML.
- **historyName**, *optional string attribute*, in the case the “DataObject” type is either a *Point* or an *History*, this XML node contains the name of the associated history needs to be placed in.  
**Note:** This attribute is a required parameter when the `<Point>` or `<History>` types are constructed from a Database (e.g.HDF5).
- **hierarchical**, *optional boolean attribute*, if True this data object will be constructed, if possible, in a hierarchical fashion.  
*Default: False*

In each XML node (e.g. `<Point>` or `<HistorySet>`), the user needs to specify the following sub-nodes:

- `<Input>`, *comma separated string, required field* lists the input parameters to which this data is connected.
- `<Output>`, *comma separated string, required field* lists the output parameters to which this data is connected.

In addition to the XML nodes `<Input>` and `<Output>` explained above, the user can optionally specify a XML node named `<options>`. The `<options>` node can contain the following optional XML sub-nodes:

- `<inputRow>`, *integer, optional field*, used to specify the row (in the CSV file or HDF5 table) from which the input space needs to be retrieved (e.g. the time-step);
- `<outputRow>`, *integer, optional field*, used to specify the row (in the CSV file or HDF5 table) from which the output space needs to be retrieved (e.g. the time-step). If this node is inputted, the nodes `<operator>` and `<outputPivotValue>` can not be inputted (mutually exclusive).  
**Note:** This XML node is available for DataObjects of type `<Point>` and `<PointSet>` only;
- `<operator>`, *string, optional field*, is aimed to perform simple operations on the data to be stored. The 3 options currently available are:
  - ‘max’
  - ‘min’

- 'average'

If this node is inputted, the nodes `<outputRow>` and `<outputPivotValue>` can not be inputted (mutually exclusive).

**Note:** This XML node is available for DataObjects of type `<Point>` and `<PointSet>` only;

- `<pivotParameter>`, *string, optional field* the name of the parameter whose values need to be used as reference for the values specified in the XML nodes `<inputPivotValue>`, `<outputPivotValue>`, or `<inputPivotValue>` (if inputted). This field can be used, for example, if the driven code output file uses a different name for the variable “time” or to specify a different reference parameter (e.g. PRESSURE). Default value is 'time' .  
**Note:** The variable specified here should be monotonic; the code does not check for eventual oscillation and is going to take the first occurrence for the values specified in the XML nodes `<inputPivotValue>`, `<outputPivotValue>`, and `<inputPivotValue>`;
- `<inputPivotValue>`, *float, optional field*, the value of the `<pivotParameter>` at which the input space needs to be retrieved. If this node is inputted, the node `<inputRow>` can not be inputted (mutually exclusive).
- `<outputPivotValue>`. This node can be either a float or a list of floats, depending on the type of DataObjects:
  - if `<History>` or `<HistorySet>`, `<outputPivotValue>`, *list of floats, optional field*, list of values of the `<pivotParameter>` at which the output space needs to be retrieved;
  - if `<Point>` or `<PointSet>`, `<outputPivotValue>`, *float, optional field*, the value of the `<pivotParameter>` at which the output space needs to be retrieved. If this node is inputted, the node `<outputRow>` can not be inputted (mutually exclusive);

It needs to be noticed that if the optional nodes in the block `<options>` are not inputted, the following default are applied:

- the Input space is retrieved from the first row in the CSVs files or HDF5 tables (if the parameters specified are not among the variables sampled by RAVEN);
- the output space defaults are as follows:
  - \* if `<Point>` or `<PointSet>`, the output space is retrieved from the last row in the CSVs files or HDF5 tables;
  - \* if `<History>` or `<HistorySet>`, the output space is represented by all the rows found in the CSVs or HDF5 tables.

```

<DataObjects>
  <PointSet name='outTPS1'>
    <options>
      <inputRow>1</inputRow>
      <outputRow>-1</outputRow>
    </options>
    <Input>pipe_Area,pipe_Dh,Dummy1</Input>
    <Output>pipe_Hw,pipe_Tw,time</Output>
  </PointSet>
  <PointSet name='outTPS2'>
    <options>
      <inputPivotValue>0.00011</inputPivotValue>
      <pivotParameter>time</pivotParameter>
      <outputPivotValue>0.00012345</outputPivotValue>
    </options>
    <Input>pipe_Area,pipe_Dh,Dummy1</Input>
    <Output>pipe_Hw,pipe_Tw,time</Output>
  </PointSet>
  <Point name='outTP'>
    <options>
      <inputRow>-1</inputRow>
      <operator>max</operator>
    </options>
    <Input>pipe_Area,pipe_Dh,Dummy1</Input>
    <Output>pipe_Hw,pipe_Tw,time</Output>
  </Point>
  <HistorySet name='stories1'>
    <options>
      <inputRow>1</inputRow>
    </options>
    <Input>pipe_Area,pipe_Dh</Input>
    <Output>pipe_Hw,pipe_Tw,time</Output>
  </HistorySet>
  <HistorySet name='stories2'>
    <options>
      <inputRow>-1</inputRow>
      <pivotParameter>time</pivotParameter>
      <outputPivotValue>0.0002 0.0003
        0.0004</outputPivotValue>
    </options>
    <Input>pipe_Area,pipe_Dh</Input>
    <Output>pipe_Hw,pipe_Tw,time</Output>

```

```
</HistorySet>
<History name='hst1'>
  <options>
    <inputRow>1</inputRow>
  </options>
  <Input>pipe_Area,pipe_Dh</Input>
  <Output>pipe_Hw,pipe_Tw,time</Output>
</History>
<History name='hst2'>
  <options>
    <inputRow>-1</inputRow>
    <pivotParameter>time</pivotParameter>
    <outputPivotValue>0.0002 0.0003
      0.0004</outputPivotValue>
  </options>
  <Input>pipe_Area,pipe_Dh</Input>
  <Output>pipe_Hw,pipe_Tw,time</Output>
</History>
</DataObjects>
```

## 13 Databases

The RAVEN framework provides the capability to store and retrieve data to/from an external database. Currently RAVEN has support for only a database type called **HDF5**. This database, depending on the data format it is receiving, will organize itself in a “parallel” or “hierarchical” fashion. The user can create as many database objects as needed. The Database objects are defined within the main XML block called **<Databases>**:

```
<Simulation>
...
  <Databases>
    ...
    <HDF5 name="aDatabaseName1" />
    <HDF5 name="aDatabaseName2" />
    ...
  </Databases>
  ...
</Simulation>
```

The specifications of each Database of type HDF5 needs to be defined within the XML block **<HDF5>**, that recognizes the following attributes:

- **name**, *required string attribute*, a user-defined identifier of this object. **Note:** As with other objects, this name can be used to reference this specific entity from other input blocks in the XML.
- **directory**, *optional string attribute*, this attribute can be used to specify a particular directory path where the database will be created, if no *filename* is specified, or where to find an already existing one if *filename* is provided.  
*Default: workingDir/framework/DatabaseStorage.* The **<workingDir>** is the one defined within the **<RunInfo>** XML block (see Section 8).
- **filename**, *optional string attribute*, specifies the filename of an HDF5 that already exists in the **directory**. This is the only way to let RAVEN know that an HDF5 should be opened and not overwritten. **Note:** When this attribute is not specified, the newer database filename will be named *name.h5*, where *name* corresponds to the **name** attribute of this object.  
*Default: None*
- **compression**, *optional string attribute*, compression algorithm to be used. Available are:
  - ‘**gzip**’, best where portability is required. Good compression, moderate speed.
  - ‘**lzf**’, Low to moderate compression, very fast.

*Default: None*

```
<Databases>
  <HDF5 name="aDatabaseName1" directory='path_to_a_dir'
        compression='lzf' />
  <HDF5 name="aDatabaseName2" filename='aDatabaseName2.h5' />
</Databases>
```



## 14 OutStream system

The PRA and UQ framework provides the capabilities to visualize and dump out the data that are generated, imported (from a system code) and post-processed during the analysis. These capabilities are contained in the “OutStream” system. Actually, two different OutStream types are available:

- **Print**, module that lets the user dump the data contained in the internal objects;
- **Plot**, module, based on Matplotlib [3], aimed to provide advanced plotting capabilities.

Both the types listed above accept as “input” a *DataObjects* object type. This choice is due to the “*DataObjects*” system (see section 12) having the main advantage of ensuring a standardized approach for exchanging the data/meta-data among the different framework entities. Every module can project its outcomes into a *DataObjects* object. This provides the user with the capability to visualize/dump all the modules’ results. Additionally, the **Print** system can accept a ROM and inquire some of its specialized methods. As already mentioned, the RAVEN framework input is based on the eXtensible Markup Language (**XML**) format. Thus, in order to activate the “*OutStream*” system, the input needs to contain a block identified by the `<OutStreamManager>` tag (as shown below).

```
<OutStreamManager>  
  <!-- "OutStream" objects that need to be created-->  
</OutStreamManager>
```

In the “OutStreamManager” block an unlimited number of “Plot” and “Print” sub-blocks can be specified. The input specifications and the main capabilities for both types are reported in the following sections.

### 14.1 Printing system

The Printing system has been created in order to let the user dump the data, contained in the internal data objects (see Section 12), out at anytime during the calculation. Additionally, the user can inquire special methods of a **ROM** after training it, through a printing step. Currently, the only available output is a Comma Separated Value (**CSV**) file for **DataObjects**, and **XML** for **ROM** objects. This will facilitate the exchanging of results and provide the possibility to dump the solution of an analysis and “restart” another one constructing a data object from scratch, as well as access advanced features of particular reduced order models.

### 14.1.1 DataObjects Printing

The XML code, that is reported below, shows different ways to request a *Print* OutStream for **DataObjects**. The user needs to provide a name for each sub-block (XML attribute). These names are then used in the *Step* blocks to activate the Printing keywords at any time. The XML node has the following available attributes:

- **name**, *required string attribute*, is a user-defined identifier for this data object. **Note:** As with other objects, this name can be used to refer to this specific entity from other input blocks in the XML.

As shown in the examples below, every **<Print>** block must contain, at least, the two required tags:

- **<type>**, the output file type (csv or xml). **Note:** Only **csv** is currently available for **<DataObjects>**
- **<source>**, the *Data* name (one of the *Data* items defined in the **<DataObjects>** block).

If only these two tags are provided (as in the “first-example” below), the output file will be filled with the whole content of the “d-name” *Data* object.

```
<OutStreamManager>
  <Print name='first-example'>
    <type>csv</type>
    <source>d-name</source>
  </Print>
  <Print name='second-example'>
    <type>csv</type>
    <source>d-name</source>
    <what>Output</what>
  </Print>
  <Print name='third-example'>
    <type>csv</type>
    <source>d-name</source>
    <what>Input</what>
  </Print>
  <Print name='fourth-example'>
    <type>csv</type>
    <source>d-name</source>
    <what>Input | var-name-in, Output | var-name-out</what>
  </Print>
</OutStreamManager>
```

If just part of the `<source>` is important for a particular analysis, the additional XML tag `<what>` can be provided. In this block, the variables that need to be dumped must be specified, in comma separated format. The available options, for the `<what>` sub-block, are listed below:

- **Output**, the output space will be dumped out (see “second-example”)
- **Input**, the input space will be dumped out (see “third-example”)
- **Input—var-name-in/Output—var-name-out**, only the particular variables “var-name-in” and “var-name-out” will be reported in the output file (see “fourth-example”)

Note all of the XML tags are case-sensitive but not their content.

### 14.1.2 ROM Printing

While all **ROMs** in RAVEN are designed to be used as surrogate models, some **ROMs** additionally offer information about the original model that isn’t accessible through another means. For instance, **HDMRRom** objects can calculate sensitivity coefficients for subsets of the input domain. The XML code shown below demonstrates the methods to request these features from a **ROM**. The user needs to provide a `<name>` for each sub-block (XML attribute). These names are then used in the *Step* blocks to activate the Printing keywords at any time. As shown in the examples below, every `<Print>` block for ROMs must contain, at least, the three required tags

- `<type>`, the output file type (csv or xml). **Note:** Only **xml** is currently available for ROMs
- `<source>`, the *ROM* name (one of the `<ROM>` items defined in the `<Models>` block).
- `<what>`, the comma-separated list of desired metrics. The list of metrics available in each ROM is listed under that ROM type in Section 15.3. Alternatively, the keyword ‘**all**’ can be provided to request all available metrics, if any.

Additionally, when printing ROMs one optional node is available,

- `<target>`, the ROM target for which to inquire data

ROM printing uses the same naming conventions as DataObjects printing. Examples:

```
<OutputStreamManager>
  <Print name='first-ROM-example'>
    <type>xml</type>
    <source>mySobolRom</source>
    <what>all</what>
  </Print>
  <Print name='second-ROM-example'>
    <type>xml</type>
    <source>myGaussPolyRom</source>
    <what>mean,variance</what>
  </Print>
</OutputStreamManager>
```

## 14.2 Plotting system

The Plotting system provides all the capabilities to visualize the analysis outcomes, in real-time or as a post-processing stage. The system is based on the Python library Matplotlib [3]. Matplotlib is a 2D/3D plotting library which produces publication quality figures in a variety of hardcopy formats and interactive environments across platforms. This external tool has been wrapped in the RAVEN framework, and is available to the user. Since it was unfeasible to support, in the source code, all the interfaces for all the available plot types, the RAVEN Plotting system directly provide a formatted input structure for 11 different plot types (2D/3D). The user may request a plot not present among the supported ones, since the RAVEN Plotting system has the capability to construct on the fly the interface for a Plot, based on XML instructions.

### 14.2.1 Plot input structure

In order to create a plot, the user needs to add, within the `<OutputStreamManager>` block, a `<Plot>` sub-block. Similar to the `<Print>` OutputStream, the user needs to specify a `name` as an attribute of the plot. This name will then be used to request the plot in the `<Steps>` block. In addition, the Plot block accepts the following attributes:

- `dim`, *required integer attribute*, defines the dimensionality of the plot: 2 (2D) or 3 (3D).
- `interactive`, *optional bool attribute*, specifies if the Plot needs to be interactively created (real-time screen visualization).

*Default: False*

- `overwrite`, *optional bool attribute*, used when the plot needs to be dumped into picture file/s. This attribute determines whether the code needs to overwrite the image files every time a new plot (with the same name) is requested.

*Default: False*

As shown, in the XML input example below, the body of the Plot XML input contains two main sub-nodes:

- `<actions>`, where general control options for the figure layout are defined (see Section 14.2.1.1).
- `<plotSettings>`, where the actual plot options are provided.

These two main sub-block are discussed in the following paragraphs.

#### 14.2.1.1 “Actions” input block

The input in the `<actions>` sub-node is common to all the Plot types, since, in it, the user specifies all the controls that need to be applied to the figure style. This block must be unique in the definition of the `<Plot>` main block. In the following list, all the predefined “actions” are reported:

- **<how>**, comma separated list of output types:
  - screen, show the figure on the screen in interactive mode
  - pdf, save the figure as a Portable Document Format file (PDF). **Note:** The pdf format does not support multiple layers that lay on the same pixel. If the user gets an error about this, he/she should move to another format.
  - png, save the figure as a Portable Network Graphics file (PNG)
  - eps, save the figure as an Encapsulated Postscript file (EPS)
  - pgf, save the figure as a LaTeX PGF Figure file (PGF)
  - ps, save the figure as a Postscript file (PS)
  - gif, save the figure as a Graphics Interchange Format (GIF)
  - svg, save the figure as a Scalable Vector Graphics file (SVG)
  - jpeg, save the figure as a jpeg file (JPEG)
  - raw, save the figure as a Raw RGBA bitmap file (RAW)
  - bmp, save the figure as a Windows bitmap file (BMP)
  - tiff, save the figure as a Tagged Image Format file (TIFF)
  - svgz, save the figure as a Scalable Vector Graphics file (SVGZ)
- **<title>**, as the name suggests, within this block the user can specify the title of the figure. In the body of this node, a few other tags are available:

- **<text>**, *string type*, title of the figure
- **<kwargs>**, within this block the user can specify optional parameters with the following format:

```

<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>
```

The kwargs block is able to convert whatever string into a python type (for example **<param1>**{ '1stKey' : 45}**</param1>** will be converted into a dictionary, **<param2>** [56, 67]**</param2>** into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.title” method in [3].

- **<labelFormat>**, within this block the default scale formatting can be modified. In the body, a few other tags are available:
  - **<axis>**, *string*, the axis where to apply the defined format, ‘x,’ ‘y,’ or ‘both’.  
*Default:* ‘both’ **Note:** If this action will be used in a 3-D plot, the user can input ‘z’ as well and ‘both’ will apply this format to all three axis.

- **<style>**, *string*, the style of the number notation, ‘sci’ or ‘scientific’ for scientific, ‘plain’ for plain notation.  
*Default: scientific*
- **<scilimits>**, *tuple, (m, n), pair of integers*, if style is ‘sci’, scientific notation will be used for numbers outside the range  $10^m$  to  $10^n$ . Use (0,0) to include all numbers.  
**Note:** The value for this keyword, needs to be specified between brackets [for example, (5,6)].  
*Default: (0,0)*
- **<useOffset>**, *bool or double*, if True, the offset will be calculated as needed; if False, no offset will be used; if a numeric offset is specified, it will be used.  
*Default: False*
- **<figureProperties>**, within this block the user specifies how to customize the figure style/quality. Thus, through this “action” the user has got full control on the quality of the figure, its dimensions, etc. This control is performed by the following keywords:
  - **<figsize>**, *tuple (optional)*, (width, height), in inches.
  - **<dpi>**, *integer*, dots per inch.
  - **<facecolor>**, *string*, set the figure background color (please refer to “matplotlib.figure.Figure” in [3] for a list of all the colors available).
  - **<edgecolor>**, *string*, the figure edge background color (please refer to “matplotlib.figure.Figure” in [3] for a list of all the colors available).
  - **<linewidth>**, *float*, the width of lines drawn on the plot.
  - **<frameon>**, *bool*, if False, suppress drawing the figure frame.
- **<range>**, the range “action” specifies the ranges of all the axis. All the keywords in the body of this block are optional:
  - **<ymin>**, *double (optional)*, lower boundary for the y axis.
  - **<ymax>**, *double (optional)*, upper boundary for the y axis.
  - **<xmin>**, *double (optional)*, lower boundary for the x axis.
  - **<xmax>**, *double (optional)*, upper boundary for the x axis.
  - **<zmin>**, *double (optional)*, lower boundary for the z axis. **Note:** This keyword is effective in 3-D plots only.
  - **<zmax>**, *double (optional)*, upper boundary for the z axis. **Note:** This keyword is effective in 3-D plots only.
- **<camera>**, the camera item is available in 3-D plots only. Through this “action,” it is possible to orientate the plot as one wishes. The controls are:
  - **<elevation>**, *double (optional)*, stores the elevation angle in the z plane.

- **<azimuth>**, *double (optional)*, stores the azimuth angle in the x,y plane.
- **<scale>**, the scale block allows the specification of the axis scales:
  - **<xscale>**, *string (optional)*, scale of the x axis. Three options are available: “linear,” “log,” or “symlog.”  
*Default: linear*
  - **<yscale>**, *string (optional)*, scale of the y axis. Three options are available: “linear,” “log,” or “symlog.”  
*Default: linear*
  - **<zscale>**, *string (optional)*, scale of the z axis. Three options are available: “linear,” “log,” or “symlog.”  
*Default: linear* **Note:** This keyword is effective in 3-D plots only.
- **<addText>**, same as title.
- **<autoscale>**, is a convenience method for simple axis view autoscaling. It turns autoscaling on or off, and then, if autoscaling for either axis is on, it performs the autoscaling on the specified axis or axes. The following sub-nodes may be specified:
  - **<enable>**, *bool (optional)*, True turns autoscaling on, False turns it off. None leaves the autoscaling state unchanged.  
*Default: True*
  - **<axis>**, *string (optional)*, determines which axis to apply the defined format, ‘x,’ ‘y,’ or ‘both.’  
*Default: ‘both’* **Note:** If this action is used in a 3-D plot, the user can input ‘z’ as well and ‘both’ will apply this format to all three axis.
  - **<tight>**, *bool (optional)*, if True, sets the view limits to the data limits; if False, let the locator and margins expand the view limits; if None, use tight scaling if the only output is an image file, otherwise treat tight as False.
- **<horizontalLine>**, this “action” provides the ability to draw a horizontal line in the current figure. This capability might be useful, for example, if the user wants to highlight a trigger function of a variable. The following sub-nodes may be specified:
  - **<y>**, *double (optional)*, sets the y-value for the line.  
*Default: 0*
  - **<xmin>**, *double (optional)*, is the starting coordinate on the x axis.  
*Default: 0*
  - **<xmax>**, *double (optional)*, is the ending coordinate on the x axis.  
*Default: 1*
  - **<kwargs>**, within this block the user can specify optional parameters with the following format:



```
<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>
```

The kwargs block is able to convert whatever string into a python type (for example `<param1>{'1stKey':45}</param1>` will be converted into a dictionary, `<param2>[56,67]</param2>` into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.axhline” method in [3].

**Note:** This capability is not available for 3-D plots.

- `<verticalLine>`, similar to the “horizontalLine” action, this block provides the ability to draw a vertical line in the current figure. This capability might be useful, for example, if the user wants to highlight a trigger function of a variable. The following sub-nodes may be specified:

- `<x>`, *double (optional)*, sets the x coordinate of the line.  
*Default: 0*
- `<ymin>`, *double (optional)*, starting coordinate on the y axis.  
*Default: 0*
- `<ymax>`, *double (optional)*, ending coordinate on the y axis.  
*Default: 1*
- `<kwargs>`, within this block the user can specify optional parameters with the following format:

```
<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>
```

The kwargs block is able to convert whatever string into a python type (for example `<param1>{'1stKey':45}</param1>` will be converted into a dictionary, `<param2>[56,67]</param2>` into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.axvline” method in [3].

**Note:** This capability is not available for 3-D plots.

- `<horizontalRectangle>`, this “action” provides the ability to draw, in the current figure, a horizontally orientated rectangle. This capability might be useful, for example, if the user wants to highlight a zone in the plot. The following sub-nodes may be specified:

- `<ymin>`, *double (required)*, starting coordinate on the y axis.
- `<ymax>`, *double (required)*, ending coordinate on the y axis.



- **<xmin>**, *double (optional)*, starting coordinate on the x axis.  
*Default: 0*
- **<xmax>**, *double (optional)*, ending coordinate on the x axis. *Default = 1*
- **<kwargs>**, within this block the user can specify optional parameters with the following format:

```
<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>
```

The kwargs block is able to convert whatever string into a python type (for example **<param1>**{ '1stKey' : 45}**</param1>** will be converted into a dictionary, **<param2>** [56, 67]**</param2>** into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.axhspan” method in [3].

**Note:** This capability is not available for 3D plots.

- **<verticalRectangle>**, this “action” provides the possibility to draw, in the current figure, a vertically orientated rectangle. This capability might be useful, for example, if the user wants to highlight a zone in the plot. The following sub-nodes may be specified:
  - **<xmin>**, *double (required)*, starting coordinate on the x axis.
  - **<xmax>**, *double (required)*, ending coordinate on the x axis.
  - **<ymin>**, *double (optional)*, starting coordinate on the y axis.  
*Default: 0*
  - **<ymax>**, *double (optional)*, ending coordinate on the y axis.  
*Default: 1*
  - **<kwargs>**, within this block the user can specify optional parameters with the following format:

```
<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>
```

The kwargs block is able to convert whatever string into a python type (for example **<param1>**{ '1stKey' : 45}**</param1>** will be converted into a dictionary, **<param2>** [56, 67]**</param2>** into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.axvspan” method in [3].

**Note:** This capability is not available for 3D plots.

- **<axesBox>**, this keyword controls the axes’ box. Its value can be ‘on’ or ‘off’.

- **<axisProperties>**, this block is used to set axis properties. There are no fixed keywords. If only a single property needs to be set, it can be specified as the body of this block, otherwise a dictionary-like string needs to be provided. For reference regarding the available keys, refer to “matplotlib.pyplot.axis” method in [3].
- **<grid>**, this block is used to define a grid that needs to be added in the plot. The following keywords can be inputted:
  - **<b>**, *boolean (required)*, toggles the grid lines on or off.
  - **<which>**, *double (required)*, ending coordinate on the x axis.
  - **<axis>**, *double (optional)*, starting coordinate on the y axis.  
*Default: 0*
  - **<kwargs>**, within this block the user can specify optional parameters with the following format:

```

<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>

```

The kwargs block is able to convert whatever string into a python type (for example **<param1>**{`1stKey` : 45}**</param1>** will be converted into a dictionary, **<param2>** [56, 67]**</param2>** into a list, etc.).

#### 14.2.1.2 “plotSettings” input block

The sub-block identified by the keyword **<plotSettings>** is used to define the plot characteristics. Within this sub-section at least a **<plot>** block must be present. the **<plot>** sub-section may not be unique within the **<plotSettings>** definition; the number of **<plot>** sub-blocks is equal to the number of plots that need to be placed in the same figure.

If sub-plots are to be defined then **<gridSpace>** needs to be present. **<gridSpace>** specifies the geometry of the grid that a subplot will be placed. The number of rows and number of columns of the grid need to be set.

For example, in the following XML cut, a “line” and a “scatter” type are combined in the same figure.

```

<OutputStreamManager>
  <Plot name='example2PlotsCombined' dim='2'>
    <actions>
      <!-- Actions -->
    </actions>
    <plotSettings>

```

```

<gridSpace>2 2</gridSpace>
<plot>
  <type>line</type>
  <x>d-type|Output|x1</x>
  <y>d-type|Output|y1</y>
  <xlabel>label X</xlabel>
  <ylabel>label Y</ylabel>
  <gridLocation>
    <x>0 2</x>
    <y>0</y>
  </gridLocation>
</plot>
<plot>
  <type>scatter</type>
  <x>d-type|Output|x2</x>
  <y>d-type|Output|y2</y>
  <xlabel>label X</xlabel>
  <ylabel>label Y</ylabel>
  <gridLocation>
    <x>0 2</x>
    <y>1</y>
  </gridLocation>
</plot>
</plot>
</plotSettings>
</Plot>
</OutputStreamManager>

```

The axis labels are conditionally optional nodes that can be defined under the `<plotSetting>`. If the plot does not contain any sub-plots, i.e. `<gridSpace>` is not defined then the axis labels are global parameters for the figure which are defined under `<plotSettings>`, otherwise the axis labels can be defined under `<plot>` for each sub-plot separately.

- `<xlabel>`, *string, optional parameter*, the x axis label.
- `<ylabel>`, *string, optional parameter*, the y axis label.
- `<zlabel>`, *string, optional parameter (3D plots only)*, the z axis label.

As already mentioned, within the `<plotSettings>` block, at least a `<plot>` sub-block needs to be specified. Independent of the plot type, some keywords are mandatory:

- **<type>**, *string, required parameter*, the plot type (for example, line, scatter, wireframe, etc.).
- **<x>**, *string, required parameter*, this parameter needs to be considered as the x coordinate.
- **<y>**, *string, required parameter*, this parameter needs to be considered as the y coordinate.
- **<z>**, *string required parameter (3D plots only)*, this parameter needs to be considered as the z coordinate.

In addition, other plot-dependent keywords, reported in Section 14.2.1.3, can be provided.

Under the **<plot>** sub-block other optional keywords can be specified, such as:

- **<xlabel>**, *string, optional parameter*, the x axis label.
- **<ylabel>**, *string, optional parameter*, the y axis label.
- **<zlabel>**, *string, optional parameter (3D plots only)*, the z axis label.
- **<gridLocation>**, *xmlNode, optional xmlNode (depending on the grid geometry)*
  - **<x>**, *integer, required parameter*, the position of the subPlot in the grid Space. if this node has a single value then the subplot occupies a single node at the specified location, otherwise the second integer represents the number of nodes that this subplot occupies, i.e. in the example above the first subplot occupies 2 nodes starting from the zero node in x direction.
  - **<y>**, *integer, required parameter*, the position of the subPlot in the grid Space. if this node has a single value then the subplot occupies a single node at the specified location, otherwise the second integer represents the number of nodes that this subplot occupies, i.e. in the example above the first subplot occupies a single node at the zero node in y direction.
- **<colorMap>**, *string, optional parameter*, this parameter will be used to define a color map.

As already mentioned, the Plot system accepts as input for the visual parameters (i.e., x, y, z, colorMap), data only from a **DataObjects** object. Considering the structure of “DataObjects,” the parameters are inputted as follows: DataObjectName | Input (or) Output | variableName.

If the “variableName” contains the symbol |, it must be surrounded by brackets: [DataObjectName | Input

### 14.2.1.3 Predefined Plotting System: 2D/3D

As already mentioned above, the Plotting system provides a specialized input structure for several different kind of plots specified in the `<type>` node:

- *2 Dimensional plots:*

- `scatter` creates a scatter plot of  $x$  vs  $y$ , where  $x$  and  $y$  are sequences of numbers of the same length.
- `line` creates a line plot of  $x$  vs  $y$ , where  $x$  and  $y$  are sequences of numbers of the same length.
- `histogram` computes and draws the histogram of  $x$ . **Note:** This plot accepts only the XML node `<x>` even if it is considered as a 2D plot type.
- `stem` plots vertical lines at each  $x$  location from the baseline to  $y$ , and places a marker there.
- `step` creates a 2 dimensional step plot.
- `pseudocolor` creates a pseudocolor plot of a two dimensional array. The two dimensional array is built creating a mesh from `<x>` and `<y>` data, in conjunction with the data specified in the `<colorMap>` node.
- `contour` builds a contour plot creating a plot from `<x>` and `<y>` data, in conjunction with the data specified in the `<colorMap>` node.
- `filledContour` creates a filled contour plot from `<x>` and `<y>` data, in conjunction with the data specified in the `<colorMap>` node.

- *3 Dimensional plots:*

- `scatter` creates a scatter plot of  $(x,y)$  vs  $z$ , where  $x$ ,  $y$ ,  $z$  are sequences of numbers of the same length.
- `line` creates a line plot of  $(x,y)$  vs  $z$ , where  $x$ ,  $y$ ,  $z$  are sequences of numbers of the same length.
- `stem` creates a 3 Dimensional stem plot of  $(x,y)$  vs  $z$ .
- `surface` creates a surface plot of  $(x,y)$  vs  $z$ . By default it will be colored in shades of a solid color, but it also supports color mapping.
- `wireframe` creates a 3D wire-frame plot. No color mapping is supported.
- `tri-surface` creates a 3D tri-surface plot. It is a surface plot with automatic triangulation.
- `contour3D` builds a 3D contour plot creating the plot from `<x>`, `<y>` and `<z>` data, in conjunction with the data specified in `<colorMap>`.

- `filledContour3D` builds a filled 3D contour plot creating the plot from `<x>`, `<y>` and `<z>` data, in conjunction with the data specified in `<colorMap>`.
- `histogram` computes and draws the histogram of x and y. **Note:** This plot accepts only the XML nodes `<x>` and `<y>` even if it is considered as 3D plot type since the frequency is mapped to the third dimension.

As already mentioned, the settings for each plot type are specified within the XML block called `<plot>`. The sub-nodes that are available depends on the plot type as each plot type has its own set of parameters that can be specified.

In the following sub-sections all the options for the plot types listed above are reported.

### 14.2.2 2D & 3D Scatter plot

In order to create a “scatter” plot, either 2D or 3D, the user needs to write in the `<type>` body the keyword “scatter.” In order to customize the plot, the user can define the following XML sub nodes:

- `<s>`, *integer, optional field*, represents the size in points<sup>2</sup>. The “points” have the same meaning of the font size (e.g. Times New Roman, pts 10). In here the user specifies the area of the marker size.  
*Default: 20*
- `<c>`, *string, optional field*, specifies the color or sequence of color to use. `<c>` can be a single color format string, a sequence of color specifications of length N, or a sequence of N numbers to be mapped to colors using the `cmap` and `norm` specified via `<kwargs>`.  
**Note:** `<c>` should not be a single numeric RGB or RGBA sequence because that is indistinguishable from an array of values to be colormapped. `<c>` can be a 2D array in which the rows are RGB or RGBA. **Note:** `<colorMap>` will overwrite `<c>`. If `<colorMap>` is defined then the color set used can be defined by `<cmap>`. If no `<cmap>` is given then the default color set of “`matplotlib.pyplot.scatter`” method in [3] is used. If `<colorMap>` is not defined then the plot is in solid color (default *blue*) as defined with `<color>` in `<kwargs>`.
- `<marker>`, *string, optional field*, specifies the type of marker to use.  
*Default: o*
- `<alpha>`, *string, optional field*, sets the alpha blending value, between 0 (transparent) and 1 (opaque).  
*Default: None*
- `<linewidths>`, *string, optional field*, widths of lines used in the plot. Note that this is a tuple, and if you set the `linewidths` argument you must set it as a sequence of floats.  
*Default: None;*

- **<kwargs>**, within this block the user can specify optional parameters with the following format:

```

<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>
```

The kwargs block is able to convert whatever string into a python type (for example **<param1>** { '1stKey' : 45 } **</param1>** will be converted into a dictionary, **<param2>** [ 56, 67 ] **</param2>** into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.scatter” method in [3].

### 14.2.3 2D & 3D Line plot

In order to create a “line” plot, either 2D or 3D, the user needs to write in the **<type>** body the keyword “line.” In order to customize the plot, the user can define the following XML sub nodes:

- **<interpolationType>**, *string, optional field*, is the type of interpolation algorithm to use for the data. Available options are “nearest,” “linear,” “cubic,” “multiquadric,” “inverse,” “gaussian,” “Rbflinear,” “Rbfcubic,” “quintic,” and “thin\_plate.”  
*Default: linear*
- **<interpPointsX>**, *integer, optional field*, sets the number of points need to be used for interpolation of the x axis.
- **<interpPointsY>**, *integer, optional field*, sets the number of points need to be used for interpolation of the y axis. (only 3D line plot). **Note:** If **<colorMap>** is used then a **scatter plot** will be plotted.

### 14.2.4 2D & 3D Histogram plot

In order to create a “histogram” plot, either 2D or 3D, the user needs to write in the **<type>** body the keyword “histogram.” In order to customize the plot, the user can define the following XML sub nodes:

- **<bins>**, *integer or array\_like, optional field*, sets the number of bins if an integer is used or a sequence of edges if a python list is used.  
*Default: 10*
- **<normed>**, *boolean, optional field*, if True then the the histogram will be normalized to 1.  
*Default: False*

- **<weights>**, *sequence, optional field*, represents an array of weights, of the same shape as *x*. Each value in *x* only contributes its associated weight towards the bin count (instead of 1). If *normed* is True, the weights are normalized, so that the integral of the density over the range remains 1.

*Default: None*

- **<cumulative>**, *boolean, optional field*, if True, then a histogram is computed where each bin gives the counts in that bin plus all bins for smaller values. The last bin gives the total number of data points. If *normed* is also True then the histogram is normalized such that the last bin equals 1. If *cumulative* evaluates to less than 0 (e.g., -1), the direction of accumulation is reversed. In this case, if *normed* is also True, then the histogram is normalized such that the first bin equals 1.

*Default: False*

- **<histtype>**, *string, optional field*, The type of histogram to draw:

- **bar** is a traditional bar-type histogram. If multiple data sets are given the bars are arranged side by side.
- **barstacked** is a bar-type histogram where multiple data sets are stacked on top of each other.
- **step** generates a line plot that is by default unfilled.
- **stepfilled** generates a line plot that is by default filled.

*Default: bar*

- **<align>**, *string, optional field*, controls how the histogram is plotted.

- **left** bars are centered on the left bin edge.
- **mid** bars are centered between the bin edges.
- **right** bars are centered on the right bin edges.

*Default: mid*

- **<orientation>**, *string, optional field*, specifies the orientation of the histogram:

- **horizontal**
- **vertical**

*Default: vertical*



- **<rwidth>**, *float, optional field*, sets the relative width of the bars as a fraction of the bin width.  
*Default: None*
- **<log>**, *boolean, optional field*, sets a log scale.  
*Default: False*
- **<color>**, *string, optional field*, specifies the color of the histogram.  
*Default: blue;*
- **<stacked>**, *boolean, optional field*, if True, multiple data elements are stacked on top of each other. If False, multiple data sets are arranged side by side if histtype is 'bar' or on top of each other if histtype is 'step.'  
*Default: False*
- **<kwargs>**, within this block the user can specify optional parameters with the following format:

```

<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>

```

The kwargs block is able to convert whatever string into a python type (for example **<param1>** { '1stKey' : 45 } **</param1>** will be converted into a dictionary, **<param2>** [56, 67] **</param2>** into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.hist” method in [3].

### 14.2.5 2D & 3D Stem plot

In order to create a “stem” plot, either 2D or 3D, the user needs to write in the **<type>** body the keyword “stem.” In order to customize the plot, the user can define the following XML sub nodes:

- **<linefmt>**, *string, optional field*, sets the line style used in the plot.  
*Default: b-*
- **<markerfmt>**, *string, optional field*, sets the type of marker format to use in the plot.  
*Default: bo*
- **<basefmt>**, *string, optional field*, sets the base format.  
*Default: r-*
- **<kwargs>**, within this block the user can specify optional parameters with the following format:

```
<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>
```

The kwargs block is able to convert whatever string into a python type (for example `<param1> { '1stKey' : 45 }</param1>` will be converted into a dictionary, `<param2> [56, 67]</param2>` into a list, etc.).

For reference regarding the available kwargs, see “matplotlib.pyplot.stem” method in [3].

### 14.2.6 2D Step plot

In order to create a 2D “step” plot, the user needs to write in the `<type>` body the keyword “step.” In order to customize the plot, the user can define the following XML sub nodes:

- `<where>`, *string, optional field*, specifies the positioning:
  - **pre**, the interval from  $x[i]$  to  $x[i+1]$  has level  $y[i+1]$
  - **post**, that interval has level  $y[i]$
  - **mid**, the jumps in  $y$  occur half-way between the  $x$ -values

*Default: mid*

- `<kwargs>`, within this block the user can specify optional parameters with the following format:

```
<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>
```

The kwargs block is able to convert whatever string into a python type (for example `<param1> { '1stKey' : 45 }</param1>` will be converted into a dictionary, `<param2> [56, 67]</param2>` into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.step” method in [3].

### 14.2.7 2D Pseudocolor plot

In order to create a 2D “pseudocolor” plot, the user needs to write in the `<type>` body the keyword “pseudocolor.” In order to customize the plot, the user can define the following XML sub nodes:

- **<interpolationType>**, *string, optional field*, is the type of interpolation algorithm to use for the data. Available options are “nearest,” “linear,” “cubic,” “multiquadric,” “inverse,” “gaussian,” “Rbflinear,” “Rbfcubic,” “quintic,” and “thin\_plate.”  
*Default:* [ linear]
- **<interpPointsX>**, *integer, optional field*, sets the number of points need to be used for interpolation of the x axis.
- **<kwargs>**, within this block the user can specify optional parameters with the following format:

```

<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>

```

The kwargs block is able to convert whatever string into a python type (for example **<param1>** { '1stKey' : 45} **</param1>** will be converted into a dictionary, **<param2>** [56, 67] **</param2>** into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.pcolor” method in [3].

#### 14.2.8 2D Contour or filledContour plots

In order to create a 2D “contour” or “filledContour” plot, the user needs to write in the **<type>** body the keyword “contour” or “filledContour,” respectively. In order to customize the plot, the user can define the following XML sub-nodes:

- **<numberBins>**, *integer, optional field*, sets the number of bins.  
*Default:* 5
- **<interpolationType>**, *string, optional field*, is the type of interpolation algorithm to use for the data. Available options are “nearest,” “linear,” “cubic,” “multiquadric,” “inverse,” “gaussian,” “Rbflinear,” “Rbfcubic,” “quintic,” and “thin\_plate.”  
*Default:* linear
- **<interpPointsX>**, *integer, optional field*, sets the number of points need to be used for interpolation of the x axis.
- **<colorMap>** vector is the array to visualize. If **<colorMap>** is defined then the color set used can be defined by **<cmap>**. If no **<cmap>** is given then the plot is in solid color (default *blue*) as defined with **<color>** in **<kwargs>**.
- **<cmap>**, *string, optional field*, defines the color map to use for this plot.  
*Default:* None

- **<kwargs>**, within this block the user can specify optional parameters with the following format:

```

<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>

```

The kwargs block is able to convert whatever string into a python type (for example **<param1>** { '1stKey' : 45 } **</param1>** will be converted into a dictionary, **<param2>** [ 56, 67 ] **</param2>** into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.contour” method in [3].

### 14.2.9 3D Surface Plot

In order to create a 3D “surface” plot, the user needs to write in the **<type>** body the keyword “surface.” In order to customize the plot, the user can define the following XML sub nodes:

- **<rstride>**, *integer, optional field*, specifies the array row stride (step size).  
*Default: 1*
- **<cstride>**, *integer, optional field*, specifies the array column stride (step size).  
*Default: 1*
- **<cmap>**, *string, optional field*, defines the color map to use for this plot.  
*Default: None* **Note:** If **<colorMap>** is defined then the plot will always use a color set even if no **<cmap>** is given. In such a case, if no **<cmap>** is given, then the default color set of “matplotlib.pyplot.surface” method in [3] is used. If **<colorMap>** and **<cmap>** are both not defined then the plot is in solid color (*default blue*) as defined with **<color>** in **<kwargs>**.
- **<antialiased>**, *boolean, optional field*, determines whether or not the rendering should be antialiased.  
*Default: False*
- **<linewidth>**, *integer, optional field*, defines the widths of lines rendered on the plot.  
*Default: 0*
- **<interpolationType>**, *string, optional field*, is the type of interpolation algorithm to use for the data. Available options are “nearest,” “linear,” “cubic,” “multiquadric,” “inverse,” “gaussian,” “Rbflinear,” “Rbfcubic,” “quintic,” and “thin\_plate.”  
*Default: linear*

- **<interpPointsX>**, *integer, optional field*, sets the number of points need to be used for interpolation of the x axis.
- **<interpPointsY>**, *integer, optional field*, sets the number of points need to be used for interpolation of the y axis.
- **<kwargs>**, within this block the user can specify optional parameters with the following format:

```

<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>

```

The kwargs block is able to convert whatever string into a python type (for example **<param1>** { '1stKey' : 45 } **</param1>** will be converted into a dictionary, **<param2>** [ 56, 67 ] **</param2>** into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.surface” method in [3].

### 14.2.10 3D Wireframe Plot

In order to create a 3D “wireframe” plot, the user needs to write in the **<type>** body the keyword “wireframe.” In order to customize the plot, the user can define the following XML sub nodes:

- **<rstride>**, *integer, optional field*, sets the array row stride (step size).  
*Default: 1*
- **<cstride>**, *integer, optional field*, sets the array column stride (step size).  
*Default: 1*
- **<cmap>**, *string, optional field*, defines the color map to use for this plot.  
*Default: None* **Note:** **<cmap>** is not applicable in the current version of Matplotlib for wireframe plots. However, if the colorMap option is set then a surface plot is plotted with a transparency of 0.4 on top of wireframe to give a visual colormap. **Note:** If **<colorMap>** is defined then the plot will always use a color set even if no **<cmap>** is given. In such a case, if no **<cmap>** is given, then the default color set of “matplotlib.pyplot.surface” method in [3] is used. If **<colorMap>** and **<cmap>** are both not defined then the plot is in solid color (*default blue*) as defined with **<color>** in **<kwargs>**.
- **<interpolationType>**, *string, optional field*, is the type of interpolation algorithm to use for the data. Available options are “nearest,” “linear,” “cubic,” “multiquadric,” “inverse,” “gaussian,” “Rbflinear,” “Rbfcubic,” “quintic,” and “thin\_plate.”  
*Default: linear*

- **<interpPointsX>**, *integer, optional field*, sets the number of points need to be used for interpolation of the x axis.
- **<interpPointsY>**, *integer, optional field*, sets the number of points need to be used for interpolation of the y axis.
- **<kwargs>**, within this block the user can specify optional parameters with the following format:

```

<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>

```

The kwargs block is able to convert whatever string into a python type (for example **<param1>** { '1stKey' : 45 } **</param1>** will be converted into a dictionary, **<param2>** [56, 67] **</param2>** into a list, etc.). For reference regarding the available kwargs, see “matplotlib.pyplot.wireframe” method in [3].

### 14.2.11 3D Tri-surface Plot

In order to create a 3D “tri-surface” plot, the user needs to write in the **<type>** body the keyword “tri-surface.” In order to customize the plot, the user can define the following XML sub nodes:

- **<color>**, *string, optional field*, sets the color of the surface patches.  
*Default: b*
- **<shade>**, *boolean, optional field*, determines whether to apply shading or not.  
*Default: False*
- **<cmap>**, *string, optional field*, defines the color map to use for this plot.  
*Default: None* **Note:** If **<colorMap>** is defined then the plot will always use a color set even if no **<cmap>** is given. In such a case, if no **<cmap>** is given, then the default color set of “matplotlib.pyplot.trisurface” method in [3] is used. If **<colorMap>** and **<cmap>** are both not defined then the plot is in solid color (*default blue*) as defined with **<color>** in **<kwargs>**.
- **<kwargs>**, within this block the user can specify optional parameters with the following format:

```

<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>

```

The `kwargs` block is able to convert whatever string into a python type (for example `<param1> { '1stKey' : 45 } </param1>` will be converted into a dictionary, `<param2> [56, 67] </param2>` into a list, etc.). For reference regarding the available `kwargs`, see “`matplotlib.pyplot.trisurface`” method in [3].

### 14.2.12 3D Contour or filledContour plots

In order to create a 3D “Contour” or “filledContour” plot, the user needs to write in the `<type>` body the keyword “`contour3D`” or “`filledContour3D`,” respectively. In order to customize these plots, the user can define the following XML sub nodes:

- `<numberBins>`, *integer, optional field*, sets the number of bins to use.  
*Default: 5*
- `<interpolationType>`, *string, optional field*, is the type of interpolation algorithm to use for the data. Available options are “nearest,” “linear,” “cubic,” “multiquadric,” “inverse,” “gaussian,” “Rbflinear,” “Rbfcubic,” “quintic,” and “thin\_plate.”  
*Default: linear*
- `<interpPointsX>`, *integer, optional field*, sets the number of points need to be used for interpolation of the x axis.
- `<interpPointsY>`, *integer, optional field*, sets the number of points need to be used for interpolation of the y axis.
- `<colorMap>` vector is the array to visualize. If `<colorMap>` is defined then the color set used can be defined by `<cmap>`. If no `<cmap>` is given then the plot is in solid color (default *blue*) as defined with `<color>` in `<kwargs>`.
- `<kwargs>`, within this block the user can specify optional parameters with the following format:

```
<kwargs>
  <param1>value1</param1>
  <param2>value2</param2>
</kwargs>
```

The `kwargs` block is able to convert whatever string into a python type (for example `<param1> { '1stKey' : 45 } </param1>` will be converted into a dictionary, `<param2> [56, 67] </param2>` into a list, etc.). For reference regarding the available `kwargs`, see “`matplotlib.pyplot.contour3d`” method in [3].

### 14.2.13 DataMining plots

In order to create a “DataMining” plot, the user needs to write in the `<type>` body the keyword “dataMining”. “DataMining” plots are based on 2D or 3D Scattering plots, depending on the method/algorithm used in the “DataMining” postprocessor [see 15.5.8]. These plots are created to ease the color labeling the clusters, etc parameters in the data. The following are the optional or required input parameters that can be used in these plots additional to the coordinate inputs `<x>`, `<y>`, or `<z>` depending on the dimension:

- `<type>`, *string, required field*, this block should read “dataMining” in order to create a data mining plot.
- `<SKLtype>`, *string, required field*, name of the algorithm used in the “dataMining” postprocessor. It is one of:
  - cluster: for clustering algorithms, such as KMeans clustering.
  - bicluster ( **Note:** not implemented yet!)
  - mixture: for Gaussian mixture algorithms, such as GMM classifier
  - manifold: for Manifold Learning algorithms, such as Spectral Embedding
  - decomposition: for decomposing signals in components algorithms, such as Principal Component Analysis (PCA)
- `<clusterLabels>`, *string, optional field*, defines the place where the labels of the clusters are located. As in the visual parameters (i.e., x,y,z and colorMap) this is also from a **DataObjects** object. Considering the structure of “DataObjects”, the labels inputted as follows: `DataObjectName|Output|DataMiningPPNameLabels`.  
*Default: None*
- `<noClusters>`, *integer, optional field*, defines the number of clusters used in the “dataMining” postprocessor  
*Default: 1*
- `<kwargs>`, within this block the user can specify optional parameters with the following format:

```
<kwargs>  
  <param1>value1</param1>  
  <param2>value2</param2>  
</kwargs>
```

The kwargs block is able to convert whatever string into a python type (for example `<param1> { '1stKey' : 45 } </param1>` will be converted into a dictionary, `<param2> [56, 67] </param2>` into a list, etc.).



For reference regarding the other available kwargs, see “matplotlib.pyplot.scatter” method in [3].

#### 14.2.14 Example XML input

```
<OutputStreamManager>
  <Plot name='2DHistoryPlot' dim='2' interactive='False'
    overwrite='False'>
    <actions>
      <how>pdf,png,eps</how>
      <title>
        <text>***</text>
      </title>
    </actions>
    <plot_settings>
      <plot>
        <type>line</type>
        <x>stories|Output|time</x>
        <y>stories|Output|pipe1_Hw</y>
        <kwargs>
          <color>green</color>
          <label>pipe1-Hw</label>
        </kwargs>
      </plot>
      <plot>
        <type>line</type>
        <x>stories|Output|time</x>
        <y>stories|Output|pipe1_aw</y>
        <kwargs>
          <color>blue</color>
          <label>pipe1-aw</label>
        </kwargs>
      </plot>
      <xlabel>time [s]</xlabel>
      <ylabel>evolution</ylabel>
    </plot_settings>
  </Plot>
</OutputStreamManager>
```

## 15 Models

In RAVEN, **Models** are important entities. A model is an object that employs a mathematical representation of a phenomenon, either of a physical or other nature (e.g. statistical operators, etc.). From a practical point of view, it can be seen, as a “black box” that, given an input, returns an output.

RAVEN has a strict classification of the different types of models. Each “class” of models is represented by the definition reported above, but it can be further classified based on its particular functionalities:

- **<Code>** represents an external system code that employs a high fidelity physical model.
- **<Dummy>** acts as “transfer” tool. The only action it performs is transferring the the information in the input space (inputs) into the output space (outputs). For example, it can be used to check the effect of a sampling strategy, since its outputs are the sampled parameters’ values (input space) and a counter that keeps track of the number of times an evaluation has been requested.
- **<ROM>**, or reduced order model, is a mathematical model trained to predict a response of interest of a physical system. Typically, ROMs trade speed for accuracy representing a faster, rough estimate of the underlying phenomenon. The “training” process is performed by sampling the response of a physical model with respect to variation of its parameters subject to probabilistic behavior. The results (outcomes of the physical model) of the sampling are fed into the algorithm representing the ROM that tunes itself to replicate those results.
- **<ExternalModel>**, as its name suggests, is an entity existing outside the RAVEN framework that is embedded in the RAVEN code at run time. This object allows the user to create a Python module that will be treated as a predefined internal model object.
- **<PostProcessor>** is a container of all the actions that can manipulate and process a data object in order to extract key information, such as statistical quantities, clustering, etc.

Before analyzing each model in detail, it is important to mention that each type needs to be contained in the main XML node **<Models>**, as reported below:

### Example:

```
<Simulation>
...
<Models>
...
<WhateverModel name='whatever'>
```

```
    ...
  </WhateverModel>
  ...
</Models>
...
</Simulation>
```

In the following sub-sections each **Model** type is fully analyzed and described.

## 15.1 Code

As already mentioned, the **Code** model represents an external system software employing a high fidelity physical model. The link between RAVEN and the driven code is performed at run-time, through coded interfaces that are the responsible for transferring information from the code to RAVEN and vice versa. In Section 18, all of the available interfaces are reported and, for advanced users, Section 19 explains how to couple a new code.

The specifications of this model must be defined within a **<Code>** XML block. This XML node accepts the following attributes:

- **name**, *required string attribute*, user-defined identifier of this model. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.
- **subType**, *required string attribute*, specifies the code that needs to be associated to this Model. **Note:** See Section 18 for a list of currently supported codes.

This model can be initialized with the following children:

- **<executable>** *string, required field* specifies the path of the executable to be used. **Note:** Either an absolute or relative path can be used.
- **<alias>** *string, optional field* specifies aliases for any variables of interest coming from the code to which this model refers. These aliases can be used anywhere in the RAVEN input to refer to the code variables. In the body of this node the user specifies the name of the variable that RAVEN will look for in the output files of the code. The actual alias, usable throughout the input, is instead defined in the **variable** attribute of this tag. **Note:** The user can specify as many aliases as needed.

*Default: None*

- **<clargs>** *string, optional field* allows addition of command-line arguments to the execution command. If the code interface specified in **<Code>** **subType** does not specify how to determine the input file(s), this node must be used to specify them. There are several types of **<clargs>**, based on the **type**:

- **type string, required field** specifies the type of command-line argument to add. Options include 'input', 'output', 'prepend', 'postpend', and 'text'.
- **arg string, optional field** specifies the flag to be used before the entry. For example, **arg='-i'** would place a -i before the entry in the execution command. Required for the 'output' type.
- **extension string, optional field** specifies the type of file extension to use (for example, -i or -o). This links the <Input> file in the <Step> to this location in the execution command. Required for 'input' type.

The execution command is combined in the order 'prepend', <executable>, 'input', 'output', 'text', 'postpend'.

- **<fileargs> string, optional field** like <clargs>, but allows editing of input files to specify the output filename and/or auxiliary file names. The location in the input files to edit using these arguments are identified in the input file using the prefix-postfix notation, which defaults to \$RAVEN-var\$ for variable keyword *var*. The variable keyword is then listed in the <fileargs> node in the attribute **arg** to couple it in Raven. If the code interface specified in <Code> **subType** does not specify how to name the output file, that must be specified either through <clargs> or <filargs>, with **type 'output'**. The attributes required for <fileargs> are as follows:
  - **type string, required field** specifies the type of entry to replace in the file. Possible values for <fileargs> **type** are 'input' and 'output'.
  - **arg string, required field** specifies the Raven variable with which to replace the file of interest. This should match the entry in the template input file; that is, if \$RAVEN-auxinp\$ is in the input file, the arg for the corresponding input file should be 'auxinp'.
  - **extension string, optional field** specifies the extension of the input file that should replace the Raven variable in the input file. This attribute is required for the 'input' type and ignored for the 'output' type. **Note:** Currently, there can only be a one-to-one pairing between input files and extensions; that is, multiple Raven-editable input files cannot have the same extension.

### Example:

```
<Simulation>
...
<Models>
...
<Code name='aUserDefinedName' subType='RAVEN_Driven_code'>
  <executable>path_to_executable</executable>
  <alias variable='internal_variable_name1'>
    External_Code_Variable_Name_1
```

```

    </alias>
    <alias variable='internal_variable_name2'>
      External_Code_Variable_Name_2
    </alias>
    <clargs type='prepend' arg='python' />
    <clargs type='input' arg='-i' extension='.i' />
    <fileargs type='input' arg='aux' extension='.two' />
    <fileargs type='output' arg='out' />
  </Code>
  ...
</Models>
...
</Simulation>

```

## 15.2 Dummy

The **Dummy** model is an object that acts as a pass-through tool. The only action it performs is transferring the information in the input space (inputs) to the output space (outputs). For example, it can be used to check the effect of a particular sampling strategy, since its outputs are the sampled parameters' values (input space) and a counter that keeps track of the number of times an evaluation has been requested.

The specifications of this model must be defined within a **<Dummy>** XML block. . This XML node accepts the following attributes:

- **name**, *required string attribute*, user-defined identifier of this model. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.
- **subType**, *required string attribute*, this attribute must be kept empty.

Given a particular *Step* using this model, if this model is linked to a *Data* with the role of **Output**, it expects one of the output parameters will be identified by the keyword “OutputPlaceholder” (see Section 17).

### Example:

```

<Simulation>
  ...
  <Models>
    ...
    <Dummy name='aUserDefinedName' subType=' ' />
    ...
  </Models>

```

```
...  
</Simulation>
```

## 15.3 ROM

A Reduced Order Model (ROM) is a mathematical model consisting of a fast solution trained to predict a response of interest of a physical system. The “training” process is performed by sampling the response of a physical model with respect to variations of its parameters subject, for example, to probabilistic behavior. The results (outcomes of the physical model) of the sampling are fed into the algorithm representing the ROM that tunes itself to replicate those results. RAVEN supports several different types of ROMs, both internally developed and imported through an external library called “scikit-learn” [4].

Currently in RAVEN, the ROMs are classified into several sub-types that, once chosen, provide access to several different algorithms. These sub-types are specified in the `subType` attribute and should be one of the following:

- `'NDspline'`
- `'GaussPolynomialRom'`
- `'HDMRRom'`
- `'NDinvDistWeight'`
- `'SciKitLearn'`
- `'MSR'`

The specifications of this model must be defined within a `<ROM>` XML block. This XML node accepts the following attributes:

- `name`, *required string attribute*, user-defined identifier of this model. **Note:** As with other objects, this identifier can be used to reference this specific entity from other input blocks in the XML.
- `subType`, *required string attribute*, defines which of the sub-types should be used, choosing among the previously reported types. This choice conditions the subsequent the required and/or optional `<ROM>` sub-nodes.

In the `<ROM>` input block, the following XML sub-nodes are required, independent of the `subType` specified:

- **<Features>**, *comma separated string, required field*, specifies the names of the features of this ROM. **Note:** These parameters are going to be requested for the training of this object (see Section 17.4);
- **<Target>**, *comma separated string, required field*, contains a comma separated list of the targets of this ROM. These parameters are the Figures of Merit (FOMs) this ROM is supposed to predict. **Note:** These parameters are going to be requested for the training of this object (see Section 17.4).

The types and meaning of the remaining sub-nodes depend on the sub-type specified in the attribute **subType**. Note that if an HistorySet is provided in the training step then a temporal ROM is created, i.e. a ROM that generates not a single value prediction of each element indicated in the **<Target>** block but its full temporal profile.

In the sections that follow, the specifications of each type are reported.

### 15.3.1 NDspline

The NDspline sub-type contains a single ROM type, based on an  $N$ -dimensional spline interpolation/extrapolation scheme. In spline interpolation, the interpolant is a special type of piecewise polynomial called a spline. The interpolation error can be made small even when using low degree polynomials for the spline. Spline interpolation avoids the problem of Runge’s phenomenon, in which oscillation can occur between points when interpolating using higher degree polynomials.

In order to use this ROM, the **<ROM>** attribute **subType** needs to be **'NDspline'** (see the example below). No further XML sub-nodes are required. **Note:** This ROM type must be trained from a regular Cartesian grid. Thus, it can only be trained from the outcomes of a grid sampling strategy.

#### Example:

```

<Simulation>
...
  <Models>
    ...
    <ROM name='aUserDefinedName' subType='NDspline'>
      <Features>var1,var2,var3</Features>
      <Target>result1,result2</Target>
    </ROM>
    ...
  </Models>
  ...
</Simulation>

```

### 15.3.2 GaussPolynomialRom

The GaussPolynomialRom sub-type contains a single ROM type, based on a characteristic Gaussian polynomial fitting scheme: generalized polynomial chaos expansion (gPC). In gPC, sets of polynomials orthogonal with respect to the distribution of uncertainty are used to represent the original model. The method converges moments of the original model faster than Monte Carlo for small-dimension uncertainty spaces ( $N < 15$ ). In order to use this ROM, the `<ROM>` attribute `subType` needs to be `'GaussPolynomialRom'` (see the example below). The GaussPolynomialRom is dependent on specific sampling; thus, this ROM cannot be trained unless a SparseGridCollocation or similar Sampler specifies this ROM in its input and is sampled in a MultiRun step. In addition to the common `<Target>` and `<Features>`, this ROM requires two more nodes and can accept multiple entries of a third optional node.

- `<IndexSet>`, *string, required field*, specifies the rules by which to construct multidimensional polynomials. The options are `'TensorProduct'`, `'TotalDegree'`, `'HyperbolicCross'`, and `'Custom'`. Total degree is efficient for uncertain inputs with a large degree of regularity, while hyperbolic cross is more efficient for low-regularity input spaces. If `'Custom'` is chosen, the `<IndexPoints>` is required.
- `<PolynomialOrder>`, *integer, required field*, indicates the maximum polynomial order in any one dimension to use in the polynomial chaos expansion. **Note:** If non-equal importance weights are supplied in the optional `<Interpolation>` node, the actual polynomial order in dimensions with high importance might exceed this value; however, this value is still used to limit the relative overall order.
- `<IndexPoints>`, *list of tuples, required field*, used to specify the index set points in a `'Custom'` index set. The tuples are entered as comma-separated values between parenthesis, with each tuple separated by a comma. Any amount of whitespace is acceptable. For example, `<IndexPoints> (0, 1), (0, 2), (1, 1), (4, 0) </IndexPoints>` **Note:** Using custom index sets does not guarantee accurate convergence.
- `<Interpolation>`, *string, optional field*, offers the option to specify quadrature, polynomials, and importance weights for the given variable name. The ROM accepts any number of `<Interpolation>` nodes up to the dimensionality of the input space. This node accepts several attributes, all of which are optional and default to the code-defined optimal choices based on the input dimension uncertainty distribution:
  - `quad`, *string, optional field*, specifies the quadrature type to use for collocation in this dimension. The default options depend on the uncertainty distribution of the input dimension, as shown in Table 1. Additionally, Clenshaw Curtis quadrature can be used for any distribution that doesn't include an infinite bound.  
*Default: see Table 1.* **Note:** For an uncertain distribution aside from the four listed on Table 1, this ROM makes use of the uniform-like range of the distribution's CDF



to apply quadrature that is suited uniform uncertainty (Legendre). It converges more slowly than the four listed, but are viable choices. Choosing polynomial type Legendre for any non-uniform distribution will enable this formulation automatically.

- **poly**, *string, optional field*, specifies the interpolating polynomial family to use for the polynomial expansion in this dimension. The default options depend on the quadrature type chosen, as shown in Table 1. Currently, no polynomials are available outside the default.

*Default: see Table 1.*

- **weight**, *float, optional field*, delineates the importance weighting of this dimension. A larger importance weight will result in increased resolution for this dimension at the cost of resolution in lower-weighted dimensions. The algorithm normalizes weights at run-time.

*Default: 1.*

Unc. Distribution	Default Quadrature	Default Polynomials
Uniform	Legendre	Legendre
Normal	Hermite	Hermite
Gamma	Laguerre	Laguerre
Beta	Jacobi	Jacobi
Other	Legendre*	Legendre*

**Table 1.** GaussPolynomialRom defaults

**Note:** This ROM type must be trained from a collocation quadrature set. Thus, it can only be trained from the outcomes of a SparseGridCollocation sampler. Also, this ROM must be referenced in the SparseGridCollocation sampler in order to accurately produce the necessary sparse grid points to train this ROM.

**Example:**

```
<Simulation>
...
<Samplers>
...
<SparseGridCollocation name="mySG" parallel="0">
  <variable name="x1">
    <distribution>myDist1</distribution>
  </variable>
  <variable name="x2">
    <distribution>myDist2</distribution>
  </variable>
  <ROM class = 'Models' type = 'ROM' >myROM</ROM>
</SparseGridCollocation>
...
```

```

</Samplers>
...
<Models>
...
  <ROM name='myRom' subType='GaussPolynomialRom'>
    <Target>ans</Target>
    <Features>x1, x2</Features>
    <IndexSet>TotalDegree</IndexSet>
    <PolynomialOrder>4</PolynomialOrder>
    <Interpolation quad='Legendre' poly='Legendre'
      weight='1'>x1</Interpolation>
    <Interpolation quad='ClenshawCurtis' poly='Jacobi'
      weight='2'>x2</Interpolation>
  </ROM>
...
</Models>
...
</Simulation>

```

When Printing this ROM via an OutStreamManager (see 14.1), the available metrics are:

- **'mean'**, the mean value of the ROM output within the input space it was trained,
- **'variance'**, the variance of the ROM output within the input space it was trained,
- **'numRuns'**, the number of distinct model runs required to construct the ROM,
- **'polyCoeffs'**, the polynomial expansion coefficients (PCE moments) of the ROM. These are listed by each polynomial combination, with the polynomial order tags listed in the order of the variables shown in the XML print.

### 15.3.3 HDMRRom

The HDMRRom sub-type contains a single ROM type, based on a Sobol decomposition scheme. In Sobol decomposition, also known as high-density model reduction (HDMR, specifically Cut-HDMR), a model is approximated as the sum of increasing-complexity interactions. At its lowest level (order 1), it treats the function as a sum of the reference case plus a functional of each input dimension separately. At order 2, it adds functionals to consider the pairing of each dimension with each other dimension. The benefit to this approach is considering several functions of small input cardinality instead of a single function with large input cardinality. This allows reduced order models like generalized polynomial chaos (see 15.3.2) to approximate the functionals accurately with few computations runs. In order to use this ROM, the `<ROM>` attribute `subType` needs to be **'HDMRRom'** (see the example below). The HDMRRom is dependent on specific sampling; thus, this ROM cannot be trained unless a Sobol or similar Sampler specifies this ROM in its input and is sampled in a MultiRun step. In addition to the common `<Target>` and `<Features>`, this

ROM requires the same nodes as the GaussPolynomialRom (see 15.3.2. Additionally, this ROM requires the `<SobolOrder>` node.

- `<SobolOrder>`, *integer, required field*, indicates the maximum cardinality of the input space used in the subset functionals. For example, order 1 includes only functionals of each independent dimension separately, while order 2 considers pair-wise interactions.

**Note:** This ROM type must be trained from a Sobol decomposition training set. Thus, it can only be trained from the outcomes of a Sobol sampler. Also, this ROM must be referenced in the Sobol sampler in order to accurately produce the necessary sparse grid points to train this ROM. Experience has shown order 2 Sobol decompositions to include the great majority of uncertainty in most models.

### Example:

```

<Samplers>
...
<Sobol name="mySobol" parallel="0">
  <variable name="x1">
    <distribution>myDist1</distribution>
  </variable>
  <variable name="x2">
    <distribution>myDist2</distribution>
  </variable>
  <ROM class = 'Models' type = 'ROM' >myHDMR</ROM>
</Sobol>
...
</Samplers>
...
<Models>
...
<ROM name='myHDMR' subType='HDMRRom'>
  <Target>ans</Target>
  <Features>x1, x2</Features>
  <SobolOrder>2</SobolOrder>
  <IndexSet>TotalDegree</IndexSet>
  <PolynomialOrder>4</PolynomialOrder>
  <Interpolation quad='Legendre' poly='Legendre'
    weight='1'>x1</Interpolation>
  <Interpolation quad='ClenshawCurtis' poly='Jacobi'
    weight='2'>x2</Interpolation>
</ROM>
...
</Models>

```

When Printing this ROM via an OutStreamManager (see 14.1), the available metrics are:

- `'mean'`, the mean value of the ROM output within the input space it was trained,

- **'variance'**, the roughly-approximated variance of the ROM output within the input space it was trained. It should be noted this is calculated as the sum of the variance of the constituent terms, which neglects cross terms since cut-HDMR terms are not orthogonal. Work is being done to fix this.
- **'numRuns'**, the number of distinct model runs required to construct the ROM,
- **'indices'**, the relative impact parameter of each subset to the total variance (these are not Sobol sensitivity indices, but relative impacts of the impact subspace on the variance of the whole expansion).

### 15.3.4 MSR

The MSR sub-type contains a class of ROMs that perform a topological decomposition of the data into approximately monotonic regions and fits weighted linear patches to the identified monotonic regions of the input space. Query points have estimated probabilities that they belong to each cluster. These probabilities can either be used to give a smooth, weighted prediction based on the associated linear models, or a hard classification to a particular local linear model which is then used for prediction. Currently, the probability prediction can be done using kernel density estimation (KDE) or through a one-versus-one support vector machine (SVM).

In order to use this ROM, the `<ROM>` attribute `subType` needs to be **'MSR'** (see the associated example). This model can be initialized with the following children:

- **<persistence>**, *string, optional field*, specifies how to define the hierarchical simplification by assigning a value to each local minimum and maximum according to the one of the strategy options below:
  - `difference` - The function value difference between the extremum and its closest-valued neighboring saddle.
  - `probability` - The probability integral computed as the sum of the probability of each point in a cluster divided by the count of the cluster.
  - `count` - The count of points that flow to or from the extremum.

*Default: difference*

- **<gradient>**, *string, optional field*, specifies the method used for estimating the gradient, available options are:
  - `steepest`

*Default: steepest*

- **<simplification>**, *float, optional field*, specifies the amount of noise reduction to apply before returning labels.

*Default: 0*

- **<graph>**, *string, optional field*, specifies the type of neighborhood graph used in the algorithm, available options are:

- beta skeleton
- relaxed beta skeleton
- approximate knn

*Default: beta skeleton*

- **<beta>**, *float in the range: (0,2], optional field*, is only used when the **<graph>** is set to beta skeleton or relaxed beta skeleton.

*Default: 1.0*

- **<knn>**, *integer, optional field*, is the number of neighbors when using the approximate knn for the **<graph>** sub-node and used to speed up the computation of other graphs by using the approximate knn graph as a starting point for pruning. -1 means use a fully connected graph.

*Default: -1*

- **<partitionPredictor>**, *string, optional*, a flag that specifies how the predictions for query point classification should be performed. Available options are:

- kde
- svm

*Default: kde*

- **<smooth>**, if this node is present, the ROM will blend the estimates of all of the local linear models weighted by the probability the query point is classified as belonging to that partition of the input space.

- **<kernel>**, *string, optional field*, this option is only used when the **<partitionPredictor>** is set to kde and specifies the type of kernel to use in the kernel density estimation. Available options are:

- uniform
- triangular
- gaussian

- epanechnikov
- biweight or quartic
- triweight
- tricube
- cosine
- logistic
- silverman
- exponential

*Default: gaussian*

- **<bandwidth>** *float or string, optional field*, this option is only used when the **<partitionPredictor>** is set to kde and specifies the scale of the fall-off. A higher bandwidth implies a smoother blending. If set to *variable*, then the bandwidth will be set to the distance of the *k*-nearest neighbor of the query point where *k* is set by the **<knn>** parameter.

*Default: 1.*

### Example:

```

<Simulation>
...
<Models>
...
</ROM>
<ROM name='aUserDefinedName' subType='MSR'>
  <Features>var1,var2,var3</Features>
  <Target>result1,result2</Target>
  <!-- <weighted>true</weighted> -->
  <simplification>0.0</simplification>
  <persistence>difference</persistence>
  <gradient>steepest</gradient>
  <graph>beta skeleton</graph>
  <beta>1</beta>
  <knn>8</knn>
  <partitionPredictor>kde</partitionPredictor>
  <kernel>gaussian</kernel>
  <smooth/>
  <bandwidth>0.2</bandwidth>
</ROM>
...

```

```
</Models>
...
</Simulation>
```

### 15.3.5 NDinvDistWeight

The NDinvDistWeight sub-type contains a single ROM type, based on an  $N$ -dimensional inverse distance weighting formulation. Inverse distance weighting (IDW) is a type of deterministic method for multivariate interpolation with a known scattered set of points. The assigned values to unknown points are calculated via a weighted average of the values available at the known points.

In order to use this Reduced Order Model, the `<ROM>` attribute `subType` needs to be `xml-StringNDinvDistWeight` (see the example below). This model can be initialized with the following child:

- `<p>`, *integer, required field*, must be greater than zero and represents the “power parameter”. For the choice of value for `<p>`, it is necessary to consider the degree of smoothing desired in the interpolation/extrapolation, the density and distribution of samples being interpolated, and the maximum distance over which an individual sample is allowed to influence the surrounding ones (lower  $p$  means greater importance for points far away).

#### Example:

```
<Simulation>
...
  <Models>
    ...
    <ROM name='aUserDefinedName' subType='NDinvDistWeight'>
      <Features>var1,var2,var3</Features>
      <Target>result1,result2</Target>
      <p>3</p>
    </ROM>
    ...
  </Models>
  ...
</Simulation>
```

### 15.3.6 SciKitLearn

The SciKitLearn sub-type represents the container of several ROMs available in RAVEN through the external library scikit-learn [4].

In order to use this Reduced Order Model, the `<ROM>` attribute `subType` needs to be `'SciKitLearn'` (i.e. `subType='SciKitLearn'`). The specifications of a `'SciKitLearn'` ROM depend on the value assumed by the following sub-node within the main `<ROM>` XML node:

- `<SKLtype>`, *vertical bar ( | ) separated string, required field*, contains a string that represents the ROM type to be used. As mentioned, its format is:  
`<SKLtype>mainSKLclass|algorithm</SKLtype>` where the first word (before the “|” symbol) represents the main class of algorithms, and the second word (after the “|” symbol) represents the specific algorithm.

Based on the `<SKLtype>` several different algorithms are available. In the following paragraphs a brief explanation and the input requirements are reported for each of them.

#### 15.3.6.1 Linear Models

The LinearModels’ algorithms implement generalized linear models. They include Ridge regression, Bayesian regression, lasso, and elastic net estimators computed with least angle regression and coordinate descent. This class also implements stochastic gradient descent related algorithms. In the following, all of the linear models available in RAVEN are reported.

##### 15.3.6.1.1 Linear Model: Automatic Relevance Determination Regression

The *Automatic Relevance Determination* (ARD) regressor is a hierarchical Bayesian approach where hyperparameters explicitly represent the relevance of different input features. These relevance hyperparameters determine the range of variation for the parameters relating to a particular input, usually by modelling the width of a zero-mean Gaussian prior on those parameters. If the width of the Gaussian is zero, then those parameters are constrained to be zero, and the corresponding input cannot have any effect on the predictions, therefore making it irrelevant. ARD optimizes these hyperparameters to discover which inputs are relevant. In order to use the *Automatic Relevance Determination regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|ARDRegression</SKLtype>.
```

In addition to this XML node, several others are available: .

- `<n_iter>`, *integer, optional field*, is the maximum number of iterations.  
*Default: 300*



- **<tol>**, *float, optional field*, stop the algorithm if the convergence error felt below the tolerance specified here.  
*Default: 1.e-3*
- **<alpha\_1>**, *float, optional field*, is a shape hyperparameter for the Gamma distribution prior over the  $\alpha$  parameter.  
*Default: 1.e-6*
- **<alpha\_2>**, *float, optional field*, inverse scale hyperparameter (rate parameter) for the Gamma distribution prior over the  $\alpha$  parameter.  
*Default: 1.e-6*
- **<lambda\_1>**, *float, optional field*, shape hyperparameter for the Gamma distribution prior over the  $\lambda$  parameter.  
*Default: 1.e-6*
- **<lambda\_2>**, *float, optional field*, inverse scale hyperparameter (rate parameter) for the Gamma distribution prior over the  $\lambda$  parameter.  
*Default: 1.e-6*
- **<compute\_score>**, *boolean, optional field*, if True, compute the objective function at each step of the model.  
*Default: False*
- **<threshold\_lambda>**, *float, optional field*, specifies the threshold for removing (pruning) weights with high precision from the computation.  
*Default: 1.e+4*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: False*
- **<verbose>**, *boolean, optional field*, use verbose mode when fitting the model.  
*Default: False*

### 15.3.6.1.2 Linear Model: Bayesian ridge regression

The *Bayesian ridge regression* estimates a probabilistic model of the regression problem as described above. The prior for the parameter  $w$  is given by a spherical Gaussian:

$$p(w|\lambda) = \mathcal{N}(w|0, \lambda^{-1}\mathbf{I}_p) \quad (1)$$

The priors over  $\alpha$  and  $\lambda$  are chosen to be gamma distributions, the conjugate prior for the precision of the Gaussian. The resulting model is called Bayesian ridge regression, and is similar to the classical ridge regression. The parameters  $w$ ,  $\alpha$ , and  $\lambda$  are estimated jointly during the fit of the model. The remaining hyperparameters are the parameters of the gamma priors over  $\alpha$  and  $\lambda$ . These are usually chosen to be non-informative. The parameters are estimated by maximizing the marginal log likelihood. In order to use the *Bayesian ridge regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|BayesianRidge</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_iter>**, *integer, optional field*, is the maximum number of iterations.  
*Default: 300*
- **<tol>**, *float, optional field*, stop the algorithm if the convergence error fell below the tolerance specified here.  
*Default: 1.e-3*
- **<alpha\_1>**, *float, optional field*, is a shape hyperparameter for the Gamma distribution prior over the  $\alpha$  parameter.  
*Default: 1.e-6*
- **<alpha\_2>**, *float, optional field*, inverse scale hyperparameter (rate parameter) for the Gamma distribution prior over the  $\alpha$  parameter.  
*Default: 1.e-6*
- **<lambda\_1>**, *float, optional field*, shape hyperparameter for the Gamma distribution prior over the  $\lambda$  parameter.  
*Default: 1.e-6*
- **<lambda\_2>**, *float, optional field*, inverse scale hyperparameter (rate parameter) for the Gamma distribution prior over the  $\lambda$  parameter.  
*Default: 1.e-6*
- **<compute\_score>**, *boolean, optional field*, if True, compute the objective function at each step of the model.  
*Default: False*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*

- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: False*
- **<verbose>**, *boolean, optional field*, use verbose mode when fitting the model.  
*Default: False*

### 15.3.6.1.3 Linear Model: Elastic Net

The *Elastic Net* is a linear regression technique with combined L1 and L2 priors as regularizers. It minimizes the objective function:

$$1/(2 * n_{samples}) * ||y - Xw||_2^2 + alpha * l1\_ratio * ||w||_1 + 0.5 * alpha * (1 - l1\_ratio) * ||w||_2^2 \quad (2)$$

In order to use the *Elastic Net regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|ElasticNet</SKLtype>.
```

In addition to this XML node, several others are available:

- **<alpha>**, *float, optional field*, specifies a constant that multiplies the penalty terms.  $alpha = 0$  is equivalent to an ordinary least square, solved by the **LinearRegression** object.  
*Default: 1.0*
- **<l1\_ratio>**, *float, optional field*, specifies the ElasticNet mixing parameter, with  $0 \leq l1\_ratio \leq 1$ . For  $l1\_ratio = 0$  the penalty is an L2 penalty. For  $l1\_ratio = 1$  it is an L1 penalty. For  $0 < l1\_ratio < 1$ , the penalty is a combination of L1 and L2.  
*Default: 0.5*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: False*
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: 1000*
- **<tol>**, *float, optional field*, specifies the tolerance for the optimization: if the updates are smaller than tol, the optimization code checks the dual gap for optimality and continues until it is smaller than tol.  
*Default: 1.e-4*

- **<warm\_start>**, *boolean, optional field*, when set to True, the model will reuse the solution of the previous call to fit as initialization, otherwise, it will just erase the previous solution.  
*Default: False*
- **<positive>**, *float, optional field*, when set to True, this forces the coefficients to be positive.  
*Default: False*

#### 15.3.6.1.4 Linear Model: Elastic Net CV

The *Elastic Net CV* is a linear regression similar to the Elastic Net model but with an iterative fitting along a regularization path. The best model is selected by cross-validation.

In order to use the *Elastic Net CV regressor*, the user needs to set the sub-node:

**<SKLtype>**linear\_model|ElasticNetCV**</SKLtype>**.

In addition to this XML node, several others are available:

- **<l1\_ratio>**, *float, optional field*, Float flag between 0 and 1 passed to ElasticNet (scaling between l1 and l2 penalties). For  $l1\_ratio = 0$  the penalty is an L2 penalty. For  $l1\_ratio = 1$  it is an L1 penalty. For  $0 < l1\_ratio < 1$ , the penalty is a combination of L1 and L2. This parameter can be a list, in which case the different values are tested by cross-validation and the one giving the best prediction score is used. Note that a good choice of list of values for  $l1\_ratio$  is often to put more values close to 1 (i.e. Lasso) and less close to 0 (i.e. Ridge), as in [.1, .5, .7, .9, .95, .99, 1].  
*Default: 0.5*
- **<eps>**, *float, optional field*, specifies the length of the path.  $eps=1e-3$  means that  $alpha\_min/alpha\_max = 1e - 3$ .  
*Default: 0.001*
- **<n\_alphas>**, *integer, optional field*, is the number of alphas along the regularization path used for each  $l1\_ratio$ .  
*Default: 100*
- **<precompute>**, *boolean or string, optional field*, determines whether to use a precomputed Gram matrix to speed up calculations. If set to 'auto,' RAVEN will decide. The Gram matrix can also be passed as an argument. Available options are [True — False — 'auto' — array-like].  
*Default: 'auto'*
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: 1000*

- **<tol>**, *float, optional field*, specifies the tolerance for the optimization: if the updates are smaller than tol, the optimization code checks the dual gap for optimality and continues until it is smaller than tol.  
*Default: 1.e-4*
- **<positive>**, *float, optional field*, when set to True, this forces the coefficients to be positive.  
*Default: False*

### 15.3.6.1.5 Linear Model: Least Angle Regression model

The *Least Angle Regression model* (LARS) is a regression algorithm for high-dimensional data. The LARS algorithm provides a means of producing an estimate of which variables to include, as well as their coefficients, when a response variable is determined by a linear combination of a subset of potential covariates.

In order to use the *Least Angle Regression model*, the user needs to set the sub-node:

```
<SKLtype>linear_model|Lars</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_nonzero\_coef>**, *integer, optional field*, represents the target number of non-zero coefficients.  
*Default: 500*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<verbose>**, *boolean, optional field*, use verbose mode when fitting the model.  
*Default: False*
- **<precompute>**, *boolean or string, optional field*, determines whether to use a precomputed Gram matrix to speed up calculations. If set to 'auto,' RAVEN will decide. The Gram matrix can also be passed as an argument. Available options are [True — False — 'auto' — array-like].  
*Default: 'auto'*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: True*

- **<eps>**, *float, optional field*, represents the machine precision regularization in the computation of the Cholesky diagonal factors. Increase this for very ill-conditioned systems. Unlike the **<tol>** parameter in some iterative optimization-based algorithms, this parameter does not control the tolerance of the optimization.  
*Default: 2.2204460492503131e-16*
- **<fit\_path>**, *boolean, optional field*, if True the full path is stored in the `coef_path` attribute. If you compute the solution for a large problem or many targets, setting `fit_path` to False will lead to a speedup, especially with a small alpha.  
*Default: True*

#### 15.3.6.1.6 Linear Model: Cross-validated Least Angle Regression model

The *Cross-validated Least Angle Regression model* is a regression algorithm for high-dimensional data. It is similar to the LARS method, but the best model is selected by cross-validation. In order to use the *Cross-validated Least Angle Regression model*, the user needs to set the sub-node:

```
<SKLtype>linear_model|LarsCV</SKLtype>.
```

In addition to this XML node, several others are available:

- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<verbose>**, *boolean, optional field*, use verbose mode when fitting the model.  
*Default: False*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: True*
- **<precompute>**, *boolean or string, optional field*, determines whether to use a precomputed Gram matrix to speed up calculations. If set to 'auto,' RAVEN will decide. The Gram matrix can also be passed as an argument. Available options are [True — False — 'auto' — array-like].  
*Default: 'auto'*
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: 500*
- **<max\_n\_alphas>**, *integer, optional field*, specifies the maximum number of points on the path used to compute the residuals in the cross-validation.  
*Default: 1000*

- **<eps>**, *float, optional field*, represents the machine-precision regularization in the computation of the Cholesky diagonal factors. Increase this for very ill-conditioned systems. Unlike the *tol* parameter in some iterative optimization-based algorithms, this parameter does not control the tolerance of the optimization.  
*Default: 2.2204460492503131e-16*

### 15.3.6.1.7 Linear Model trained with L1 prior as regularizer (aka the Lasso)

The *Linear Model trained with L1 prior as regularizer (Lasso)* is a shrinkage and selection method for linear regression. It minimizes the usual sum of squared errors, with a bound on the sum of the absolute values of the coefficients. In order to use the *Linear Model trained with L1 prior as regularizer (Lasso)*, the user needs to set the sub-node:

```
<SKLtype>linear_model|Lasso</SKLtype>.
```

In addition to this XML node, several others are available:

- **<alpha>**, *float, optional field*, sets a constant multiplier for the L1 term.  $\alpha = 0$  is equivalent to an ordinary least square, solved by the LinearRegression object. For numerical reasons, using  $\alpha = 0$  with the Lasso object is not advised and you should instead use the LinearRegression object.  
*Default: 1.0*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: False*
- **<precompute>**, *boolean or string, optional field*, determines whether to use a precomputed Gram matrix to speed up calculations. If set to 'auto,' RAVEN will decide. The Gram matrix can also be passed as an argument. Available options are [True — False — 'auto' — array-like].  
*Default: False* **Note:** For sparse input this option is always True to preserve sparsity.
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: 1000*
- **<tol>**, *float, optional field*, specifies the tolerance for the optimization: if the updates are smaller than tol, the optimization code checks the dual gap for optimality and continues until it is smaller than tol.  
*Default: 1.e-4*

- **<warm\_start>**, *boolean, optional field*, when set to True, the model will reuse the solution of the previous call to fit as initialization, otherwise, it will just erase the previous solution.  
*Default: False*
- **<positive>**, *float, optional field*, when set to True, this forces the coefficients to be positive.  
*Default: False*

### 15.3.6.1.8 Lasso linear model with iterative fitting along a regularization path

The *Lasso linear model with iterative fitting along a regularization path* is an algorithm of the Lasso family, that computes the linear regressor weights, identifying the regularization path in an iterative fitting (see <http://www.jstatsoft.org/v33/i01/paper>)

In order to use the *Lasso linear model with iterative fitting along a regularization path regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|LassoCV</SKLtype>.
```

In addition to this XML node, several others are available:

- **<eps>**, *float, optional field*, represents the length of the path.  $\text{eps}=1\text{e-}3$  means that  $\text{alpha\_min} / \text{alpha\_max} = 1\text{e-}3$ .  
*Default: 1.0e-3*
- **<n\_alphas>**, *int, optional field*, sets the number of alphas along the regularization path.  
*Default: 100*
- **<alphas>**, *numpy array, optional field*, lists the locations of the alphas used to compute the models.  
*Default: None* If None, alphas are set automatically.
- **<precompute>**, *boolean or string, optional field*, determines whether to use a precomputed Gram matrix to speed up calculations. If set to 'auto,' RAVEN will decide. The Gram matrix can also be passed as an argument. Available options are [True — False — 'auto' — array-like].  
*Default: 'auto'*
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: 1000*
- **<tol>**, *float, optional field*, specifies the tolerance for the optimization: if the updates are smaller than tol, the optimization code checks the dual gap for optimality and continues until it is smaller than tol.  
*Default: 1.e-4*



- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: False*
- **<positive>**, *float, optional field*, when set to True, this forces the coefficients to be positive.  
*Default: False*

### 15.3.6.1.9 Lasso model fit with Least Angle Regression

*Lasso model fit with Least Angle Regression* (aka Lars) It is a Linear Model trained with an L1 prior as regularizer. In order to use the *Least Angle Regression model regressor*, the user needs to set the sub-node `linear_model`. In order to use the *Least Angle Regression model regressor*, the user needs to set the sub-node `LassoLars`:

```
<SKLtype>linear_model|LassoLars</SKLtype>.
```

In addition to this XML node, several others are available:

- **<alpha>**, *float, optional field*, specifies a constant that multiplies the penalty terms.  $\alpha = 0$  is equivalent to an ordinary least square, solved by the **LinearRegression** object.  
*Default: 1.0*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: False*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: True*
- **<precompute>**, *boolean or string, optional field*, determines whether to use a precomputed Gram matrix to speed up calculations. If set to 'auto,' RAVEN will decide. The Gram matrix can also be passed as an argument. Available options are [True — False — 'auto' — array-like].  
*Default: 'auto'*
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: 500*

- **<eps>**, *float, optional field*, sets the machine precision regularization in the computation of the Cholesky diagonal factors. Increase this for very ill-conditioned systems.  
*Default: 2.2204460492503131e-16*

#### 15.3.6.1.10 Cross-validated Lasso, using the LARS algorithm

The *Cross-validated Lasso, using the LARS algorithm* is a cross-validated Lasso, using the LARS algorithm.

In order to use the *Cross-validated Lasso, using the LARS algorithm regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|LassoLarsCV</SKLtype>.
```

In addition to this XML node, several others are available:

- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: False*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: True*
- **<precompute>**, *boolean or string, optional field*, determines whether to use a precomputed Gram matrix to speed up calculations. If set to 'auto,' RAVEN will decide. The Gram matrix can also be passed as an argument. Available options are [True — False — 'auto' — array-like].  
*Default: 'auto'*
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: 500*
- **<max\_n\_alphas>**, *integer, optional field*, specifies the maximum number of points on the path used to compute the residuals in the cross-validation.  
*Default: 1000*
- **<eps>**, *float, optional field*, specifies the machine precision regularization in the computation of the Cholesky diagonal factors. Increase this for very ill-conditioned systems.  
*Default: 2.2204460492503131e-16*

### 15.3.6.1.11 Lasso model fit with Lars using BIC or AIC for model selection

The *Lasso model fit with Lars using BIC or AIC for model selection* is a Lasso model fit with Lars using BIC or AIC for model selection. The optimization objective for Lasso is:  $(1/(2 * n\_samples)) * ||y - Xw||_2^2 + alpha * ||w||_1$  AIC is the Akaike information criterion and BIC is the Bayes information criterion. Such criteria are useful in selecting the value of the regularization parameter by making a trade-off between the goodness of fit and the complexity of the model. A good model explains the data well while maintaining simplicity. In order to use the *Lasso model fit with Lars using BIC or AIC for model selection regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|LassoLarsIC</SKLtype>.
```

In addition to this XML node, several others are available:

- **<criterion>**, *'bic' — 'aic'*, specifies the type of criterion to use.  
*Default: 'aic'*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: False*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: True*
- **<precompute>**, *boolean or string, optional field*, determines whether to use a precomputed Gram matrix to speed up calculations. If set to 'auto,' RAVEN will decide. The Gram matrix can also be passed as an argument. Available options are [True — False — 'auto' — array-like].  
*Default: 'auto'*
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: 500*
- **<eps>**, *float, optional field*, represents the machine precision regularization in the computation of the Cholesky diagonal factors. Increase this for very ill-conditioned systems. Unlike the tol parameter in some iterative optimization-based algorithms, this parameter does not control the tolerance of the optimization.  
*Default: 2.2204460492503131e-16*

### 15.3.6.1.12 Ordinary least squares Linear Regression

The *Ordinary least squares Linear Regression* is a method for estimating the unknown parameters in a linear regression model, with the goal of minimizing the differences between the observed responses in some arbitrary dataset and the responses predicted by the linear approximation of the data. In order to use the *Ordinary least squares Linear Regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|LinearRegression</SKLtype>.
```

In addition to this XML node, several others are available:

- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: False*

### 15.3.6.1.13 Logistic Regression

The *Logistic Regression* implements L1 and L2 regularized logistic regression using the liblinear library. It can handle both dense and sparse input. This regressor uses C-ordered arrays or CSR matrices containing 64-bit floats for optimal performance; any other input format will be converted (and copied). In order to use the *Logistic Regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|LogisticRegression</SKLtype>.
```

In addition to this XML node, several others are available:

- **<penalty>**, *string, 'l1' or 'l2'*, specifies the norm used in the penalization.  
*Default: 'l2'*
- **<dual>**, *boolean*, specifies the dual or primal formulation. Dual formulation is only implemented for the l2 penalty. Prefer dual=False when n\_samples > n\_features.  
*Default: False*
- **<C>**, *float, optional field*, is the inverse of the regularization strength; must be a positive float. Like in support vector machines, smaller values specify stronger regularization.  
*Default: 1.0*
- **<fit\_intercept>**, *boolean*, specifies if a constant (a.k.a. bias or intercept) should be added to the decision function.  
*Default: True*

- **<intercept\_scaling>**, *float, optional field*, when `self.fit_intercept` is `True`, instance vector `x` becomes `[x, self.intercept_scaling]`, i.e. a “synthetic” feature with constant value equal to `intercept_scaling` is appended to the instance vector. The intercept becomes `intercept_scaling * synthetic feature weight`. **Note:** The synthetic feature weight is subject to L1/L2 regularization as are all other features. To lessen the effect of regularization on synthetic feature weight (and therefore on the intercept) `intercept_scaling` has to be increased.  
*Default: 1.0*
- **<class\_weight>**, *dict, or 'balanced', optional* Weights associated with classes in the form `{class_label: weight}`. If not given, all classes are supposed to have weight one. The “balanced” mode uses the values of `y` to automatically adjust weights inversely proportional to class frequencies in the input data as `n_samples / (n_classes * np.bincount(y))` Note that these weights will be multiplied with `sample_weight` (passed through the fit method) if `sample_weight` is specified. New in version 0.17: `class_weight='balanced'` instead of deprecated `class_weight='auto'`.  
*Default: None*
- **<random\_state>**, *int seed, RandomState instance, or None*, sets the seed of the pseudo random number generator to use when shuffling the data.  
*Default: None*
- **<tol>**, *float, optional field*, specifies the tolerance for stopping criteria.  
*Default: 0.0001*

#### 15.3.6.1.14 Multi-task Lasso model trained with L1/L2 mixed-norm as regularizer

The *Multi-task Lasso model trained with L1/L2 mixed-norm as regularizer* is a regressor where the optimization objective for Lasso is:  $(1/(2*n\_samples))*||Y - XW||_{Fro}^2 + alpha*||W||_{21}$  Where:  $||W||_{21} = \sum_i \sqrt{\sum_j w_{ij}^2}$  i.e. the sum of norm of each row. In order to use the *Multi-task Lasso model trained with L1/L2 mixed-norm as regularizer regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|MultiTaskLasso</SKLtype>.
```

In addition to this XML node, several others are available:

- **<alpha>**, *float, optional field*, sets the constant multiplier for the L1/L2 term.  
*Default: 1.0*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to `False`, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*

- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: False*
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: 1000*
- **<tol>**, *float, optional field*, specifies the tolerance for the optimization: if the updates are smaller than tol, the optimization code checks the dual gap for optimality and continues until it is smaller than tol.  
*Default: 1.e-4*
- **<warm\_start>**, *boolean, optional field*, when set to True, the model will reuse the solution of the previous call to fit as initialization, otherwise, it will just erase the previous solution.  
*Default: False*

### 15.3.6.1.15 Multi-task Elastic Net model trained with L1/L2 mixed-norm as regularizer

The *Multi-task Elastic Net model trained with L1/L2 mixed-norm as regularizer* is a regressor where the optimization objective for MultiTaskElasticNet is:  $(1/(2 * n\_samples)) * ||Y - XW||_2^{Fro} + alpha * l1\_ratio * ||W||_{21} + 0.5 * alpha * (1 - l1\_ratio) * ||W||_{Fro}^2$  Where:  $||W||_{21} = \sum_i \sqrt{\sum_j w_{ij}^2}$  i.e. the sum of norm of each row. In order to use the *Multi-task ElasticNet model trained with L1/L2 mixed-norm as regularizer regressor*, the user needs to set the sub-node:

**<SKLtype>**linear\_model|MultiTaskElasticNet**</SKLtype>**.

In addition to this XML node, several others are available:

- **<alpha>**, *float, optional field*, represents a constant multiplier for the L1/L2 term.  
*Default: 1.0*
- **<l1\_ratio>**, *float*, represents the Elastic Net mixing parameter, with  $0 < l1\_ratio \leq 1$ . For  $l1\_ratio = 0$  the penalty is an L1/L2 penalty. For  $l1\_ratio = 1$  it is an L1 penalty. For  $0 < l1\_ratio < 1$ , the penalty is a combination of L1/L2 and L2.  
*Default: 0.5*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: False*

- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.
- **<tol>**, *float, optional field*, specifies the tolerance for the optimization: if the updates are smaller than tol, the optimization code checks the dual gap for optimality and continues until it is smaller than tol.  
*Default: 1.e-4*
- **<warm\_start>**, *boolean, optional field*, when set to True, the model will reuse the solution of the previous call to fit as initialization, otherwise, it will just erase the previous solution.  
*Default: False*

### 15.3.6.1.16 Orthogonal Matching Pursuit model (OMP)

The *Orthogonal Matching Pursuit model (OMP)* is a type of sparse approximation which involves finding the “best matching” projections of multidimensional data onto an over-complete dictionary,  $D$ . In order to use the *Orthogonal Matching Pursuit model (OMP) regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|OrthogonalMatchingPursuit</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_nonzero\_coefs>**, *int, optional field*, represents the desired number of non-zero entries in the solution. If None, this value is set to 10% of n\_features.  
*Default: None*
- **<tol>**, *float, optional field*, specifies the maximum norm of the residual. If not None, overrides n\_nonzero\_coefs.  
*Default: None*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: True*
- **<precompute>**, *{True, False, 'auto'}*, specifies whether to use a precomputed Gram and Xy matrix to speed up calculations. Improves performance when n\_targets or n\_samples is very large. **Note:** If you already have such matrices, you can pass them directly to the fit method.  
*Default: 'auto'*



### 15.3.6.1.17 Cross-validated Orthogonal Matching Pursuit model (OMP)

The *Cross-validated Orthogonal Matching Pursuit model (OMP)* is a regressor similar to OMP which has good performance in sparse recovery. In order to use the *Cross-validated Orthogonal Matching Pursuit model (OMP) regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|OrthogonalMatchingPursuitCV</SKLtype>.
```

In addition to this XML node, several others are available:

- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: True*
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: None* Maximum number of iterations to perform, therefore maximum features to include 10% of n\_features but at least 5 if available.
- **<cv>**, *cross-validation generator, optional*, see sklearn.cross\_validation.  
*Default: None*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: False*

### 15.3.6.1.18 Passive Aggressive Classifier

The *Passive Aggressive Classifier* is a principled approach to linear classification that advocates minimal weight updates i.e., the least required to correctly classify the current training instance. In order to use the *Passive Aggressive Classifier*, the user needs to set the sub-node:

```
<SKLtype>linear_model|PassiveAggressiveClassifier</SKLtype>.
```

In addition to this XML node, several others are available:

- **<C>**, *float*, specifies the maximum step size (regularization).  
*Default: 1.0*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*



- **<n\_iter>**, *int, optional field*, specifies the number of passes over the training data (aka epochs).  
*Default: 5*
- **<shuffle>**, *boolean, optional field*, specifies whether or not the training data should be shuffled after each epoch.  
*Default: True*
- **<random\_state>**, *int seed, RandomState instance, or None*, sets the seed of the pseudo random number generator to use when shuffling the data.  
*Default: None*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: 0*
- **<loss>**, *string, optional field*, the loss function to be used:
  - hinge: equivalent to PA-I (<http://jmlr.csail.mit.edu/papers/volume7/crammer06a/crammer06a.pdf>)
  - squared\_hinge: equivalent to PA-II (<http://jmlr.csail.mit.edu/papers/volume7/crammer06a/crammer06a.pdf>)  
*Default: 'hinge'*
- **<warm\_start>**, *boolean, optional field*, when set to True, the model will reuse the solution of the previous call to fit as initialization, otherwise, it will just erase the previous solution.  
*Default: False*

### 15.3.6.1.19 Passive Aggressive Regressor

The *Passive Aggressive Regressor* is similar to the Perceptron in that it does not require a learning rate. However, contrary to the Perceptron, this regressor includes a regularization parameter,  $C$ .

In order to use the *Passive Aggressive Regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|PassiveAggressiveRegressor</SKLtype>.
```

In addition to this XML node, several others are available:

- **<C>**, *float*, sets the maximum step size (regularization).  
*Default: 1.0*
- **<epsilon>**, *float*, if the difference between the current prediction and the correct label is below this threshold, the model is not updated.  
*Default: 0.1*

- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<n\_iter>**, *int, optional field*, specifies the number of passes over the training data (aka epochs).  
*Default: 5*
- **<shuffle>**, *boolean, optional field*, specifies whether or not the training data should be shuffled after each epoch.  
*Default: True*
- **<random\_state>**, *int seed, RandomState instance, or None*, sets the seed of the pseudo random number generator to use when shuffling the data.  
*Default: None*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: 0*
- **<loss>**, *string, optional field*, specifies the loss function to be used:
  - `epsilon_insensitive`: equivalent to PA-I in the reference paper (<http://jmlr.csail.mit.edu/papers/volume7/crammer06a/crammer06a.pdf>).
  - `squared_epsilon_insensitive`: equivalent to PA-II in the reference paper (<http://jmlr.csail.mit.edu/papers/volume7/crammer06a/crammer06a.pdf>).  
*Default: 'epsilon\_insensitive'*
- **<warm\_start>**, *boolean, optional field*, when set to True, the model will reuse the solution of the previous call to fit as initialization, otherwise, it will just erase the previous solution.  
*Default: False*

### 15.3.6.1.20 Perceptron

The *Perceptron* method is an algorithm for supervised classification of an input into one of several possible non-binary outputs. It is a type of linear classifier, i.e. a classification algorithm that makes its predictions based on a linear predictor function combining a set of weights with the feature vector. The algorithm allows for online learning, in that it processes elements in the training set one at a time. In order to use the *Perceptron classifier*, the user needs to set the sub-node:

```
<SKLtype>linear_model|Perceptron</SKLtype>.
```

In addition to this XML node, several others are available:

- **<penalty>**, *None, 'l2' or 'l1' or 'elasticnet'*, defines the penalty (aka regularization term) to be used.  
*Default: None*
- **<alpha>**, *float*, sets the constant multiplier for the regularization term if regularization is used.  
*Default: 0.0001*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<n\_iter>**, *int, optional field*, specifies the number of passes over the training data (aka epochs).  
*Default: 5*
- **<shuffle>**, *boolean, optional field*, specifies whether or not the training data should be shuffled after each epoch.  
*Default: True*
- **<random\_state>**, *int seed, RandomState instance, or None*, sets the seed of the pseudo random number generator to use when shuffling the data.  
*Default: 0*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: 0*
- **<eta0>**, *double, optional field*, defines the constant multiplier for the updates.  
*Default: 1.0*
- **<class\_weight>**, *dict, {class\_label: weight} or "balanced" or None, optional* Preset for the class\_weight fit parameter. Weights associated with classes. If not given, all classes are supposed to have weight one. The "balanced" mode uses the values of y to automatically adjust weights inversely proportional to class frequencies in the input data as  $n\_samples / (n\_classes * np.bincount(y))$
- **<warm\_start>**, *boolean, optional field*, when set to True, the model will reuse the solution of the previous call to fit as initialization, otherwise, it will just erase the previous solution.  
*Default: False*

#### 15.3.6.1.21 Randomized Lasso

The *Randomized Lasso* works by resampling the training data and computing a Lasso on each resampling. In short, the features selected more often are good features. It is also known as

stability selection. In order to use the *Randomized Lasso regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|RandomizedLasso</SKLtype>.
```

In addition to this XML node, several others are available:

- **<alpha>**, *float, 'aic', or 'bic', optional*, defines the regularization parameter  $\alpha$  in the Lasso.  
*Default: 'aic'*
- **<scaling>**, *float, optional field*, specifies the alpha parameter used to randomly scale the. Should be between 0 and 1.  
*Default: 0.5*
- **<sample\_fraction>**, *float, optional field*, defines the fraction of samples to be used in each randomized design. Should be between 0 and 1. If 1, all samples are used.  
*Default: 0.75*
- **<n\_resampling>**, *int, optional field*, is the number of randomized models used.  
*Default: 200*
- **<selection\_threshold>**, *float, optional field*, sets the score above for which features should be selected.  
*Default: 0.25*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: False*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: True*
- **<precompute>**, *boolean or string, optional field*, determines whether to use a precomputed Gram matrix to speed up calculations. If set to 'auto,' RAVEN will decide. The Gram matrix can also be passed as an argument. Available options are [True — False — 'auto' — array-like].  
*Default: 'auto'*
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: 500*

- **<eps>**, *float, optional field*, defines the machine precision regularization in the computation of the Cholesky diagonal factors. Increase this for very ill-conditioned systems. Unlike the ‘tol’ parameter in some iterative optimization-based algorithms, this parameter does not control the tolerance of the optimization.  
*Default: 2.2204460492503131e-16*
- **<random\_state>**, *int, RandomState instance or None, optional*, if int, random\_state is the seed used by the random number generator; If RandomState instance, random\_state is the random number generator; If None, the random number generator is the RandomState instance used by np.random.  
*Default: None*

### 15.3.6.1.22 Randomized Logistic Regression

The *Randomized Logistic Regression* works by resampling the training data and computing a logistic regression on each resampling. In short, the features selected more often are good features. It is also known as stability selection. In order to use the *Randomized Logistic Regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|RandomizedLogisticRegression</SKLtype>.
```

In addition to this XML node, several others are available:

- **<C>**, *float, optional field*, specifies the regularization parameter  $C$  in the LogisticRegression.  
*Default: 1*
- **<scaling>**, *float, optional field (default=0.5)*, sets the alpha parameter used to randomly scale the features. Should be between 0 and 1.  
*Default: 0.5*
- **<sample\_fraction>**, *float, optional field*, is the fraction of samples to be used in each randomized design. Should be between 0 and 1. If 1, all samples are used.  
*Default: 0.75*
- **<n\_resampling>**, *int, optional field*, sets the number of randomized models.  
*Default: 200*
- **<selection\_threshold>**, *float, optional field*, sets the score above which features should be selected.  
*Default: 0.25*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is

expected to be already centered).

*Default: True*

- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: False*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: True*
- **<tol>**, *float, optional field*, specifies the tolerance for stopping criteria.  
*Default: 1=e-3*
- **<random\_state>**, *int, RandomState instance or None, optional*, if int, random\_state is the seed used by the random number generator; If RandomState instance, random\_state is the random number generator; If None, the random number generator is the RandomState instance used by np.random.  
*Default: None*

#### 15.3.6.1.23 Linear least squares with l2 regularization

The *Linear least squares with l2 regularization* solves a regression model where the loss function is the linear least squares function and the regularization is given by the l2-norm. Also known as Ridge Regression or Tikhonov regularization. This estimator has built-in support for multivariate regression (i.e., when y is a 2d-array of shape [n\_samples, n\_targets]). In order to use the *Linear least squares with l2 regularization*, the user needs to set the sub-node:

```
<SKLtype>linear_model|Ridge</SKLtype>.
```

In addition to this XML node, several others are available:

- **<alpha>**, *float, array-like*, shape = [n\_targets] Small positive values of alpha improve the conditioning of the problem and reduce the variance of the estimates. Alpha corresponds to  $(2 * C)^{-1}$  in other linear models such as LogisticRegression or LinearSVC. If an array is passed, penalties are assumed to be specific to the targets. Hence they must correspond in number.  
*Default: 1.0*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: determined by scipy.sparse.linalg.*

- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: False*
- **<solver>**, {‘auto’, ‘svd’, ‘cholesky’, ‘lsqr’, ‘sparse\_cg’}, specifies the solver to use in the computational routines:
  - ‘auto’ chooses the solver automatically based on the type of data.
  - ‘svd’ uses a singular value decomposition of X to compute the ridge coefficients. More stable for singular matrices than ‘cholesky.’
  - ‘cholesky’ uses the standard `scipy.linalg.solve` function to obtain a closed-form solution.
  - ‘sparse\_cg’ uses the conjugate gradient solver as found in `scipy.sparse.linalg.cg`. As an iterative algorithm, this solver is more appropriate than ‘cholesky’ for large-scale data (possibility to set `tol` and `max_iter`).
  - ‘lsqr’ uses the dedicated regularized least-squares routine `scipy.sparse.linalg.lsqr`. It is the fastest but may not be available in old `scipy` versions. It also uses an iterative procedure.

All three solvers support both dense and sparse data.

*Default: ‘auto’*

#### 15.3.6.1.24 Classifier using Ridge regression

The *Classifier using Ridge regression* is a classifier based on linear least squares with l2 regularization. In order to use the *Classifier using Ridge regression*, the user needs to set the sub-node:

**<SKLtype>**`linear_model|RidgeClassifier`**</SKLtype>**.

In addition to this XML node, several others are available:

- **<alpha>**, *float*, small positive values of alpha improve the conditioning of the problem and reduce the variance of the estimates. Alpha corresponds to  $(2 * C)^{-1}$  in other linear models such as `LogisticRegression` or `LinearSVC`.  
*Default: 1.0*
- **<class\_weight>**, *dict, optional field*, specifies weights associated with classes in the form `class_label: weight`. If not given, all classes are assumed to have weight one.  
*Default: None*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is

expected to be already centered).

*Default: True*

- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: determined by scipy.sparse.linalg.*
- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: False*
- **<solver>**, {*'auto', 'svd', 'cholesky', 'lsqr', 'sparse\_cg'*}, specifies the solver to use in the computational routines:
  - 'auto' chooses the solver automatically based on the type of data.
  - 'svd' uses a singular value decomposition of X to compute the ridge coefficients. More stable for singular matrices than 'cholesky.'
  - 'cholesky' uses the standard `scipy.linalg.solve` function to obtain a closed-form solution.
  - 'sparse\_cg' uses the conjugate gradient solver as found in `scipy.sparse.linalg.cg`. As an iterative algorithm, this solver is more appropriate than 'cholesky' for large-scale data (possibility to set `tol` and `max_iter`).
  - 'lsqr' uses the dedicated regularized least-squares routine `scipy.sparse.linalg.lsqr`. It is the fastest but may not be available in old `scipy` versions. It also uses an iterative procedure.

All three solvers support both dense and sparse data.

*Default: 'auto'*

- **<tol>**, *float*, defines the required precision of the solution.  
*Default: 0.001*

#### 15.3.6.1.25 Ridge classifier with built-in cross-validation

The *Ridge classifier with built-in cross-validation* performs Generalized Cross-Validation, which is a form of efficient leave-one-out cross-validation. Currently, only the `n_features > n_samples` case is handled efficiently. In order to use the *Ridge classifier with built-in cross-validation classifier*, the user needs to set the sub-node:

```
<SKLtype>linear_model|RidgeClassifierCV</SKLtype>.
```

In addition to this XML node, several others are available:

- **<alphas>**, *numpy array of shape [n\_alphas]*, is an array of alpha values to try. Small positive values of alpha improve the conditioning of the problem and reduce the variance of the



estimates. Alpha corresponds to  $(2 * C)^{-1}$  in other linear models such as LogisticRegression or LinearSVC.

*Default: (0.1, 1.0, 10.0)*

- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).

*Default: True*

- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.

*Default: False*

- **<scoring>**, *string, callable or None, optional*, is a string (see model evaluation documentation) or a scorer callable object / function with signature `scorer(estimator, X, y)`.

*Default: None*

- **<cv>**, *cross-validation generator, optional*, If None, Generalized Cross-Validation (efficient leave-one-out) will be used.

*Default: None*

- **<class\_weight>**, *dic, optional field*, specifies weights associated with classes in the form `class_label:weight`. If not given, all classes are supposed to have weight one.

*Default: None*

#### 15.3.6.1.26 Ridge regression with built-in cross-validation

The *Ridge regression with built-in cross-validation* performs Generalized Cross-Validation, which is a form of efficient leave-one-out cross-validation. In order to use the *Ridge regression with built-in cross-validation regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|RidgeCV</SKLtype>.
```

In addition to this XML node, several others are available:

- **<alphas>**, *numpy array of shape [n\_alphas]*, specifies an array of alpha values to try. Small positive values of alpha improve the conditioning of the problem and reduce the variance of the estimates. Alpha corresponds to  $(2 * C)^{-1}$  in other linear models such as LogisticRegression or LinearSVC.

*Default: (0.1, 1.0, 10.0)*

- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to False, no intercept will be used in the calculations (e.g. data is expected to be already centered).

*Default: True*

- **<normalize>**, *boolean, optional field*, if True, the regressors X will be normalized before regression.  
*Default: False*
- **<scoring>**, *string, callable or None, optional*, is a string (see model evaluation documentation) or a scorer callable object / function with signature `scorer(estimator, X, y)`.  
*Default: None*
- **<cv>**, *cross-validation generator, optional field*, if None, Generalized Cross-Validation (efficient leave-one-out) will be used.  
*Default: None*
- **<gcv\_mode>**, *{None, 'auto,' 'svd,' 'eigen'}, optional field*, is a flag indicating which strategy to use when performing Generalized Cross-Validation. Options are:
  - ‘auto:’ use svd if `n_samples < n_features` or when X is a sparse matrix, otherwise use eigen
  - ‘svd:’ force computation via singular value decomposition of X (does not work for sparse matrices)
  - ‘eigen:’ force computation via eigendecomposition of  $X^T X$

The ‘auto’ mode is the default and is intended to pick the cheaper option of the two depending upon the shape and format of the training data.

*Default: None*

- **<store\_cv\_values>**, *boolean*, is a flag indicating if the cross-validation values corresponding to each alpha should be stored in the `cv_values` attribute (see below). This flag is only compatible with `cv=None` (i.e. using Generalized Cross-Validation).  
*Default: False*

### 15.3.6.1.27 Linear classifiers (SVM, logistic regression, a.o.) with SGD training

The *Linear classifiers (SVM, logistic regression, a.o.) with SGD training* implements regularized linear models with stochastic gradient descent (SGD) learning: the gradient of the loss is estimated for each sample at a time and the model is updated along the way with a decreasing strength schedule (aka learning rate). SGD allows minibatch (online/out-of-core) learning, see the `partial_fit` method. This implementation works with data represented as dense or sparse arrays of floating point values for the features. The model it fits can be controlled with the `loss` parameter; by default, it fits a linear support vector machine (SVM). The regularizer is a penalty added to the loss function that shrinks model parameters towards the zero vector using either the squared Euclidean norm L2 or the absolute norm L1 or a combination of both (Elastic Net). If the parameter update crosses the 0.0 value because of the regularizer, the update is truncated to 0.0 to allow for learning sparse models and achieves online feature selection. In order to use the *Linear classifiers (SVM, logistic regression, a.o.) with SGD training*, the user needs to set the sub-node:

`<SKLtype>linear_model | SGDClassifier</SKLtype>`.

In addition to this XML node, several others are available:

- **<loss>**, *str*, *'hinge'*, *'log'*, *'modified\_huber'*, *'squared\_hinge'*, *'perceptron'*, or a *regression loss*: *'squared\_loss'*, *'huber'*, *'epsilon\_insensitive'*, or *'squared\_epsilon\_insensitive'*, dictates the loss function to be used. The available options are:
  - *'hinge'* gives a linear SVM.
  - *'log'* loss gives logistic regression, a probabilistic classifier.
  - *'modified\_huber'* is another smooth loss that brings tolerance to outliers as well as probability estimates.
  - *'squared\_hinge'* is like hinge but is quadratically penalized.
  - *'perceptron'* is the linear loss used by the perceptron algorithm.

The other losses are designed for regression but can be useful in classification as well; see SGDRegressor for a description.

*Default: 'hinge'*

- **<penalty>**, *str*, *'l2'* or *'l1'* or *'elasticnet'*, defines the penalty (aka regularization term) to be used. *'l2'* is the standard regularizer for linear SVM models. *'l1'* and *'elasticnet'* might bring sparsity to the model (feature selection) not achievable with *'l2'*.

*Default: 'l2'*

- **<alpha>**, *float*, is the constant multiplier for the regularization term.

*Default: 0.0001*

- **<l1\_ratio>**, *float*, is the Elastic Net mixing parameter, with  $0 \leq \text{l1\_ratio} \leq 1$ . *l1\_ratio=0* corresponds to L2 penalty, *l1\_ratio=1* to L1.

*Default: 0.15*

- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to *False*, no intercept will be used in the calculations (e.g. data is expected to be already centered).

*Default: True*

- **<n\_iter>**, *int, optional field*, specifies the number of passes over the training data (aka epochs).

*Default: 5*

- **<shuffle>**, *boolean, optional field*, specifies whether or not the training data should be shuffled after each epoch.

*Default: True*

- **<random\_state>**, *int seed, RandomState instance, or None*, represents the seed of the pseudo random number generator to use when shuffling the data for probability estimation.  
*Default: None*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: 0*
- **<epsilon>**, *int, optional field*, is the number of CPUs to use to do the OVA (One Versus All, for multi-class problems) computation. -1 means ‘all CPUs.’  
*Default: 0.1*
- **<learning\_rate>**, *string, optional field*, specifies the learning rate:
  - ‘constant:’  $\eta = \eta_0$
  - ‘optimal:’  $\eta = 1.0 / (t + t_0)$
  - ‘invscaling:’  $\eta = \eta_0 / \text{pow}(t, \text{power}_t)$  
*Default: ‘optimal’*
- **<eta0>**, *double*, specifies the initial learning rate for the ‘constant’ or ‘invscaling’ schedules. The default value is 0.0 as  $\eta_0$  is not used by the default schedule ‘optimal.’  
*Default: 0.0*
- **<power\_t>**, *double*, represents the exponent for the inverse scaling learning rate.  
*Default: 0.5*
- **<class\_weight>**, *dict, class\_label*, is the preset for the class\_weight fit parameter. Weights associated with classes. If not given, all classes are assumed to have weight one. The “auto” mode uses the values of y to automatically adjust weights inversely proportional to class frequencies.  
*Default: None*
- **<warm\_start>**, *boolean, optional field*, when set to True, the model will reuse the solution of the previous call to fit as initialization, otherwise, it will just erase the previous solution.  
*Default: False*

#### 15.3.6.1.28 Linear model fitted by minimizing a regularized empirical loss with SGD

The *Linear model fitted by minimizing a regularized empirical loss with SGD* is a model where SGD stands for Stochastic Gradient Descent: the gradient of the loss is estimated each sample at a time and the model is updated along the way with a decreasing strength schedule (aka learning rate). The regularizer is a penalty added to the loss function that shrinks model parameters towards the zero vector using either the squared euclidean norm L2 or the absolute norm L1 or a combination of both (Elastic Net). If the parameter update crosses the 0.0 value because of the regularizer,

the update is truncated to 0.0 to allow for learning sparse models and achieving online feature selection. This implementation works with data represented as dense numpy arrays of floating point values for the features. In order to use the *Linear model fitted by minimizing a regularized empirical loss with SGD*, the user needs to set the sub-node:

```
<SKLtype>linear_model|SGDRegressor</SKLtype>.
```

In addition to this XML node, several others are available:

- **<loss>**, *str*, *'squared\_loss'*, *'huber'*, *'epsilon\_insensitive'*, or *'squared\_epsilon\_insensitive'*, specifies the loss function to be used. Defaults to *'squared\_loss'* which refers to the ordinary least squares fit. *'huber'* modifies *'squared\_loss'* to focus less on getting outliers correct by switching from squared to linear loss past a distance of epsilon. *'epsilon\_insensitive'* ignores errors less than epsilon and is linear past that; this is the loss function used in SVR. *'squared\_epsilon\_insensitive'* is the same but becomes squared loss past a tolerance of epsilon.  
*Default: 'squared\_loss'*
- **<penalty>**, *str*, *'l2'* or *'l1'* or *'elasticnet'*, sets the penalty (aka regularization term) to be used. Defaults to *'l2'* which is the standard regularizer for linear SVM models. *'l1'* and *'elasticnet'* might bring sparsity to the model (feature selection) not achievable with *'l2'*.  
*Default: 'l2'*
- **<alpha>**, *float*, Constant that multiplies the regularization term. Defaults to 0.0001
- **<l1\_ratio>**, *float*, is the Elastic Net mixing parameter, with  $0 \leq l1\_ratio \leq 1$ . *l1\_ratio=0* corresponds to L2 penalty, *l1\_ratio=1* to L1.  
*Default: 0.15*
- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to *False*, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<n\_iter>**, *int, optional field*, specifies the number of passes over the training data (aka epochs).  
*Default: 5*
- **<shuffle>**, *boolean, optional field*, specifies whether or not the training data should be shuffled after each epoch.  
*Default: True*
- **<random\_state>**, *int seed, RandomState instance, or None*, sets the seed of the pseudo random number generator to use when shuffling the data.  
*Default: None*

- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: 0*
- **<epsilon>**, *float*, sets the epsilon in the epsilon-insensitive loss functions; only if loss is 'huber,' 'epsilon\_insensitive,' or 'squared\_epsilon\_insensitive.' For 'huber', determines the threshold at which it becomes less important to get the prediction exactly right. For epsilon-insensitive, any differences between the current prediction and the correct label are ignored if they are less than this threshold.  
*Default: 0.1*
- **<learning\_rate>**, *string, optional field*, Learning rate:
  - constant:  $\eta = \eta_0$
  - optimal:  $\eta = 1.0/(t+t_0)$
  - invscaling:  $\eta = \eta_0 / \text{pow}(t, \text{power}_t)$

*Default: invscaling*

- **<eta0>**, *double*, specifies the initial learning rate.  
*Default: 0.01*
- **<power\_t>**, *double, optional field*, specifies the exponent for inverse scaling learning rate.  
*Default: 0.25*
- **<warm\_start>**, *boolean, optional field*, when set to True, the model will reuse the solution of the previous call to fit as initialization, otherwise, it will just erase the previous solution.  
*Default: False*

### 15.3.6.1.29 Compute Least Angle Regression or Lasso path using LARS algorithm

The *Compute Least Angle Regression or Lasso path using LARS algorithm* is a regressor where the optimization objective for the case `method='lasso'` is:  $(1/(2 * n\_samples)) * ||y - Xw||_2^2 + \alpha * ||w||_1$  in the case of `method='lars'`, the objective function is only known in the form of an implicit equation. In order to use the *Compute Least Angle Regression or Lasso path using LARS algorithm regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|lars_path</SKLtype>.
```

In addition to this XML node, several others are available:

- **<x>**, *array, shape: (n\_samples, n\_features)*, is the input data.
- **<y>**, *array, shape: (n\_samples)*, is the set of input targets.

- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
Default: 500
- **<gram>**, *None, 'auto', array, shape: (n\_features, n\_features), optional*, Precomputed Gram matrix ( $X^T * X$ ), if 'auto', the Gram matrix is precomputed from the given X, if there are more samples than features.  
Default: None
- **<alpha\_min>**, *float, optional field*, sets the minimum correlation along the path. It corresponds to the regularization parameter alpha parameter in the Lasso.  
Default: 0
- **<method>**, *'lar', 'lasso', optional*, specifies the returned model. Select 'lar' for Least Angle Regression, 'lasso' for the Lasso.  
Default: 'lar'
- **<eps>**, *float, optional*, sets the machine-precision regularization in the computation of the Cholesky diagonal factors. Increase this for very ill-conditioned systems.  
Default: "np.finfo(np.float).eps"
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
Default: 0

### 15.3.6.1.30 Compute Lasso path with coordinate descent

The *Compute Lasso path with coordinate descent* is a regressor where the Lasso optimization function varies for mono and multi-outputs. For mono-output tasks it is:  $(1/(2 * n\_samples)) * ||y - Xw||_2^2 + alpha * ||w||_1$  For multi-output tasks it is:  $(1/(2 * n\_samples)) * ||Y - XW||_F^2 + alpha * ||W||_21$  Where:  $||W||_21 = \sum_i \sqrt{\sum_j w_{ij}^2}$  i.e. the sum of norm of each row. In order to use the *Compute Lasso path with coordinate descent regressor*, the user needs to set the sub-node:

**<SKLtype>**linear\_model|lasso\_path**</SKLtype>**.

In addition to this XML node, several others are available:

- **<X>**, *array-like, sparse matrix, shape (n\_samples, n\_features)*, represents the training data. Pass directly as Fortran-contiguous data to avoid unnecessary memory duplication. If y is mono-output then X can be sparse.
- **<y>**, *ndarray, shape = (n\_samples,) or (n\_samples, n\_outputs)*, represents the target values.
- **<eps>**, *float, optional field*, specifies the length of the path. eps=1e-3 means that alpha\_min / alpha\_max = 1e-3.
- **<n\_alphas>**, *int, optional field*, is the number of alphas along the regularization path.  
Default: 100



- **<alphas>**, *ndarray, optional field*, lists the alphas where the models will be computed. If None, alphas are set automatically.  
*Default: None*
- **<precompute>**, *True — False — ‘auto’ — array-like*, decides whether to use a precomputed Gram matrix to speed up calculations or not. If set to ‘auto,’ RAVEN will decide. The Gram matrix can also be passed as argument.  
*Default: ‘auto’*
- **<Xy>**, *array-like, optional field*,  $Xy = np.dot(X.T, y)$  that can be precomputed. It is useful only when the Gram matrix is precomputed.  
*Default: None*
- **<copy X>**, *boolean, optional field*, if True, X will be copied; else, it may be overwritten.  
*Default: True*
- **<coef\_init>**, *array, shape (n\_features, ) — None*, is the initial set of values of the coefficients.  
*Default: None*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: False*

**15.3.6.1.31 Stability path based on randomized Lasso estimates** In order to use the *Stability path based on randomized Lasso estimates regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|lasso_stability_path</SKLtype>.
```

In addition to this XML node, several others are available:

- **<X>**, *array-like, shape = [n\_samples, n\_features]*, represents the training data.
- **<y>**, *array-like, shape = [n\_samples]*, represents the target values.
- **<scaling>**, *float, optional, default=0.5*, sets the alpha parameter used to randomly scale the features. Should be between 0 and 1.  
*Default: 0.5*
- **<random\_state>**, *integer or numpy.random.RandomState, optional*, sets the generator used to randomize the design.  
*Default: None*
- **<n\_resampling>**, *int, optional*, sets the number of randomized models.  
*Default: 200*



- **<n\_grid>**, *int, optional*, specifies the number of grid points. The path is linearly reinterpolated on a grid between 0 and 1 before computing the scores.  
*Default: 100*
- **<sample\_fraction>**, *float, optional*, determines the fraction of samples to be used in each randomized design. Should be between 0 and 1. If 1, all samples are used.  
*Default: 0.75*
- **<eps>**, *float, optional field*, is the smallest value of alpha / alpha\_max considered.  
*Default: 8.8817841970012523e-16*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: False*

### 15.3.6.1.32 Gram Orthogonal Matching Pursuit (OMP)

The *Gram OMP* solves  $n\_targets$  Orthogonal Matching Pursuit problems using only the Gram matrix,  $X.T * X$ , and the product,  $X.T * y$ . In order to use the *Gram OMP regressor*, the user needs to set the sub-node:

```
<SKLtype>linear_model|orthogonal_mp_gram</SKLtype>.
```

In addition to this XML node, several others are available:

- **<Gram>**, *array, shape (n\_features, n\_features)*, contains the Gram matrix of the input data:  $X.T * X$
- **<Xy>**, *array, shape (n\_features,) or (n\_features, n\_targets)*, specifies the input targets multiplied by X:  $X.T * y$ .
- **<n\_nonzero\_coefs>**, *int*, sets the desired number of non-zero entries in the solution. If None this value is set to 10% of  $n\_features$ .  
*Default: None*
- **<tol>**, *float*, sets the maximum norm of the residual. If not None, overrides  $n\_nonzero\_coefs$ .  
*Default: None*
- **<norms\_squared>**, *array-like, shape(n\_targets)*,  
*Default: None*

### 15.3.6.2 Support Vector Machines

In machine learning, **Support Vector Machines** (SVMs, also support vector networks) are supervised learning models with associated learning algorithms that analyze data and recognize patterns,

used for classification and regression analysis. Given a set of training examples, each marked as belonging to one of two categories, an SVM training algorithm builds a model that assigns new examples into one category or the other, making it a non-probabilistic binary linear classifier. An SVM model is a representation of the examples as points in space, mapped so that the examples of the separate categories are divided by a clear gap that is as wide as possible. New examples are then mapped into that same space and predicted to belong to a category based on which side of the gap they fall on. In addition to performing linear classification, SVMs can efficiently perform a non-linear classification using what is called the kernel trick, implicitly mapping their inputs into high-dimensional feature spaces. In the following, all the SVM models available in RAVEN are reported.

### 15.3.6.2.1 Linear Support Vector Classifier

The *Linear Support Vector Classifier* is similar to SVC with parameter kernel='linear', but implemented in terms of liblinear rather than libsvm, so it has more flexibility in the choice of penalties and loss functions and should scale better (to large numbers of samples). This class supports both dense and sparse input and the multiclass support is handled according to a one-vs-the-rest scheme. In order to use the *Linear Support Vector Classifier*, the user needs to set the sub-node:

```
<SKLtype>svm|LinearSVC</SKLtype>.
```

In addition to this XML node, several others are available:

- **<C>**, *boolean, optional field*, sets the penalty parameter C of the error term.  
*Default: 1.0*
- **<loss>**, *string, 'hinge' or 'squared\_hinge'*, specifies the loss function. 'hinge' is the hinge loss (standard SVM) while 'squared\_hinge' is the squared hinge loss.  
*Default: 'squared\_hinge'*
- **<penalty>**, *string, 'l1' or 'l2'*, specifies the norm used in the penalization. The 'l2' penalty is the standard used in SVC. The 'l1' leads to coef\_vectors that are sparse.  
*Default: 'l2'*
- **<dual>**, *boolean*, selects the algorithm to either solve the dual or primal optimization problem. Prefer dual=False when n\_samples > n\_features.  
*Default: True*
- **<tol>**, *float, optional field*, specifies the tolerance for stopping criteria.  
*Default: 1e-4*
- **<multi\_class>**, *string, 'ovr' or 'crammer\_singer'*, Determines the multi-class strategy if y contains more than two classes. ovr trains n\_classes one-vs-rest classifiers, while crammer\_singer optimizes a joint objective over all classes. While crammer\_singer is interesting from a theoretical perspective as it is consistent, it is seldom used in practice and rarely leads

to better accuracy and is more expensive to compute. If `crammer_singer` is chosen, the options loss, penalty and dual will be ignored.

*Default: 'ovr'*

- **<fit\_intercept>**, *boolean, optional field*, determines whether to calculate the intercept for this model. If set to `False`, no intercept will be used in the calculations (e.g. data is expected to be already centered).  
*Default: True*
- **<intercept\_scaling>**, *float, optional field*, when `True`, the instance vector `x` becomes `[x,self.intercept_scaling]`, i.e. a “synthetic” feature with constant value equals to `intercept_scaling` is appended to the instance vector. The intercept becomes `intercept_scaling * synthetic feature weight`. **Note:** The synthetic feature weight is subject to  $l1/l2$  regularization as are all other features. To lessen the effect of regularization on the synthetic feature weight (and therefore on the intercept) `intercept_scaling` has to be increased.  
*Default: 1*
- **<class\_weight>**, *dict, 'auto', optional*, sets the parameter `C` of class `i` to `class_weight[i]*C` for SVC. If not given, all classes are assumed to have weight one. The ‘auto’ mode uses the values of `y` to automatically adjust weights inversely proportional to class frequencies.  
*Default: None*
- **<verbose>**, *boolean or integer, optional field*, use verbose mode when fitting the model.  
*Default: 0* **Note:** This setting takes advantage of a per-process runtime setting in `liblinear` that, if enabled, may not work properly in a multithreaded context.
- **<random\_state>**, *int seed, RandomState instance, or None*, represents the seed of the pseudo random number generator to use when shuffling the data for probability estimation.  
*Default: None*

### 15.3.6.2.2 C-Support Vector Classification

The *C-Support Vector Classification* is based on `libsvm`. The fit time complexity is more than quadratic with the number of samples which makes it hard to scale to datasets with more than a couple of 10000 samples. The multiclass support is handled according to a one-vs-one scheme. In order to use the *C-Support Vector Classifier*, the user needs to set the sub-node:

```
<SKLtype> svm | SVC </SKLtype>.
```

In addition to this XML node, several others are available:

- **<C>**, *boolean, optional field*, sets the penalty parameter `C` of the error term.  
*Default: 1.0*

- **<kernel>**, *string, optional*, specifies the kernel type to be used in the algorithm. It must be one of:
  - ‘linear’
  - ‘poly’
  - ‘rbf’
  - ‘sigmoid’
  - ‘precomputed’
  - a callable object

If a callable is given it is used to pre-compute the kernel matrix.

*Default: ‘rbf’*

- **<degree>**, *int, optional field*, determines the degree of the polynomial kernel function (‘poly’). Ignored by all other kernels.
 

*Default: 3.0*
- **<gamma>**, *float, optional field*, sets the kernel coefficient for the kernels ‘rbf,’ ‘poly,’ and ‘sigmoid.’ If gamma is ‘auto’ then  $1/n\_features$  will be used instead.
 

*Default: ‘auto’*
- **<coef0>**, *float, optional field*, is an independent term in kernel function. It is only significant in ‘poly’ and ‘sigmoid.’
 

*Default: 0.0*
- **<probability>**, *boolean, optional field*, determines whether or not to enable probability estimates. This must be enabled prior to calling fit, and will slow down that method.
 

*Default: False*
- **<shrinking>**, *boolean, optional field*, determines whether or not to use the shrinking heuristic.
 

*Default: True*
- **<tol>**, *float, optional field*, specifies the tolerance for stopping criteria.
 

*Default: 1e-3*
- **<cache\_size>**, *float, optional field*, specifies the size of the kernel cache (in MB).
- **<class\_weight>**, *dict, ‘auto’, optional*, sets the parameter C of class i to  $class\_weight[i]*C$  for SVC. If not given, all classes are assumed to have weight one. The ‘auto’ mode uses the values of y to automatically adjust weights inversely proportional to class frequencies.
 

*Default: None*

- **<verbose>**, *boolean, optional field*, use verbose mode when fitting the model.  
*Default: False* **Note:** This setting takes advantage of a per-process runtime setting in libsvm that, if enabled, may not work properly in a multithreaded context.
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: -1*
- **<random\_state>**, *int seed, RandomState instance, or None*, represents the seed of the pseudo random number generator to use when shuffling the data for probability estimation.  
*Default: None*

### 15.3.6.2.3 Nu-Support Vector Classification

The *Nu-Support Vector Classification* is similar to SVC but uses a parameter to control the number of support vectors. The implementation is based on libsvm. In order to use the *Nu-Support Vector Classifier*, the user needs to set the sub-node:

```
<SKLtype>svm|NuSVC</SKLtype>.
```

In addition to this XML node, several others are available:

- **<nu>**, *float, optional field*, is an upper bound on the fraction of training errors and a lower bound of the fraction of support vectors. Should be in the interval (0, 1].  
*Default: 0.5*
- **<kernel>**, *string, optional*, specifies the kernel type to be used in the algorithm. It must be one of:
  - ‘linear’
  - ‘poly’
  - ‘rbf’
  - ‘sigmoid’
  - ‘precomputed’
  - a callable object

If a callable is given it is used to pre-compute the kernel matrix.

*Default: ‘rbf’*

- **<degree>**, *int, optional field*, determines the degree of the polynomial kernel function (‘poly’). Ignored by all other kernels.  
*Default: 3*

- **<gamma>**, *float, optional field*, sets the kernel coefficient for the kernels ‘rbf,’ ‘poly,’ and ‘sigmoid.’ If gamma is ‘auto’ then 1/n\_features will be used instead.  
*Default: ‘auto’*
- **<coef0>**, *float, optional field*, is an independent term in kernel function. It is only significant in ‘poly’ and ‘sigmoid.’  
*Default: 0.0*
- **<probability>**, *boolean, optional field*, determines whether or not to enable probability estimates. This must be enabled prior to calling fit, and will slow down that method.  
*Default: False*
- **<shrinking>**, *boolean, optional field*, determines whether or not to use the shrinking heuristic.  
*Default: True*
- **<tol>**, *float, optional field*, specifies the tolerance for stopping criteria.  
*Default: 1e-3*
- **<cache\_size>**, *float, optional field*, specifies the size of the kernel cache (in MB).
- **<verbose>**, *boolean, optional field*, use verbose mode when fitting the model.  
*Default: False* **Note:** This setting takes advantage of a per-process runtime setting in libsvm that, if enabled, may not work properly in a multithreaded context.
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: -1*
- **<random\_state>**, *int seed, RandomState instance, or None*, represents the seed of the pseudo random number generator to use when shuffling the data for probability estimation.  
*Default: None*

#### 15.3.6.2.4 Support Vector Regression

The *Support Vector Regression* is an epsilon-Support Vector Regression. The free parameters in this model are C and epsilon. The implementation is based on libsvm. In order to use the *Support Vector Regressor*, the user needs to set the sub-node:

```
<SKLtype> svm | SVR</SKLtype>.
```

In addition to this XML node, several others are available:

- **<C>**, *boolean, optional field*, sets the penalty parameter C of the error term.  
*Default: 1.0*

- **<epsilon>**, *float, optional field*, specifies the epsilon-tube within which no penalty is associated in the training loss function with points predicted within a distance epsilon from the actual value.  
*Default: 0.1*
- **<kernel>**, *string, optional*, specifies the kernel type to be used in the algorithm. It must be one of:
  - ‘linear’
  - ‘poly’
  - ‘rbf’
  - ‘sigmoid’
  - ‘precomputed’
  - a callable object

If a callable is given it is used to pre-compute the kernel matrix.

*Default: ‘rbf’*

- **<degree>**, *int, optional field*, determines the degree of the polynomial kernel function (‘poly’). Ignored by all other kernels.  
*Default: 3.0*
- **<gamma>**, *float, optional field*, sets the kernel coefficient for the kernels ‘rbf,’ ‘poly,’ and ‘sigmoid.’ If gamma is ‘auto’ then  $1/n_{\text{features}}$  will be used instead.  
*Default: ‘auto’*
- **<coef0>**, *float, optional field*, is an independent term in kernel function. It is only significant in ‘poly’ and ‘sigmoid.’  
*Default: 0.0*
- **<shrinking>**, *boolean, optional field*, determines whether or not to use the shrinking heuristic.  
*Default: True*
- **<tol>**, *float, optional field*, specifies the tolerance for stopping criteria.  
*Default: 1e-3*
- **<cache\_size>**, *float, optional field*, specifies the size of the kernel cache (in MB).
- **<verbose>**, *boolean, optional field*, use verbose mode when fitting the model.  
*Default: False* **Note:** This setting takes advantage of a per-process runtime setting in libsvm that, if enabled, may not work properly in a multithreaded context.
- **<max\_iter>**, *integer, optional field*, specifies the maximum number of iterations.  
*Default: -1*

### 15.3.6.3 Multi Class

Multiclass classification means a classification task with more than two classes; e.g., classify a set of images of fruits which may be oranges, apples, or pears. Multiclass classification makes the assumption that each sample is assigned to one and only one label: a fruit can be either an apple or a pear but not both at the same time. In the following, all the multi-class models available in RAVEN are reported.

#### 15.3.6.3.1 One-vs-the-rest (OvR) multiclass/multilabel strategy

The *One-vs-the-rest (OvR) multiclass/multilabel strategy*, also known as one-vs-all, consists in fitting one classifier per class. For each classifier, the class is fitted against all the other classes. In addition to its computational efficiency (only  $n_{\text{classes}}$  classifiers are needed), one advantage of this approach is its interpretability. Since each class is represented by one and one classifier only, it is possible to gain knowledge about the class by inspecting its corresponding classifier. This is the most commonly used strategy and is a fair default choice.

In order to use the *One-vs-the-rest (OvR) multiclass/multilabel classifier*, the user needs to set the sub-node:

```
<SKLtype>multiClass|OneVsRestClassifier</SKLtype>.
```

In addition to this XML node, several others are available:

- **<estimator>**, *boolean, required field*, An estimator object implementing fit and one of decision\_function or predict\_proba. This XML node needs to contain the following attribute:
  - **estimatorType**, *required string attribute*, this attribute is another reduced order mode type that needs to be used for the construction of the multi-class algorithms. Each sub-sequential node depends on the chosen ROM.

#### 15.3.6.3.2 One-vs-one multiclass strategy

The *One-vs-one multiclass strategy* consists in fitting one classifier per class pair. At prediction time, the class which received the most votes is selected. Since it requires to fit  $n_{\text{classes}} * (n_{\text{classes}} - 1) / 2$  classifiers, this method is usually slower than one-vs-the-rest, due to its  $O(n_{\text{classes}}^2)$  complexity. However, this method may be advantageous for algorithms such as kernel algorithms which do not scale well with  $n_{\text{samples}}$ . This is because each individual learning problem only involves a small subset of the data whereas, with one-vs-the-rest, the complete dataset is used  $n_{\text{classes}}$  times.

In order to use the *One-vs-one multiclass classifier*, the user needs to set the sub-node:



`<SKLtype>multiClass|OneVsOneClassifier</SKLtype>`.

In addition to this XML node, several others are available:

- **<estimator>**, *boolean, required field*, An estimator object implementing fit and one of decision\_function or predict\_proba. This XML node needs to contain the following attribute:
  - **estimatorType**, *required string attribute*, this attribute is another reduced order mode type that needs to be used for the construction of the multi-class algorithms. Each sub-sequential node depends on the chosen ROM.

#### 15.3.6.3.3 Error-Correcting Output-Code multiclass strategy

The *Error-Correcting Output-Code multiclass strategy* consists in representing each class with a binary code (an array of 0s and 1s). At fitting time, one binary classifier per bit in the code book is fitted. At prediction time, the classifiers are used to project new points in the class space and the class closest to the points is chosen. The main advantage of these strategies is that the number of classifiers used can be controlled by the user, either for compressing the model ( $0 < code\_size < 1$ ) or for making the model more robust to errors ( $code\_size > 1$ ).

In order to use the *Error-Correcting Output-Code multiclass classifier*, the user needs to set the sub-node:

`<SKLtype>multiClass|OutputCodeClassifier</SKLtype>`.

In addition to this XML node, several others are available:

- **<estimator>**, *boolean, required field*, An estimator object implementing fit and one of decision\_function or predict\_proba. This XML node needs to contain the following attribute:
  - **estimatorType**, *required string attribute*, this attribute is another reduced order mode type that needs to be used for the construction of the multi-class algorithms. Each sub-sequential node depends on the chosen ROM.
- **<code\_size>**, *float, required field*, represents the percentage of the number of classes to be used to create the code book. A number between 0 and 1 will require fewer classifiers than one-vs-the-rest. A number greater than 1 will require more classifiers than one-vs-the-rest.

#### 15.3.6.4 Naive Bayes

Naive Bayes methods are a set of supervised learning algorithms based on applying Bayes' theorem with the "naive" assumption of independence between every pair of features. Given a class variable

$y$  and a dependent feature vector  $x_1$  through  $x_n$ , Bayes' theorem states the following relationship:

$$P(y | x_1, \dots, x_n) = \frac{P(y)P(x_1, \dots, x_n | y)}{P(x_1, \dots, x_n)} \quad (3)$$

Using the naive independence assumption that

$$P(x_i | y, x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = P(x_i | y), \quad (4)$$

for all  $i$ , this relationship is simplified to

$$P(y | x_1, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i | y)}{P(x_1, \dots, x_n)} \quad (5)$$

Since  $P(x_1, \dots, x_n)$  is constant given the input, we can use the following classification rule:

$$P(y | x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i | y) \quad (6)$$

$$\hat{y} = \arg \max_y P(y) \prod_{i=1}^n P(x_i | y), \quad (7)$$

and we can use Maximum A Posteriori (MAP) estimation to estimate  $P(y)$  and  $P(x_i | y)$ ; the former is then the relative frequency of class  $y$  in the training set. The different naive Bayes classifiers differ mainly by the assumptions they make regarding the distribution of  $P(x_i | y)$ .

In spite of their apparently over-simplified assumptions, naive Bayes classifiers have worked quite well in many real-world situations, famously document classification and spam filtering. They require a small amount of training data to estimate the necessary parameters. (For theoretical reasons why naive Bayes works well, and on which types of data it does, see the references below.) Naive Bayes learners and classifiers can be extremely fast compared to more sophisticated methods. The decoupling of the class conditional feature distributions means that each distribution can be independently estimated as a one dimensional distribution. This in turn helps to alleviate problems stemming from the curse of dimensionality.

On the flip side, although naive Bayes is known as a decent classifier, it is known to be a bad estimator, so the probability outputs from `predict_proba` are not to be taken too seriously. In the following, all the Naive Bayes available in RAVEN are reported.

#### 15.3.6.4.1 Gaussian Naive Bayes

The *Gaussian Naive Bayes strategy* implements the Gaussian Naive Bayes algorithm for classification. The likelihood of the features is assumed to be Gaussian:

$$P(x_i | y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}\right) \quad (8)$$

The parameters  $\sigma_y$  and  $\mu_y$  are estimated using maximum likelihood.

In order to use the *Gaussian Naive Bayes strategy*, the user needs to set the sub-node:

```
<SKLtype>naiveBayes|GaussianNB</SKLtype>.
```

There are no additional sub-nodes available for this method.

#### 15.3.6.4.2 Multinomial Naive Bayes

The *Multinomial Naive Bayes* implements the naive Bayes algorithm for multinomially distributed data, and is one of the two classic naive Bayes variants used in text classification (where the data is typically represented as word vector counts, although tf-idf vectors are also known to work well in practice). The distribution is parametrized by vectors  $\theta_y = (\theta_{y1}, \dots, \theta_{yn})$  for each class  $y$ , where  $n$  is the number of features (in text classification, the size of the vocabulary) and  $\theta_{yi}$  is the probability  $P(x_i | y)$  of feature  $i$  appearing in a sample belonging to class  $y$ . The parameters  $\theta_y$  are estimated by a smoothed version of maximum likelihood, i.e. relative frequency counting:

$$\hat{\theta}_{yi} = \frac{N_{yi} + \alpha}{N_y + \alpha n} \quad (9)$$

where  $N_{yi} = \sum_{x \in T} x_i$  is the number of times feature  $i$  appears in a sample of class  $y$  in the training set  $T$ , and  $N_y = \sum_{i=1}^{|T|} N_{yi}$  is the total count of all features for class  $y$ . The smoothing priors  $\alpha \geq 0$  account for features not present in the learning samples and prevents zero probabilities in further computations. Setting  $\alpha = 1$  is called Laplace smoothing, while  $\alpha < 1$  is called Lidstone smoothing. In order to use the *Multinomial Naive Bayes strategy*, the user needs to set the sub-node:

```
<SKLtype>naiveBayes|MultinomialNB</SKLtype>.
```

In addition to this XML node, several others are available:

- **<alpha>**, *float, optional field*, specifies an additive (Laplace/Lidstone) smoothing parameter (0 for no smoothing).  
*Default: 1.0*
- **<fit\_prior>**, *boolean, required field*, determines whether to learn class prior probabilities or not. If false, a uniform prior will be used.  
*Default: True*
- **<class\_prior>**, *array-like float (n\_classes), optional field*, specifies prior probabilities of the classes. If specified, the priors are not adjusted according to the data.  
*Default: None*

#### 15.3.6.4.3 Bernoulli Naive Bayes

The *Bernoulli Naive Bayes* implements the naive Bayes training and classification algorithms for data that is distributed according to multivariate Bernoulli distributions; i.e., there may be multiple features but each one is assumed to be a binary-valued (Bernoulli, boolean) variable. Therefore, this class requires samples to be represented as binary-valued feature vectors; if handed any other kind of data, a *Bernoulli Naive Bayes* instance may binarize its input (depending on the binarize parameter). The decision rule for Bernoulli naive Bayes is based on

$$P(x_i | y) = P(i | y)x_i + (1 - P(i | y))(1 - x_i) \quad (10)$$

which differs from multinomial NB's rule in that it explicitly penalizes the non-occurrence of a feature  $i$  that is an indicator for class  $y$ , where the multinomial variant would simply ignore a non-occurring feature. In the case of text classification, word occurrence vectors (rather than word count vectors) may be used to train and use this classifier. *Bernoulli Naive Bayes* might perform better on some datasets, especially those with shorter documents. It is advisable to evaluate both models, if time permits. In order to use the *Bernoulli Naive Bayes strategy*, the user needs to set the sub-node:

```
<SKLtype>naiveBayes|BernoulliNB</SKLtype>.
```

In addition to this XML node, several others are available:

- **<alpha>**, *float, optional field*, specifies an additive (Laplace/Lidstone) smoothing parameter (0 for no smoothing).  
*Default: 1.0*
- **<binarize>**, *float, optional field*, Threshold for binarizing (mapping to booleans) of sample features. If None, input is presumed to already consist of binary vectors.  
*Default: 0.0*
- **<fit\_prior>**, *boolean, required field*, determines whether to learn class prior probabilities or not. If false, a uniform prior will be used.  
*Default: True*
- **<class\_prior>**, *array-like float (n\_classes), optional field*, specifies prior probabilities of the classes. If specified, the priors are not adjusted according to the data.  
*Default: None*

#### 15.3.6.5 Neighbors

The *Neighbors* class provides functionality for unsupervised and supervised neighbor-based learning methods. The unsupervised nearest neighbors method is the foundation of many other learning

methods, notably manifold learning and spectral clustering. Supervised neighbors-based learning comes in two flavors: classification for data with discrete labels, and regression for data with continuous labels.

The principle behind nearest neighbor methods is to find a predefined number of training samples closest in distance to the new point, and predict the label from these. The number of samples can be a user-defined constant (k-nearest neighbor learning), or vary based on the local density of points (radius-based neighbor learning). The distance can, in general, be any metric measure: standard Euclidean distance is the most common choice. Neighbor-based methods are known as non-generalizing machine learning methods, since they simply “remember” all of its training data (possibly transformed into a fast indexing structure such as a Ball Tree or KD Tree.).

In the following, all the Neighbors’ models available in RAVEN are reported.

### 15.3.6.5.1 Nearest Neighbors

The *Nearest Neighbors* implements unsupervised nearest neighbors learning. It acts as a uniform interface to three different nearest neighbors algorithms: BallTree, KDTree, and a brute-force algorithm. In order to use the *Nearest Neighbors strategy*, the user needs to set the sub-node:

```
<SKLtype>neighbors|NearestNeighbors</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_neighbors>**, *integer, optional field*, specifies the number of neighbors to use by default for ‘k\_neighbors’ queries.  
*Default: 5*
- **<radius>**, *float, optional field*, specifies the range of parameter space to use by default for ‘radius\_neighbors’ queries.  
*Default: 1.0*
- **<algorithm>**, *string, optional field*, specifies the algorithm used to compute the nearest neighbors:
  - *ball\_tree* will use BallTree.
  - *kd\_tree* will use KDtree.
  - *brute* will use a brute-force search.
  - *auto* will attempt to decide the most appropriate algorithm based on the values passed to fit method.

**Note:** Fitting on sparse input will override the setting of this parameter, using brute force.  
*Default: auto*

- **<leaf\_size>**, *integer, optional field*, sets the leaf size passed to the BallTree or KDTree. This can affect the speed of the construction and query, as well as the memory required to store the tree. The optimal value depends on the nature of the problem.  
*Default: 30*
- **<p>**, *integer, optional field*, is a parameter for the Minkowski metric. When  $p = 1$ , this is equivalent to using manhattan distance (L1), and euclidean distance (L2) for  $p = 2$ . For arbitrary  $p$ , minkowski distance (L<sub>p</sub>) is used.  
*Default: 2*

#### 15.3.6.5.2 K Neighbors Classifier

The *K Neighbors Classifier* is a type of instance-based learning or non-generalizing learning: it does not attempt to construct a general internal model, but simply stores instances of the training data. Classification is computed from a simple majority vote of the nearest neighbors of each point: a query point is assigned the data class which has the most representatives within the nearest neighbors of the point. It implements learning based on the  $k$  nearest neighbors of each query point, where  $k$  is an integer value specified by the user.

In order to use the *K Neighbors Classifier*, the user needs to set the sub-node:

```
<SKLtype>neighbors|KNeighborsClassifier</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_neighbors>**, *integer, optional field*, specifies the number of neighbors to use by default for 'k\_neighbors' queries.  
*Default: 5*
- **<weights>**, *string, optional field*, specifies the weight function used in prediction. Possible values:
  - *uniform* : uniform weights. All points in each neighborhood are weighted equally;
  - *distance* : weight points by the inverse of their distance. In this case, closer neighbors of a query point will have a greater influence than neighbors which are further away.

*Default: uniform*

- **<algorithm>**, *string, optional field*, specifies the algorithm used to compute the nearest neighbors:
  - *ball\_tree* will use BallTree.
  - *kd\_tree* will use KDtree.

- *brute* will use a brute-force search.
- *auto* will attempt to decide the most appropriate algorithm based on the values passed to fit method.

**Note:** Fitting on sparse input will override the setting of this parameter, using brute force.

*Default: auto*

- **<leaf\_size>**, *integer, optional field*, sets the leaf size passed to the BallTree or KDTree. This can affect the speed of the construction and query, as well as the memory required to store the tree. The optimal value depends on the nature of the problem.

*Default: 30*

- **<metric>**, *string, optional field*, sets the distance metric to use for the tree. The Minkowski metric with  $p=2$  is equivalent to the standard Euclidean metric.

*Default: minkowski*

- **<p>**, *integer, optional field*, is a parameter for the Minkowski metric. When  $p = 1$ , this is equivalent to using manhattan distance (L1), and euclidean distance (L2) for  $p = 2$ . For arbitrary  $p$ , minkowski distance (L<sub>p</sub>) is used.

*Default: 2*

### 15.3.6.5.3 Radius Neighbors Classifier

The *Radius Neighbors Classifier* is a type of instance-based learning or non-generalizing learning: it does not attempt to construct a general internal model, but simply stores instances of the training data. Classification is computed from a simple majority vote of the nearest neighbors of each point: a query point is assigned the data class which has the most representatives within the nearest neighbors of the point. It implements learning based on the number of neighbors within a fixed radius  $r$  of each training point, where  $r$  is a floating-point value specified by the user.

In order to use the *Radius Neighbors Classifier*, the user needs to set the sub-node:

**<SKLtype>**neighbors|RadiusNeighbors**</SKLtype>**.

In addition to this XML node, several others are available:

- **<radius>**, *float, optional field*, specifies the range of parameter space to use by default for ‘radius\_neighbors’ queries.

*Default: 1.0*

- **<weights>**, *string, optional field*, specifies the weight function used in prediction. Possible values:

- *uniform* : uniform weights. All points in each neighborhood are weighted equally;

- *distance* : weight points by the inverse of their distance. In this case, closer neighbors of a query point will have a greater influence than neighbors which are further away.

*Default: uniform*

- **<algorithm>**, *string, optional field*, specifies the algorithm used to compute the nearest neighbors:
  - *ball\_tree* will use BallTree.
  - *kd\_tree* will use KDtree.
  - *brute* will use a brute-force search.
  - *auto* will attempt to decide the most appropriate algorithm based on the values passed to fit method.

**Note:** Fitting on sparse input will override the setting of this parameter, using brute force.

*Default: auto*

- **<leaf\_size>**, *integer, optional field*, sets the leaf size passed to the BallTree or KDTree. This can affect the speed of the construction and query, as well as the memory required to store the tree. The optimal value depends on the nature of the problem.

*Default: 30*

- **<metric>**, *string, optional field*, sets the distance metric to use for the tree. The Minkowski metric with  $p=2$  is equivalent to the standard Euclidean metric.

*Default: minkowski*

- **<p>**, *integer, optional field*, is a parameter for the Minkowski metric. When  $p = 1$ , this is equivalent to using manhattan distance (L1), and euclidean distance (L2) for  $p = 2$ . For arbitrary  $p$ , minkowski distance (L<sub>p</sub>) is used.

*Default: 2*

- **<outlier\_label>**, *integer, optional field*, is a label, which is given for outlier samples (samples with no neighbors on a given radius). If set to None, ValueError is raised, when an outlier is detected.

*Default: None*

#### 15.3.6.5.4 K Neighbors Regressor

The *K Neighbors Regressor* can be used in cases where the data labels are continuous rather than discrete variables. The label assigned to a query point is computed based on the mean of the labels of its nearest neighbors. It implements learning based on the  $k$  nearest neighbors of each query point, where  $k$  is an integer value specified by the user.



In order to use the *K Neighbors Regressor*, the user needs to set the sub-node:

`<SKLtype>neighbors|KNeighborsRegressor</SKLtype>`.

In addition to this XML node, several others are available:

- **<n\_neighbors>**, *integer, optional field*, specifies the number of neighbors to use by default for 'k\_neighbors' queries.

*Default: 5*

- **<weights>**, *string, optional field*, specifies the weight function used in prediction. Possible values:
  - *uniform* : uniform weights. All points in each neighborhood are weighted equally;
  - *distance* : weight points by the inverse of their distance. In this case, closer neighbors of a query point will have a greater influence than neighbors which are further away.

*Default: uniform*

- **<algorithm>**, *string, optional field*, specifies the algorithm used to compute the nearest neighbors:

- *ball\_tree* will use BallTree.
- *kd\_tree* will use KDtree.
- *brute* will use a brute-force search.
- *auto* will attempt to decide the most appropriate algorithm based on the values passed to fit method.

**Note:** Fitting on sparse input will override the setting of this parameter, using brute force.

*Default: auto*

- **<leaf\_size>**, *integer, optional field*, sets the leaf size passed to the BallTree or KDTree. This can affect the speed of the construction and query, as well as the memory required to store the tree. The optimal value depends on the nature of the problem.

*Default: 30*

- **<metric>**, *string, optional field*, sets the distance metric to use for the tree. The Minkowski metric with  $p=2$  is equivalent to the standard Euclidean metric.

*Default: minkowski*

- **<p>**, *integer, optional field*, is a parameter for the Minkowski metric. When  $p = 1$ , this is equivalent to using manhattan distance (L1), and euclidean distance (L2) for  $p = 2$ . For arbitrary  $p$ , minkowski distance (L<sub>p</sub>) is used.

*Default: 2*

### 15.3.6.5.5 Radius Neighbors Regressor

The *Radius Neighbors Regressor* can be used in cases where the data labels are continuous rather than discrete variables. The label assigned to a query point is computed based on the mean of the labels of its nearest neighbors. It implements learning based on the neighbors within a fixed radius  $r$  of the query point, where  $r$  is a floating-point value specified by the user.

In order to use the *Radius Neighbors Regressor*, the user needs to set the sub-node:

```
<SKLtype>neighbors|RadiusNeighborsRegressor</SKLtype>.
```

In addition to this XML node, several others are available:

- **<radius>**, *float, optional field*, specifies the range of parameter space to use by default for 'radius\_neighbors' queries.  
*Default: 1.0*
- **<weights>**, *string, optional field*, specifies the weight function used in prediction. Possible values:
  - *uniform* : uniform weights. All points in each neighborhood are weighted equally;
  - *distance* : weight points by the inverse of their distance. In this case, closer neighbors of a query point will have a greater influence than neighbors which are further away.

*Default: uniform*

- **<algorithm>**, *string, optional field*, specifies the algorithm used to compute the nearest neighbors:
  - *ball\_tree* will use BallTree.
  - *kd\_tree* will use KDtree.
  - *brute* will use a brute-force search.
  - *auto* will attempt to decide the most appropriate algorithm based on the values passed to fit method.

**Note:** Fitting on sparse input will override the setting of this parameter, using brute force.

*Default: auto*

- **<leaf\_size>**, *integer, optional field*, sets the leaf size passed to the BallTree or KDTree. This can affect the speed of the construction and query, as well as the memory required to store the tree. The optimal value depends on the nature of the problem.  
*Default: 30*

- **<metric>**, *string, optional field*, sets the distance metric to use for the tree. The Minkowski metric with  $p=2$  is equivalent to the standard Euclidean metric.  
*Default: minkowski*
- **<p>**, *integer, optional field*, is a parameter for the Minkowski metric. When  $p = 1$ , this is equivalent to using Manhattan distance (L1), and Euclidean distance (L2) for  $p = 2$ . For arbitrary  $p$ , Minkowski distance (L<sub>p</sub>) is used.  
*Default: 2*

#### 15.3.6.5.6 Nearest Centroid Classifier

The *Nearest Centroid classifier* is a simple algorithm that represents each class by the centroid of its members. It also has no parameters to choose, making it a good baseline classifier. It does, however, suffer on non-convex classes, as well as when classes have drastically different variances, as equal variance in all dimensions is assumed.

In order to use the *Nearest Centroid Classifier*, the user needs to set the sub-node:

```
<SKLtype>neighbors|NearestCentroid</SKLtype>.
```

In addition to this XML node, several others are available:

- **<shrink\_threshold>**, *float, optional field*, defines the threshold for shrinking centroids to remove features.  
*Default: None*

#### 15.3.6.6 Quadratic Discriminant Analysis

The *Quadratic Discriminant Analysis* is a classifier with a quadratic decision boundary, generated by fitting class conditional densities to the data and using Bayes' rule. The model fits a Gaussian density to each class.

In order to use the *Quadratic Discriminant Analysis Classifier*, the user needs to set the sub-node:

```
<SKLtype>qda|QDA</SKLtype>.
```

In addition to this XML node, several others are available:

- **<priors>**, *array-like (n\_classes), optional field*, specifies the priors on the classes.  
*Default: None*

- **<reg\_param>**, *float, optional field*, regularizes the covariance estimate as  $(1 - \text{reg\_param}) * \text{Sigma} + \text{reg\_param} * \text{Identity}(n\_features)$ .  
*Default: 0.0*

### 15.3.6.7 Tree

Decision Trees (DTs) are a non-parametric supervised learning method used for classification and regression. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features.

- Some advantages of decision trees are:
- Simple to understand and to interpret. Trees can be visualized.
- Requires little data preparation. Other techniques often require data normalization, dummy variables need to be created and blank values to be removed. Note however, that this module does not support missing values.
- The cost of using the tree (i.e., predicting data) is logarithmic in the number of data points used to train the tree.
- Able to handle both numerical and categorical data. Other techniques are usually specialized in analyzing datasets that have only one type of variable.
- Able to handle multi-output problems.
- Uses a white box model. If a given situation is observable in a model, the explanation for the condition is easily explained by boolean logic. By contrast, in a black box model (e.g., in an artificial neural network), results may be more difficult to interpret.
- Possible to validate a model using statistical tests. That makes it possible to account for the reliability of the model.
- Performs well even if its assumptions are somewhat violated by the true model from which the data were generated.

The disadvantages of decision trees include:

- Decision-tree learners can create over-complex trees that do not generalise the data well. This is called overfitting. Mechanisms such as pruning (not currently supported), setting the minimum number of samples required at a leaf node or setting the maximum depth of the tree are necessary to avoid this problem.

- Decision trees can be unstable because small variations in the data might result in a completely different tree being generated. This problem is mitigated by using decision trees within an ensemble.
- The problem of learning an optimal decision tree is known to be NP-complete under several aspects of optimality and even for simple concepts. Consequently, practical decision-tree learning algorithms are based on heuristic algorithms such as the greedy algorithm where locally optimal decisions are made at each node. Such algorithms cannot guarantee to return the globally optimal decision tree. This can be mitigated by training multiple trees in an ensemble learner, where the features and samples are randomly sampled with replacement.
- There are concepts that are hard to learn because decision trees do not express them easily, such as XOR, parity or multiplexer problems.
- Decision tree learners create biased trees if some classes dominate. It is therefore recommended to balance the dataset prior to fitting with the decision tree.

In the following, all the linear models available in RAVEN are reported.

#### 15.3.6.7.1 Decision Tree Classifier

The *Decision Tree Classifier* is a classifier that is based on the decision tree logic.

In order to use the *Decision Tree Classifier*, the user needs to set the sub-node:

```
<SKLtype>tree|DecisionTreeClassifier</SKLtype>.
```

In addition to this XML node, several others are available:

- **<criterion>**, *string, optional field*, specifies the function used to measure the quality of a split. Supported criteria are “gini” for the Gini impurity and “entropy” for the information gain.  
*Default: gini*
- **<splitter>**, *string, optional field*, specifies the strategy used to choose the split at each node. Supported strategies are “best” to choose the best split and “random” to choose the best random split.  
*Default: best*
- **<max\_features>**, *int, float or string, optional field*, sets the number of features to consider when looking for the best split:
  - If int, then consider max\_features features at each split.
  - If float, then max\_features is a percentage and int(max\_features \* n\_features) features are considered at each split.

- If “auto,” then `max_features=sqrt(n_features)`.
- If “sqrt,” then `max_features=sqrt(n_features)`.
- If “log2,” then `max_features=log2(n_features)`.
- If None, then `max_features=n_features`.

**Note:** The search for a split does not stop until at least one valid partition of the node samples is found, even if it requires to effectively inspect more than `max_features` features.

*Default: None*

- **<max\_depth>**, *integer, optional field*, determines the maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than `min_samples_split` samples. Ignored if `max_samples_leaf` is not None.

*Default: None*

- **<min\_samples\_split>**, *integer, optional field*, sets the minimum number of samples required to split an internal node.

*Default: 2*

- **<min\_samples\_leaf>**, *integer, optional field*, sets the minimum number of samples required to be at a leaf node.

*Default: 1*

- **<max\_leaf\_nodes>**, *integer, optional field*, grow a tree with `max_leaf_nodes` in best-first fashion. Best nodes are defined by relative reduction in impurity. If None then unlimited number of leaf nodes. If not None then `max_depth` will be ignored.

*Default: None*

### 15.3.6.7.2 Decision Tree Regressor

The *Decision Tree Regressor* is a Regressor that is based on the decision tree logic. In order to use the *Decision Tree Regressor*, the user needs to set the sub-node:

```
<SKLtype>tree|DecisionTreeRegressor</SKLtype>.
```

In addition to this XML node, several others are available:

- **<criterion>**, *string, optional field*, specifies the function used to measure the quality of a split. The only supported criterion is “mse” for mean squared error.

*Default: mse*

- **<splitter>**, *string, optional field*, specifies the strategy used to choose the split at each node. Supported strategies are “best” to choose the best split and “random” to choose the best random split.

*Default: best*

- **<max\_features>**, *int, float or string, optional field*, sets the number of features to consider when looking for the best split:
  - If int, then consider max\_features features at each split.
  - If float, then max\_features is a percentage and int(max\_features \* n\_features) features are considered at each split.
  - If “auto,” then max\_features=sqrt(n\_features).
  - If “sqrt,” then max\_features=sqrt(n\_features).
  - If “log2,” then max\_features=log2(n\_features).
  - If None, then max\_features=n\_features.

**Note:** The search for a split does not stop until at least one valid partition of the node samples is found, even if it requires to effectively inspect more than max\_features features.

*Default: None*

- **<max\_depth>**, *integer, optional field*, determines the maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min\_samples\_split samples. Ignored if max\_samples\_leaf is not None.

*Default: None*

- **<min\_samples\_split>**, *integer, optional field*, sets the minimum number of samples required to split an internal node.

*Default: 2*

- **<min\_samples\_leaf>**, *integer, optional field*, sets the minimum number of samples required to be at a leaf node.

*Default: 1*

- **<max\_leaf\_nodes>**, *integer, optional field*, grow a tree with max\_leaf\_nodes in best-first fashion. Best nodes are defined by relative reduction in impurity. If None then unlimited number of leaf nodes. If not None then max\_depth will be ignored.

*Default: None*

### 15.3.6.7.3 Extra Tree Classifier

The *Extra Tree Classifier* is an extremely randomized tree classifier. Extra-trees differ from classic decision trees in the way they are built. When looking for the best split to separate the samples of a node into two groups, random splits are drawn for each of the max\_features randomly selected features and the best split among those is chosen. When max\_features is set 1, this amounts to building a totally random decision tree.

In order to use the *Extra Tree Classifier*, the user needs to set the sub-node:

```
<SKLtype>tree|ExtraTreeClassifier</SKLtype>.
```

In addition to this XML node, several others are available:

- **<criterion>**, *string, optional field*, specifies the function used to measure the quality of a split. Supported criteria are “gini” for the Gini impurity and “entropy” for the information gain.  
*Default: gini*
- **<splitter>**, *string, optional field*, specifies the strategy used to choose the split at each node. Supported strategies are “best” to choose the best split and “random” to choose the best random split.  
*Default: random*
- **<max\_features>**, *int, float or string, optional field*, sets the number of features to consider when looking for the best split:
  - If int, then consider max\_features features at each split.
  - If float, then max\_features is a percentage and  $\text{int}(\text{max\_features} * \text{n\_features})$  features are considered at each split.
  - If “auto,” then  $\text{max\_features} = \text{sqrt}(\text{n\_features})$ .
  - If “sqrt,” then  $\text{max\_features} = \text{sqrt}(\text{n\_features})$ .
  - If “log2,” then  $\text{max\_features} = \text{log}_2(\text{n\_features})$ .
  - If None, then  $\text{max\_features} = \text{n\_features}$ .

**Note:** The search for a split does not stop until at least one valid partition of the node samples is found, even if it requires to effectively inspect more than max\_features features.

*Default: auto*

- **<max\_depth>**, *integer, optional field*, determines the maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min\_samples\_split samples. Ignored if max\_samples\_leaf is not None.  
*Default: None*
- **<min\_samples\_split>**, *integer, optional field*, sets the minimum number of samples required to split an internal node.  
*Default: 2*
- **<min\_samples\_leaf>**, *integer, optional field*, sets the minimum number of samples required to be at a leaf node.  
*Default: 1*
- **<max\_leaf\_nodes>**, *integer, optional field*, grow a tree with max\_leaf\_nodes in best-first fashion. Best nodes are defined by relative reduction in impurity. If None then unlimited number of leaf nodes. If not None then max\_depth will be ignored.  
*Default: None*



#### 15.3.6.7.4 Extra Tree Regressor

The *Extra Tree Regressor* is an extremely randomized tree regressor. Extra-trees differ from classic decision trees in the way they are built. When looking for the best split to separate the samples of a node into two groups, random splits are drawn for each of the `max_features` randomly selected features and the best split among those is chosen. When `max_features` is set 1, this amounts to building a totally random decision tree.

In order to use the *Extra Tree Regressor*, the user needs to set the sub-node:

```
<SKLtype>tree|ExtraTreeRegressor</SKLtype>.
```

In addition to this XML node, several others are available:

- **<criterion>**, *string, optional field*, specifies the function used to measure the quality of a split. The only supported criterion is “mse” for mean squared error.  
*Default: mse*
- **<splitter>**, *string, optional field*, specifies the strategy used to choose the split at each node. Supported strategies are “best” to choose the best split and “random” to choose the best random split.  
*Default: random*
- **<max\_features>**, *int, float or string, optional field*, sets the number of features to consider when looking for the best split:
  - If int, then consider `max_features` features at each split.
  - If float, then `max_features` is a percentage and `int(max_features * n_features)` features are considered at each split.
  - If “auto,” then `max_features=sqrt(n_features)`.
  - If “sqrt,” then `max_features=sqrt(n_features)`.
  - If “log2,” then `max_features=log2(n_features)`.
  - If None, then `max_features=n_features`.

**Note:** The search for a split does not stop until at least one valid partition of the node samples is found, even if it requires to effectively inspect more than `max_features` features.

*Default: auto*

- **<max\_depth>**, *integer, optional field*, determines the maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than `min_samples_split` samples. Ignored if `max_samples_leaf` is not None.  
*Default: None*

- **<min\_samples\_split>**, *integer, optional field*, sets the minimum number of samples required to split an internal node.  
*Default: 2*
- **<min\_samples\_leaf>**, *integer, optional field*, sets the minimum number of samples required to be at a leaf node.  
*Default: 1*
- **<max\_leaf\_nodes>**, *integer, optional field*, grow a tree with max\_leaf\_nodes in best-first fashion. Best nodes are defined by relative reduction in impurity. If None then unlimited number of leaf nodes. If not None then max\_depth will be ignored.  
*Default: None*

### 15.3.6.8 Gaussian Process

Gaussian Processes for Machine Learning (GPML) is a generic supervised learning method primarily designed to solve regression problems. The advantages of Gaussian Processes for Machine Learning are:

- The prediction interpolates the observations (at least for regular correlation models).
- The prediction is probabilistic (Gaussian) so that one can compute empirical confidence intervals and exceedance probabilities that might be used to refit (online fitting, adaptive fitting) the prediction in some region of interest.
- Versatile: different linear regression models and correlation models can be specified. Common models are provided, but it is also possible to specify custom models provided they are stationary.

The disadvantages of Gaussian Processes for Machine Learning include:

- It is not sparse. It uses the whole samples/features information to perform the prediction.
- It loses efficiency in high dimensional spaces – namely when the number of features exceeds a few dozens. It might indeed give poor performance and it loses computational efficiency.
- Classification is only a post-processing, meaning that one first needs to solve a regression problem by providing the complete scalar float precision output  $y$  of the experiment one is attempting to model.

In order to use the *Gaussian Process Regressor*, the user needs to set the sub-node:

`<SKLtype>GaussianProcess | GaussianProcess</SKLtype>`.

In addition to this XML node, several others are available:

- **<regr>**, *string, optional field*, is a regression function returning an array of outputs of the linear regression functional basis. The number of observations `n_samples` should be greater than the size `p` of this basis. Available built-in regression models are ‘constant,’ ‘linear,’ and ‘quadratic.’  
*Default: constant*
- **<corr>**, *string, optional field*, is a stationary autocorrelation function returning the autocorrelation between two points  $x$  and  $x'$ . Default assumes a squared-exponential autocorrelation model. Built-in correlation models are ‘absolute\_exponential,’ ‘squared\_exponential,’ ‘generalized\_exponential,’ ‘cubic,’ and ‘linear.’  
*Default: squared\_exponential*
- **<beta0>**, *float, array-like, optional field*, specifies the regression weight vector to perform Ordinary Kriging (OK).  
*Default: None*
- **<storage\_mode>**, *string, optional field*, specifies whether the Cholesky decomposition of the correlation matrix should be stored in the class (`storage_mode = ‘full’`) or not (`storage_mode = ‘light’`).  
*Default: full*
- **<verbose>**, *boolean, optional field*, use verbose mode when fitting the model.  
*Default: False*
- **<theta0>**, *float, array-like, optional field*, is an array with shape `(n_features, )` or `(1, )`. This represents the parameters in the autocorrelation model. If `thetaL` and `thetaU` are also specified, `theta0` is considered as the starting point for the maximum likelihood estimation of the best set of parameters.  
*Default: [1e-1]*
- **<thetaL>**, *float, array-like, optional field*, is an array with shape matching that defined by **<theta0>**. Lower bound on the autocorrelation parameters for maximum likelihood estimation.  
*Default: None*
- **<thetaU>**, *float, array-like, optional field*, is an array with shape matching that defined by **<theta0>**. Upper bound on the autocorrelation parameters for maximum likelihood estimation.  
*Default: None*

- **<normalize>**, *boolean, optional field*, if True, the input  $X$  and observations  $y$  are centered and reduced w.r.t. means and standard deviations estimated from the `n_samples` observations provided.  
*Default: True*
- **<nugget>**, *float, optional field*, specifies a nugget effect to allow smooth predictions from noisy data. The nugget is added to the diagonal of the assumed training covariance. In this way it acts as a Tikhonov regularization in the problem. In the special case of the squared exponential correlation function, the nugget mathematically represents the variance of the input values.  
*Default: 10 \* MACHINE\_EPSILON*
- **<optimizer>**, *string, optional field*, specifies the optimization algorithm to be used. Available optimizers are: 'fmin\_cobyla', 'Welch'.  
*Default: fmin\_cobyla*
- **<random\_start>**, *integer, optional field*, sets the number of times the Maximum Likelihood Estimation should be performed from a random starting point. The first MLE always uses the specified starting point (`theta0`), the next starting points are picked at random according to an exponential distribution (log-uniform on [`thetaL`, `thetaU`]).  
*Default: 1*
- **<random\_state>**, *integer, optional field*, is the seed of the internal random number generator.  
*Default: None*

### Example:

```

<Simulation>
...
<Models>
...
<ROM name='aUserDefinedName' subType='SciKitLearn'>
  <Features>var1,var2,var3</Features>
  <Target>result1</Target>
  <SKLtype>linear_model|LinearRegression</SKLtype>
  <fit_intercept>True</fit_intercept>
  <normalize>False</normalize>
</ROM>
...
</Models>
...
</Simulation>

```

## 15.4 External Model

As the name suggests, an external model is an entity that is embedded in the RAVEN code at run time. This object allows the user to create a python module that is going to be treated as a predefined internal model object. In other words, the **External Model** is going to be treated by RAVEN as a normal external Code (e.g. it is going to be called in order to compute an arbitrary quantity based on arbitrary input).

The specifications of an External Model must be defined within the XML block `<ExternalModel>`. This XML node needs to contain the attributes:

- **name**, *required string attribute*, user-defined name of this External Model. **Note:** As with the other objects, this is the name that can be used to refer to this specific entity from other input blocks in the XML.
- **subType**, *required string attribute*, must be kept empty.
- **ModuleToLoad**, *required string attribute*, file name with its absolute or relative path. **Note:** If a relative path is specified, the code first checks relative to the working directory, then it checks with respect to where the user runs the code. Using the relative path with respect to where the code is run is not recommended.

In order to make the RAVEN code aware of the variables the user is going to manipulate/use in her/his own python Module, the variables need to be specified in the `<ExternalModel>` input block. The user needs to input, within this block, only the variables that RAVEN needs to be aware of (i.e. the variables are going to directly be used by the code) and not the local variables that the user does not want to, for example, store in a RAVEN internal object. These variables are specified within consecutive `<variable>` blocks:

- `<variable>`, *string, required parameter*. In the body of this XML node, the user needs to specify the name of the variable. This variable needs to match a variable used/defined in the external python model.

When the external function variables are defined, at run time, RAVEN initializes them and tracks their values during the simulation. Each variable defined in the `<ExternalModel>` block is available in the module (each method implemented) as a python “self.”

In the External Python module, the user can implement all the methods that are needed for the functionality of the model, but only the following methods, if present, are called by the framework:

- **def \_readMoreXML**, *OPTIONAL METHOD*, can be implemented by the user if the XML input that belongs to this External Model needs to be extended to contain other information. The information read needs to be stored in “self” in order to be available to all the other methods (e.g. if the user needs to add a couple of newer XML nodes with information needed by the algorithm implemented in the “run” method).

- **def initialize**, *OPTIONAL METHOD*, can implement all the actions need to be performed at the initialization stage.
- **def createNewInput**, *OPTIONAL METHOD*, creates a new input with the information coming from the RAVEN framework. In this function the user can retrieve the information coming from the RAVEN framework, during the employment of a calculation flow, and use them to construct a new input that is going to be transferred to the “run” method.
- **def run**, *REQUIRED METHOD*, is the actual location where the user needs to implement the model action (e.g. resolution of a set of equations, etc.). This function is going to receive the Input (or Inputs) generated either by the External Model “createNewInput” method or the internal RAVEN one.

In the following sub-sections, all the methods are going to be analyzed in detail.

#### 15.4.1 Method: `def _readMoreXML`

As already mentioned, the `readMoreXML` method can be implemented by the user if the XML input that belongs to this External Model needs to be extended to contain other information. The information read needs to be stored in “self” in order to be available to all the other methods (e.g. if the user needs to add a couple of newer XML nodes with information needed by the algorithm implemented in the “run” method). If this method is implemented in the **External Model**, RAVEN is going to call it when the node `<ExternalModel>` is found parsing the XML input file. The method receives from RAVEN an attribute of type “xml.etree.ElementTree”, containing all the sub-nodes and attribute of the XML block `<ExternalModel>`.

Example XML:

```
<Simulation>
...
<Models>
...
<ExternalModel name='AnExtModule' subType=''
  ModuleToLoad='path_to_external_module'>
  <variable>sigma</variable>
  <variable>rho</variable>
  <variable>outcome</variable>
  <!--
    here we define other XML nodes RAVEN does not read
    automatically.
    We need to implement, in the external module
    'AnExtModule' the readMoreXML method
  -->
```

```

        <newNodeWeNeedToRead>
            whatNeedsToBeRead
        </newNodeWeNeedToRead>
    </ExternalModel>
    ...
</Models>
...
</Simulation>

```

Corresponding Python function:

```

def _readMoreXML(self, xmlNode):
    # the xmlNode is passed in by RAVEN framework
    # <newNodeWeNeedToRead> is unknown (in the RAVEN framework)
    # we have to read it on our own
    # get the node
    ourNode = xmlNode.find('newNodeWeNeedToRead')
    # get the information in the node
    self.ourNewVariable = ourNode.text
    # end function

```

#### 15.4.2 Method: `def initialize`

The **initialize** method can be implemented in the **External Model** in order to initialize some variables needed by it. For example, it can be used to compute a quantity needed by the “run” method before performing the actual calculation). If this method is implemented in the **External Model**, RAVEN is going to call it at the initialization stage of each “Step” (see section 17. RAVEN will communicate, through a set of method attributes, all the information that are generally needed to perform a initialization:

- `runInfo`, a dictionary containing information regarding how the calculation is set up (e.g. number of processors, etc.). It contains the following attributes:
  - `DefaultInputFile` – default input file to use
  - `SimulationFiles` – the xml input file
  - `ScriptDir` – the location of the pbs script interfaces
  - `FrameworkDir` – the directory where the framework is located
  - `WorkingDir` – the directory where the framework should be running
  - `TempWorkingDir` – the temporary directory where a simulation step is run

- NumMPI – the number of mpi process by run
  - NumThreads – number of threads by run
  - numProcByRun – total number of core used by one run (number of threads by number of mpi)
  - batchSize – number of contemporaneous runs
  - ParallelCommand – the command that should be used to submit jobs in parallel (mpi)
  - numNode – number of nodes
  - procByNode – number of processors by node
  - totalNumCoresUsed – total number of cores used by driver
  - queueingSoftware – queueing software name
  - stepName – the name of the step currently running
  - precommand – added to the front of the command that is run
  - postcommand – added after the command that is run
  - delSucLogFiles – if a simulation (code run) has not failed, delete the relative log file (if True)
  - deleteOutExtension – if a simulation (code run) has not failed, delete the relative output files with the listed extension (comma separated list, for example: ‘e,r,txt’)
  - mode – running mode, curenly the only mode supported is mpi (but custom modes can be created)
  - *expectedTime* – how long the complete input is expected to run
  - *logfileBuffer* – logfile buffer size in bytes
- inputs, a list of all the inputs that have been specified in the “Step” using this model.

In the following an example is reported:

```
def initialize(self, runInfo, inputs):
    # Let's suppose we just need to initialize some variables
    self.sigma = 10.0
    self.rho    = 28.0
    # end function
```

### 15.4.3 Method: **def createNewInput**

The **createNewInput** method can be implemented by the user to create a new input with the information coming from the RAVEN framework. In this function, the user can retrieve the information



coming from the RAVEN framework, during the employment of a calculation flow, and use them to construct a new input that is going to be transferred to the “run” method. The new input created needs to be returned to RAVEN (i.e. “return NewInput”). RAVEN communicates, through a set of method attributes, all the information that are generally needed to create a new input:

- `inputs`, *python list*, a list of all the inputs that have been defined in the “Step” using this model.
- `samplerType`, *string*, the type of Sampler, if a sampling strategy is employed; will be `None` otherwise.
- `Kwargs`, *dictionary*, a dictionary containing several pieces of information (that can change based on the “Step” type). If a sampling strategy is employed, this dictionary contains another dictionary identified by the keyword “SampledVars”, in which the variables perturbed by the sampler are reported.

**Note:** If the “Step” that is using this Model has as input(s) an object of main class type “DataObjects” (see Section 12), the internal “createNewInput” method is going to convert it in a dictionary of values.

Here we present an example:

```
def createNewInput(self, inputs, samplerType, **Kwargs):
    # in here the actual createNewInput of the
    # model is implemented
    if samplerType == 'MonteCarlo':
        avariable = inputs['something']*inputs['something2']
    else:
        avariable = inputs['something']/inputs['something2']
    return avariable*Kwargs['SampledVars']['aSampledVar']
```

#### 15.4.4 Method: `def run`

As stated previously, the only method that *must* be present in an External Module is the **run** function. In this function, the user needs to implement the algorithm that RAVEN will execute. The **run** method is generally called after having inquired the “createNewInput” method (either the internal or the user-implemented one). The only attribute this method is going to receive is a Python list of inputs (the inputs coming from the `createNewInput` method). If the user wants RAVEN to collect the results of this method, the outcomes of interest need to be stored in “self.” **Note:** RAVEN is trying to collect the values of the variables listed only in the `<ExternalModel>` XML block.

In the following an example is reported:

```
def run(self, Input) :
    # in here the actual run of the
    # model is implemented
    input = Input[0]
    self.outcome = self.sigma*self.rho*input[``whatEver'']
```

## 15.5 PostProcessor

A Post-Processor (PP) can be considered as an action performed on a set of data or other type of objects. Most of the post-processors contained in RAVEN, employ a mathematical operation on the data given as “input”. RAVEN supports several different types of PPs.

Currently, the following PPs are available in RAVEN:

- **BasicStatistics**
- **ComparisonStatistics**
- **SafestPoint**
- **LimitSurface**
- **LimitSurfaceIntegral**
- **External**
- **TopologicalDecomposition**
- **DataMining**

The specifications of these PPs must be defined within the XML block **<PostProcessor>**. This XML node needs to contain the attributes:

- **name**, *required string attribute*, user-defined identifier of this post-processor. **Note:** As with other objects, this is the name that can be used to refer to this specific entity from other input XML blocks.
- **subType**, *required string attribute*, defines which of the post-processors needs to be used, choosing among the previously reported types. This choice conditions the subsequent required and/or optional **<PostProcessor>** sub nodes.

As already mentioned, all the types and meaning of the remaining sub-nodes depend on the post-processor type specified in the attribute **subType**. In the following sections the specifications of each type are reported.

## 15.5.1 BasicStatistics

The **BasicStatistics** post-processor is the container of the algorithms to compute many of the most important statistical quantities. In order to use the *BasicStatistics post-processor* PP, the user needs to set the **subType** of a **<PostProcessor>** node:

```
<PostProcessor subType='BasicStatistics' />
```

Several sub-nodes are available:

- **<what>**, *comma separated string, required field*, List of quantities to be computed. Currently the quantities available are:
  - **covariance**: covariance matrix
  - **NormalizedSensitivity**: matrix of normalized sensitivity coefficients
  - **VarianceDependentSensitivity**: matrix of sensitivity coefficients dependent on the variance of the variables
  - **sensitivity**: matrix of sensitivity coefficients of output space to input space
  - **pearson**: matrix of correlation coefficients
  - **expectedValue**: expected value or mean
  - **sigma**: standard deviation
  - **variationCoefficient**: coefficient of variation (sigma/expected value)
  - **variance**: variance
  - **skewness**: skewness
  - **kurtosis**: excess kurtosis (also known as Fisher’s kurtosis)
  - **median**: median
  - **percentile**: the percentile. If this quantity is inputted as *percentile* the 5% and 95% percentile(s) are going to be computed. Otherwise the user can specify this quantity as *percentile\_X%*, where *X* represents the requested percentile (an integer value between 1 and 100).

**Note:** If the weights are present in the system then weighted quantities are calculated automatically.

If all the quantities need to be computed, the user can input in the body of **<what>** the string “all”.

- **<biased>**, *string (boolean), optional field*, if *True* biased quantities are going to be calculated, if *False* unbiased.  
*Default: False*

- **<parameters>**, *comma separated string, required field*, lists the parameters on which the previous operations need to be applied (e.g., massFlow, Temperature)
- **<methodsToRun>**, *comma separated string, optional field*, specifies the method names of an external Function that need to be run before computing any of the predefined quantities. If this XML node is specified, the **<Function>** node must be present.  
*Default: None*
- **Assembler Objects** These objects are either required or optional depending on the functionality of the Adaptive Sampler. The objects must be listed with a rigorous syntax that, except for the xml node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to map those within the simulation framework:
  - **class**, *required string attribute*, it is the main “class” the listed object is from;
  - **type**, *required string attribute*, it is the object identifier or sub-type.

The **BasicStatistics** post-processor approach optionally accepts the following object type:

- **<Function>**, *string, required field*, The body of this xml block needs to contain the name of an External Function defined within the **<Functions>** main block (see section 16). This object needs to contain the methods listed in the node **<methodsToRun>**.

### Example:

```

<Simulation>
...
<Models>
...
<PostProcessor name='aUserDefinedName'
  subType='BasicStatistics' verbosity='debug'>
  <!-- Here you can specify what type of figure of merit
    you need to compute: expectedValue, sigma, variance,
      kurtosis, pearson, covariance, etc.
  -->
  <what>expectedValue</what>
  <parameters>x01,x02</parameters>
  <methodsToRun>failureProbability</methodsToRun>
</PostProcessor>
...
</Models>
...
</Simulation>

```

## 15.5.2 ComparisonStatistics

The **ComparisonStatistics** post-processor computes statistics for comparing two different dataObjects. This is an experimental post-processor, and it will definitely change as it is further developed.

There are four nodes that are used in the post-processor.

- **<kind>**: specifies information to use for comparing the data that is provided. This takes either `uniformBins` which makes the bin width uniform or `equalProbability` which makes the number of counts in each bin equal. It can take the following attributes:
  - **numBins** which takes a number that directly specifies the number of bins
  - **binMethod** which takes a string that specifies the method used to calculate the number of bins. This can be either `square-root` or `sturges`.
- **<compare>**: specifies the data to use for comparison. This can either be a normal distribution or a dataObjects:
  - **<data>**: This will specify the data that is used. The different parts are separated by `|`'s.
  - **<reference>**: This specifies a reference distribution to be used. It takes distribution to use that is defined in the distributions block. A name parameter is used to tell which distribution is used.
- **<fz>**: If the text is true, then extra comparison statistics for using the  $f_z$  function are generated. These take extra time, so are not on by default.
- **<interpolation>**: This switches the interpolation used for the cdf and the pdf functions between the default of `quadratic` or `linear`.

The **ComparisonStatistics** post-processor generates a variety of data. First for each data provided, it calculates bin boundaries, and counts the numbers of data points in each bin. From the numbers in each bin, it creates a cdf function numerically, and from the cdf takes the derivative to generate a pdf. It also calculates statistics of the data such as mean and standard deviation. The post-processor can generate either a CSV file or a PointSet.

The post-processor uses the generated pdf and cdf function to calculate various statistics. The first is the cdf area difference which is:

$$cdf\_area\_difference = \int_{-\infty}^{\infty} \|CDF_a(x) - CDF_b(x)\| dx \quad (11)$$

This gives an idea about how far apart the two pieces of data are, and it will have units of  $x$ .

The common area between the two pdfs is calculated. If there is perfect overlap, this will be 1.0, if there is no overlap, this will be 0.0. The formula used is:

$$pdf\_common\_area = \int_{-\infty}^{\infty} \min(PDF_a(x), PDF_b(x)) dx \quad (12)$$

The difference pdf between the two pdfs is calculated. This is calculated as:

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(x - z) dx \quad (13)$$

This produces a pdf that contains information about the difference between the two pdfs. The mean can be calculated as (and will be calculated only if fz is true):

$$\bar{z} = \int_{-\infty}^{\infty} z f_Z(z) dz \quad (14)$$

The mean can be used to get an signed difference between the pdfs, which shows how their means compare.

The variance of the difference pdf can be calculated as (and will be calculated only if fz is true):

$$var = \int_{-\infty}^{\infty} (z - \bar{z})^2 f_Z(z) dz \quad (15)$$

The sum of the difference function is calculated if fz is true, and is:

$$sum = \int_{-\infty}^{\infty} f_z(z) dz \quad (16)$$

This should be 1.0, and if it is different that points to approximations in the calculation.

### Example:

```
<Simulation>
...
<Models>
...
<PostProcessor name="stat_stuff"
  subType="ComparisonStatistics">
<kind binMethod='sturges'>uniformBins</kind>
<compare>
  <data>OriData|Output|tsin_TEMPERATURE</data>
  <reference name='normal_410_2' />
</compare>
</compare>
```

```

    <data>OriData|Output|tsin_TEMPERATURE</data>
    <data>OriData|Output|tsout_TEMPERATURE</data>
  </compare>
</PostProcessor>
<PostProcessor name="stat_stuff2"
  subType="ComparisonStatistics">
  <kind numBins="6">equalProbability</kind>
  <compare>
    <data>OriData|Output|tsin_TEMPERATURE</data>
  </compare>
  <Distribution class='Distributions'
    type='Normal'>normal_410_2</Distribution>
  </PostProcessor>
  ...
</Models>
...
<Distributions>
  <Normal name='normal_410_2'>
    <mean>410.0</mean>
    <sigma>2.0</sigma>
  </Normal>
</Distributions>
</Simulation>

```

### 15.5.3 SafestPoint

The **SafestPoint** post-processor provides the coordinates of the farthest point from the limit surface that is given as an input. The safest point coordinates are expected values of the coordinates of the farthest points from the limit surface in the space of the “controllable” variables based on the probability distributions of the “non-controllable” variables.

The term “controllable” identifies those variables that are under control during the system operation, while the “non-controllable” variables are stochastic parameters affecting the system behaviour randomly.

The “SafestPoint” post-processor requires the set of points belonging to the limit surface, which must be given as an input. The probability distributions as “Assembler Objects” are required in the “Distribution” section for both “controllable” and “non-controllable” variables.

The sampling method used by the “SafestPoint” is a “value” or “CDF” grid. At present only the “equal” grid type is available.

In order to use the *Safest Point* PP, the user needs to set the `subType` of a `<PostProcessor>` node:

```
<PostProcessor subType=' SafestPoint' />.
```

Several sub-nodes are available:

- `<Distribution>`, *Required*, represents the probability distributions of the “controllable” and “non-controllable” variables. These are **Assembler Objects**, each of these nodes must contain 2 attributes that are used to identify those within the simulation framework:
  - `class`, *required string attribute*, is the main “class” the listed object is from.
  - `type`, *required string attribute*, is the object identifier or sub-type.
- `<controllable>` lists the controllable variables. Each variable is associated with its name and the two items below:
  - `<distribution>` names the probability distribution associated with the controllable variable.
  - `<grid>` specifies the `type`, `steps`, and tolerance of the sampling grid.
- `<non-controllable>` lists the non-controllable variables. Each variable is associated with its name and the two items below:
  - `<distribution>` names the probability distribution associated with the non-controllable variable.
  - `<grid>` specifies the `type`, `steps`, and tolerance of the sampling grid.

#### Example:

```
<Simulation>
...
  <Models>
    ...
    <PostProcessor name='SP' subType='SafestPoint'>
      <Distribution class='Distributions'
        type='Normal'>x1_dst</Distribution>
      <Distribution class='Distributions'
        type='Normal'>x2_dst</Distribution>
      <Distribution class='Distributions'
        type='Normal'>gammay_dst</Distribution>
      <controllable>
        <variable name='x1'>
```



```

    <distribution>x1_dst</distribution>
    <grid type='value' steps='20'>1</grid>
  </variable>
  <variable name='x2'>
    <distribution>x2_dst</distribution>
    <grid type='value' steps='20'>1</grid>
  </variable>
</controllable>
<non-controllable>
  <variable name='gammay'>
    <distribution>gammay_dst</distribution>
    <grid type='value' steps='20'>2</grid>
  </variable>
</non-controllable>
</PostProcessor>
...
</Models>
...
</Simulation>

```

#### 15.5.4 LimitSurface

The **LimitSurface** post-processor is aimed to identify the transition zones that determine a change in the status of the system (Limit Surface).

In order to use the *LimitSurface* PP, the user needs to set the **subType** of a **<PostProcessor>** node:

```
<PostProcessor subType='LimitSurface' />.
```

Several sub-nodes are available:

- **<parameters>**, *comma separated string, required field*, lists the parameters that define the uncertain domain and from which the LS needs to be computed.
- **<tolerance>**, *float, optional field*, sets the absolute value (in CDF) of the convergence tolerance. This value defines the coarseness of the evaluation grid.  
*Default: 1.0e-4*
- **<side>**, *string, optional field*, in this node the user can specify which side of the limit surface needs to be computed. Three options are available:  
*negative*, Limit Surface corresponding to the goal function value of “-1”;

*positive*, Limit Surface corresponding to the goal function value of “1”;  
*both*, either positive and negative Limit Surface is going to be computed.  
*Default: negative*

- **Assembler Objects** These objects are either required or optional depending on the functionality of the Adaptive Sampler. The objects must be listed with a rigorous syntax that, except for the xml node tag, is common among all the objects. Each of these nodes must contain 2 attributes that are used to map those within the simulation framework:
  - **class**, *required string attribute*, is the main “class” of the listed object. For example, it can be “Models,” “Functions,” etc.
  - **type**, *required string attribute*, is the object identifier or sub-type. For example, it can be “ROM,” “External,” etc.

The **LimitSurface** post-processor requires or optionally accepts the following objects’ types:

- **<ROM>**, *string, optional field*, body of this xml node must contain the name of a ROM defined in the **<Models>** block (see section 15.3).
- **<Function>**, *string, required field*, the body of this xml block needs to contain the name of an External Function defined within the **<Functions>** main block (see section 16). This object represents the boolean function that defines the transition boundaries. This function must implement a method called `_residuumSign(self)`, that returns either -1 or 1, depending on the system conditions (see section 16).

#### Example:

```
<Simulation>
...
<Models>
...
  <PostProcessor name="computeLimitSurface"
    subType='LimitSurface' verbosity='debug'>
    <parameters>x0,y0</parameters>
    <ROM class='Models' type='ROM'>Acc</ROM>
    <!-- Here, you can add a ROM defined in Models block.
         If it is not Present, a nearest neighbor algorithm
         will be used.
    -->
    <Function class='Functions' type='External'>
      goalFunctionForLimitSurface
    </Function>
  </PostProcessor>
...
```

```
</Models>
...
</Simulation>
```

### 15.5.5 LimitSurfaceIntegral

The **LimitSurfaceIntegral** post-processor is aimed to compute the likelihood (probability) of the event, whose boundaries are represented by the Limit Surface (either from the LimitSurface post-processor or Adaptive sampling strategies). The inputted Limit Surface needs to be, in the **Post-Process** step, of type **PointSet** and needs to contain both boundary sides (-1.0, +1.0).

The **LimitSurfaceIntegral** post-processor accepts as outputs both files (CSV) and/or **PointSets**.

In order to use the *LimitSurfaceIntegral* PP, the user needs to set the **subType** of a **<PostProcessor>** node:

```
<PostProcessor subType='LimitSurfaceIntegral' />.
```

Several sub-nodes are available:

- **<variable>**, *XML node, required parameter* will specify one attribute:
  - **name**, *required string attribute*, user-defined name of this variable.

This **<variable>** recognizes the following child node:

- **<distribution>**, *string, optional field*, name of the distribution that is associated to this variable. Its name needs to be contained in the **<Distributions>** block explained in Section 10. If this node is not present, the **<lowerBound>** and **<upperBound>** XML nodes must be inputted.
  - **<lowerBound>**, *float, optional field*, lower limit of integration domain for this dimension (variable). If this node is not present, the **<distribution>** XML node must be inputted.
  - **<upperBound>**, *float, optional field*, upper limit of integration domain for this dimension (variable). If this node is not present, the **<distribution>** XML node must be inputted.
- **<tolerance>**, *float, optional field*, specifies the tolerance for numerical integration confidence.  
*Default: 1.0e-4*
  - **<integralType>**, *string, optional field*, specifies the type of integrations that need to be used. Currently only MonteCarlo integration is available  
*Default: MonteCarlo*

- **<seed>**, *integer, optional field*, specifies the random number generator seed.  
*Default: 20021986*
- **<target>**, *string, optional field*, specifies the target name that represents the  $f(\bar{x})$  that needs to be integrated.  
*Default: last output found in the inputted PointSet*

### Example:

```

<Simulation>
...
<Models>
...
  <PostProcessor name="LimitSurfaceIntegralDistributions"
    subType='LimitSurfaceIntegral'>
      <tolerance>0.0001</tolerance>
      <integralType>MonteCarlo</integralType>
      <seed>20021986</seed>
      <target>goalFunctionOutput</target>
      <variable name='x0'>
        <distribution>x0_distrib</distribution>
      </variable>
      <variable name='y0'>
        <distribution>y0_distrib</distribution>
      </variable>
    </PostProcessor>
    <PostProcessor name="LimitSurfaceIntegralLowerUpperBounds"
      subType='LimitSurfaceIntegral'>
        <tolerance>0.0001</tolerance>
        <integralType>MonteCarlo</integralType>
        <seed>20021986</seed>
        <target>goalFunctionOutput</target>
        <variable name='x0'>
          <lowerBound>-2.0</lowerBound>
          <upperBound>12.0</upperBound>
        </variable>
        <variable name='y0'>
          <lowerBound>-1.0</lowerBound>
          <upperBound>11.0</upperBound>
        </variable>
      </PostProcessor>
    ...
  </Models>

```

```
...
</Simulation>
```

### 15.5.6 External

The **External** post-processor will execute an arbitrary python function defined externally using the *Functions* interface (see Section 16 for more details).

In order to use the *External PP*, the user needs to set the **subType** of a **<PostProcessor>** node:

```
<PostProcessor subType='External' />.
```

Several sub-nodes are available:

- **<method>**, *comma separated string, required field*, lists the method names of an external Function that will be computed (each returning a post-processing value).
- **<Function>**, *xml node, required string field*, specifies the name of a Function where the *methods* listed above are defined. **Note:** This name should match one of the Functions defined in the **<Functions>** block of the input file. The objects must be listed with a rigorous syntax that, except for the xml node tag, is common among all the objects. Each of these sub-nodes must contain 2 attributes that are used to map them within the simulation framework:
  - **class**, *required string attribute*, is the main “class” the listed object is from, the only acceptable class for this post-processor is **'Functions'**;
  - **type**, *required string attribute*, is the object identifier or sub-type, the only acceptable type for this post-processor is **'External'**.

#### Example:

```
<Simulation>
...
<Models>
...
  <PostProcessor name="externalPP" subType='External'
    verbosity='debug' >
    <method>Delta,Sum</method>
    <Function class='Functions'
      type='External'>operators</Function>
```

```

        <!--Here, you can add a Function defined in the
            Functions block. This should be present or
            else RAVEN will not know where to find the
            defined methods.
        -->
    </PostProcessor>
    ...
</Models>
...
</Simulation>

```

### 15.5.7 TopologicalDecomposition

The **TopologicalDecomposition** post-processor will compute an approximated hierarchical Morse-Smale complex which will add two columns to a dataset, namely `minLabel` and `maxLabel` that can be used to decompose a dataset.

In order to use the **TopologicalDecomposition** post-processor, the user needs to set the attribute `subType`: `<PostProcessor subType='TopologicalDecomposition'>`. The following is a list of acceptable sub-nodes:

- `<graph>`, *string, optional field*, specifies the type of neighborhood graph used in the algorithm, available options are:
  - beta skeleton
  - relaxed beta skeleton
  - approximate knn

*Default: beta skeleton*

- `<gradient>`, *string, optional field*, specifies the method used for estimating the gradient, available options are:
  - steepest

*Default: steepest*

- `<beta>`, *float in the range: (0,2], optional field*, is only used when the `<graph>` is set to beta skeleton or relaxed beta skeleton.
 

*Default: 1.0*

- **<knn>**, *integer, optional field*, is the number of neighbors when using the ' **approximate knn** ' for the **<graph>** sub-node and used to speed up the computation of other graphs by using the approximate knn graph as a starting point for pruning. -1 means use a fully connected graph.

*Default: -1*

- **<weighted>**, *boolean, optional*, a flag that specifies whether the regression models should be probability weighted.

*Default: False*

- **<persistence>**, *string, optional field*, specifies how to define the hierarchical simplification by assigning a value to each local minimum and maximum according to the one of the strategy options below:

- `difference` - The function value difference between the extremum and its closest-valued neighboring saddle.
- `probability` - The probability integral computed as the sum of the probability of each point in a cluster divided by the count of the cluster.
- `count` - The count of points that flow to or from the extremum.

*Default: difference*

- **<simplification>**, *float, optional field*, specifies the amount of noise reduction to apply before returning labels.

*Default: 0*

- **<parameters>**, *comma separated string, required field*, lists the parameters defining the input space.

- **<response>**, *string, required field*, is a single variable name defining the scalar output space.

### Example:

```
<Simulation>
...
<Models>
...
<PostProcessor name="***" subType='TopologicalDecomposition'>
  <graph>beta skeleton</graph>
  <gradient>steepest</gradient>
  <beta>1</beta>
  <knn>8</knn>
  <normalization>None</normalization>
```

```

    <parameters>X,Y</parameters>
    <response>Z</response>
    <weighted>>true</weighted>
    <simplification>0.3</simplification>
    <persistence>difference</persistence>
  </PostProcessor>
  ...
<Models>
  ...
<Simulation>

```

### 15.5.8 DataMining (*aka: Knowledge Discovery in DataBases*)

Knowledge discovery in databases (KDD) is the process of discovering useful knowledge from a collection of data. This widely used data mining technique is a process that includes data preparation and selection, data cleansing, incorporating prior knowledge on data sets and interpreting accurate solutions from the observed results. Major KDD application areas include marketing, fraud detection, telecommunication and manufacturing.

DataMining is the analysis step of the KDD process. The overall of the data mining process is to extract information from a data set and transform it into an understandable structure for further use. The actual data mining task is the automatic or semi-automatic analysis of large quantities of data to extract previously unknown, interesting patterns such as groups of data records (cluster analysis), unusual records (anomaly detection), and dependencies (association rule mining). In order to use the **DataMining** post-processor, the user needs to set the attribute **subType**:

```
<PostProcessor subType= 'DataMiningPostProcessor' >.
```

The following is a list of acceptable sub-nodes:

- **<KDD>** *string,required field*, the subnodes specifies the necessary information for the algorithm to be used in the postprocessor. the **<KDD>** has the required attribute: **lib**, the name of the library the algorithm belongs to. Current algorithms applied in the KDD model is based on SciKit-Learn library. Thus currently there is only one library:
  - **'SciKitLearn'**

The algorithm for the dataMining is chosen by the subnode **<SKLType>** under the parent node **<KDD>**. The format is same as in 15.3.6. However, for the completeness sake , it is repeated here.

The data that are used in the training of the **DataMining** postprocessor are supplied with subnode **<Features>** in the parent node **<KDD>**.



- **<SKLtype>**, *vertical bar ( | ) separated string, required field*, contains a string that represents the data mining algorithm to be used. As mentioned, its format is: **<SKLtype>**mainSKLclass|algorithm**</SKLtype>** where the first word (before the “|” symbol) represents the main class of algorithms, and the second word (after the “|” symbol) represents the specific algorithm.
- **<Features>**, *string, required field*, defines the data to be used for training the data mining algorithm. It can be:
  - the name of the variable in the defined dataObject entity
  - the location (i.e. input or output). In this case the data mining is applied to all the variables in the defined space.

The **<KDD>** node can have either optional or required subnodes depending on the dataMining algorithm used. The possible subnodes will be described separately for each algorithm below.

All the available algorithms are described in the following sections.

### 15.5.8.1 Gaussian mixture models

A Gaussian mixture model is a probabilistic model that assumes all the data points are generated from a mixture of a finite number of Gaussian distributions with unknown parameters. Scikit-learn implements different classes to estimate Gaussian mixture models, that correspond to different estimation strategies, detailed below.

#### 15.5.8.1.1 GMM classifier

The GMM object implements the expectation-maximization (EM) algorithm for fitting mixture-of-Gaussian models. The GMM comes with different options to constrain the covariance of the difference classes estimated: spherical, diagonal, tied or full covariance.

In order to use the *Gaussian Mixture Model*, the user needs to set the sub-node:

**<SKLtype>**mixture|GMM**</SKLtype>**.

In addition to this XML node, several others are available:

- **<n\_components>**, *integer, optional field* Number of mixture components.  
*Default: 1*
- **<covariance\_type>**, *string, optional field*, describes the type of covariance parameters to use. Must be one of ‘spherical’, ‘tied’, ‘diag’, ‘full’.  
*Default: diag*

- **<random\_state>**, *integer seed or random number generator instance, optional field*, A random number generator instance  
*Default: 0 or None*
- **<min\_covar>**, *float, optional field*, Floor on the diagonal of the covariance matrix to prevent overfitting.  
*Default: 1e-3.*
- **<thresh>**, *float, optional field*, convergence threshold.  
*Default: 0.01*
- **<n\_iter>**, *integer, optional field*, Number of EM iterations to perform.  
*Default: 100*
- **<n\_init>**, *integer, optional*, Number of initializations to perform. the best results is kept.  
*Default: 1*
- **<params>**, *string, optional field*, Controls which parameters are updated in the training process. Can contain any combination of ‘w’ for weights, ‘m’ for means, and ‘c’ for covars.  
*Default: ‘wmc’*
- **<init\_params>**, *string, optional field*, Controls which parameters are updated in the initialization process. Can contain any combination of ‘w’ for weights, ‘m’ for means, and ‘c’ for covars.  
*Default: ‘wmc’.*

### Example:

```

<Simulation>
...
<Models>
...
  <PostProcessor name='PostProcessorName'
    subType='DataMiningPostProcessor'>
    <KDD lib='SciKitLearn'>
      <Features>variableName</Features>
      <SKLtype>mixture|GMM</SKLtype>
      <n_components>2</n_components>
      <covariance_type>spherical</covariance_type>
    </KDD>
  </PostProcessor>
...
<Models>
...

```

### 15.5.8.1.2 Dirichlet Process GMM Classifier (DPGMM)

The DPGMM implements a variant of the Gaussian mixture model with a variable (but bounded) number of components using the Dirichlet Process. The API is identical to GMM. This class doesn't require the user to choose the number of components, and at the expense of extra computational time the user only needs to specify a loose upper bound on this number and a concentration parameter.

In order to use the *Dirichlet Process Gaussian Mixture Model*, the user needs to set the subnode:

```
<SKLtype>mixture|DPGMM</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_components>**, *integer, optional field* Number of mixture components.  
*Default: 1*
- **<covariance\_type>**, *string, optional field*, describes the type of covariance parameters to use. Must be one of 'spherical', 'tied', 'diag', 'full'.  
*Default: diag*
- **<alpha>**, *float, optional field*, represents the concentration parameter of the dirichlet process. Intuitively, the Dirichlet Process is as likely to start a new cluster for a point as it is to add that point to a cluster with alpha elements. A higher alpha means more clusters, as the expected number of clusters is  $\alpha * \log(N)$ .  
*Default: 1.*
- **<thresh>**, *float, optional field*, convergence threshold.  
*Default: 0.01*
- **<n\_iter>**, *integer, optional field*, Number of EM iterations to perform.  
*Default: 100*
- **<params>**, *string, optional field*, Controls which parameters are updated in the training process. Can contain any combination of 'w' for weights, 'm' for means, and 'c' for covars.  
*Default: 'wmc'*
- **<init\_params>**, *string, optional field*, Controls which parameters are updated in the initialization process. Can contain any combination of 'w' for weights, 'm' for means, and 'c' for covars.  
*Default: 'wmc'.*

### 15.5.8.1.3 Variational GMM Classifier (VBGMM)

The VBGMM object implements a variant of the Gaussian mixture model with variational inference algorithms. The API is identical to GMM. It is essentially a middle-ground between GMM and DPGMM, as it has some of the properties of the Dirichlet process.

In order to use the *Variational Gaussian Mixture Model*, the user needs to set the sub-node:

```
<SKLtype>mixture|VBGMM</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_components>**, *integer, optional field* Number of mixture components.  
*Default: 1*
- **<covariance\_type>**, *string, optional field*, describes the type of covariance parameters to use. Must be one of 'spherical', 'tied', 'diag', 'full'.  
*Default: diag*
- **<alpha>**, *float, optional field*, represents the concentration parameter of the dirichlet process. Intuitively, the Dirichlet Process is as likely to start a new cluster for a point as it is to add that point to a cluster with alpha elements. A higher alpha means more clusters, as the expected number of clusters is  $\alpha * \log(N)$ .  
*Default: 1.*

### 15.5.8.2 Clustering

Clustering of unlabeled data can be performed with this subType of the DataMining PostProcessor.

An overview of the different clustering algorithms is given in Table2.

**Table 2.** Overview of Clustering Methods

Method name	Parameters	Scalability	Usecase	Geometry (metric used)
K-Means	number of clusters	Very large n_samples, medium n_clusters with MiniBatch code	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
Affinity propagation	damping, sample preference	Not scalable with n_samples	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with n_samples	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral clustering	number of clusters	Medium n_samples, small n_clusters	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters	Large n_samples and n_clusters	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters, linkage type, distance	Large n_samples and n_clusters	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large n_samples, medium n_clusters	Non-flat geometry, uneven cluster sizes	Distances between nearest points
Gaussian mixtures	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers

### 15.5.8.2.1 K-Means Clustering

The KMeans algorithm clusters data by trying to separate samples in  $n$  groups of equal variance, minimizing a criterion known as the inertia or within-cluster sum-of-squares. This algorithm requires the number of clusters to be specified. It scales well to large number of samples and has been used across a large range of application areas in many different fields

In order to use the *K-Means Clustering*, the user needs to set the sub-node:

```
<SKLtype>cluster|KMeans</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_clusters>**, *integer, optional field* The number of clusters to form as well as the number of centroids to generate.  
*Default: 8*
- **<max\_iter>**, *integer, optional field*, Maximum number of iterations of the k-means algorithm for a single run.  
*Default: 300*
- **<n\_init>**, *integer, optional field*, Number of time the k-means algorithm will be run with different centroid seeds. The final results will be the best output of `n_init` consecutive runs in terms of inertia.  
*Default: 3*
- **<init>**, *string, optional*, Method for initialization, 'k-means++', 'random' or an ndarray:
  - 'k-means++' : selects initial cluster centers for k-mean clustering in a smart way to speed up convergence.
  - 'random': choose  $k$  observations (rows) at random from data for the initial centroids.
  - If an ndarray is passed, it should be of shape (`n_clusters`, `n_features`) and gives the initial centers.
- **<precompute\_distances>**, *boolean, optional field*, Precompute distances (if true faster but takes more memory).  
*Default: true*
- **<tol>**, *float, optional field*, Relative tolerance with regards to inertia to declare convergence.  
*Default: 1e-4*
- **<n\_jobs>**, *integer, optional field*, The number of jobs to use for the computation. This works by breaking down the pairwise matrix into  $n$  jobs even slices and computing them in parallel. If -1 all CPUs are used. If 1 is given, no parallel computing code is used at all,

which is useful for debugging. For `n_jobs` below -1,  $(n\_cpus + 1 + n\_jobs)$  are used. Thus for `n_jobs = -2`, all CPUs but one are used.

*Default: 1*

- **<random\_state>**, *integer or numpy.RandomState, optional field* The generator used to initialize the centers. If an integer is given, it fixes the seed.

*Default: the global numpy random number generator.*

### Example:

```
<Simulation>
...
<Models>
...
  <PostProcessor name='PostProcessorName'
    subType='DataMiningPostProcessor'>
    <KDD lib='SciKitLearn'>
      <Features>variableName</Features>
      <SKLtype>cluster|KMeans</SKLtype>
      <n_clusters>2</n_clusters>
      <tol>0.0001</tol>
      <init>random</init>
    </KDD>
  </PostProcessor>
...
<Models>
...
</Simulation>
```

#### 15.5.8.2.2 Mini Batch K-Means

The `MiniBatchKMeans` is a variant of the `KMeans` algorithm which uses mini-batches to reduce the computation time, while still attempting to optimise the same objective function. Mini-batches are subsets of the input data, randomly sampled in each training iteration.

`MiniBatchKMeans` converges faster than `KMeans`, but the quality of the results is reduced. In practice this difference in quality can be quite small.

In order to use the *Mini Batch K-Means Clustering*, the user needs to set the sub-node:

```
<SKLtype>cluster|MiniBatchKMeans</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_clusters>**, *integer, optional field* The number of clusters to form as well as the number of centroids to generate.  
*Default: 8*
- **<max\_iter>**, *integer, optional field*, Maximum number of iterations of the k-means algorithm for a single run.  
*Default: 100*
- **<max\_no\_improvement>**, *integer, optional field*, Control early stopping based on the consecutive number of mini batches that does not yield an improvement on the smoothed inertia. To disable convergence detection based on inertia, set `max_no_improvement` to `None`.  
*Default: 10*
- **<tol>**, *float, optional field*, Control early stopping based on the relative center changes as measured by a smoothed, variance-normalized of the mean center squared position changes. This early stopping heuristic is closer to the one used for the batch variant of the algorithms but induces a slight computational and memory overhead over the inertia heuristic. To disable convergence detection based on normalized center change, set `tol` to `0.0` (default).  
*Default: 0.0*
- **<batch\_size>**, *integer, optional field*, Size of the mini batches.  
*Default: 100*
- `init_size`, *integer, optional field*, Number of samples to randomly sample for speeding up the initialization (sometimes at the expense of accuracy): the only algorithm is initialized by running a batch KMeans on a random subset of the data. *This needs to be larger than k.*,  
*Default: 3 \* <batch\_size>*
- **<init>**, *string, optional*, Method for initialization, 'k-means++', 'random' or an ndarray:
  - 'k-means++' : selects initial cluster centers for k-mean clustering in a smart way to speed up convergence.
  - 'random': choose k observations (rows) at random from data for the initial centroids.
  - If an ndarray is passed, it should be of shape (n\_clusters, n\_features) and gives the initial centers.
- **<precompute\_distances>**, *boolean, optional field*, Precompute distances (if true faster but takes more memory).  
*Default: true*
- **<n\_init>**, *integer, optional field*, Number of time the k-means algorithm will be run with different centroid seeds. The final results will be the best output of `n_init` consecutive runs in terms of inertia.  
*Default: 3*



- **<compute\_labels>**, *boolean, optional field*, Compute label assignment and inertia for the complete dataset once the minibatch optimization has converged in fit.  
*Default: True*
- **<random\_state>**, *integer or numpy.RandomState, optional field* The generator used to initialize the centers. If an integer is given, it fixes the seed.  
*Default: the global numpy random number generator.*
- **reassignment\_ratio**, **<float, optional field>**, Control the fraction of the maximum number of counts for a center to be reassigned. A higher value means that low count centers are more easily reassigned, which means that the model will take longer to converge, but should converge in a better clustering.  
*Default: 0.01*

### 15.5.8.2.3 Affinity Propagation

AffinityPropagation creates clusters by sending messages between pairs of samples until convergence. A dataset is then described using a small number of exemplars, which are identified as those most representative of other samples. The messages sent between pairs represent the suitability for one sample to be the exemplar of the other, which is updated in response to the values from other pairs. This updating happens iteratively until convergence, at which point the final exemplars are chosen, and hence the final clustering is given.

In order to use the *AffinityPropogation Clustering*, the user needs to set the sub-node:

```
<SKLtype>cluster|AffinityPropogation</SKLtype>.
```

In addition to this XML node, several others are available:

- **<damping>**, *float, optional field*, Damping factor between 0.5 and 1.  
*Default: 0.5*
- **<convergence\_iter>**, *integer, optional field*, Number of iterations with no change in the number of estimated clusters that stops the convergence.  
*Default: 15*
- **<max\_iter>**, *integer, optional field*, Maximum number of iterations.  
*Default: 200*
- **<copy>**, *boolean, optional field*, Make a copy of input data or not.  
*Default: True*
- **<preference>**, *array-like, shape (n\_samples,) or float, optional field*, Preferences for each point - points with larger values of preferences are more likely to be chosen as exem-

plars. The number of exemplars, ie of clusters, is influenced by the input preferences value.  
*Default: If the preferences are not passed as arguments, they will be set to the median of the input similarities.*

- **<affinity>**, *string, optional field*, Which affinity to use. At the moment precomputed and euclidean are supported. euclidean uses the negative squared euclidean distance between points.  
*Default: "euclidean"*
- **<verbose>**, *boolean, optional field*, Whether to be verbose.  
*Default: False*

#### 15.5.8.2.4 Mean Shift

MeanShift clustering aims to discover blobs in a smooth density of samples. It is a centroid based algorithm, which works by updating candidates for centroids to be the mean of the points within a given region. These candidates are then filtered in a post-processing stage to eliminate near-duplicates to form the final set of centroids.

In order to use the *Mean Shift Clustering*, the user needs to set the sub-node:

```
<SKLtype>cluster|MeanShift</SKLtype>.
```

In addition to this XML node, several others are available:

- **<bandwidth>**, *float, optional field*, Bandwidth used in the RBF kernel. If not given, the bandwidth is estimated using *sklearn.cluster.estimate\_bandwidth*; see the documentation for that function for hints on scalability.
- **<seeds>**, *array, shape=[n\_samples, n\_features], optional field*, Seeds used to initialize kernels. If not set, the seeds are calculated by *clustering.get\_bin\_seeds* with bandwidth as the grid size and default values for other parameters.
- **<bin\_seeding>**, *boolean, optional field*, If true, initial kernel locations are not locations of all points, but rather the location of the discretized version of points, where points are binned onto a grid whose coarseness corresponds to the bandwidth. Setting this option to True will speed up the algorithm because fewer seeds will be initialized.  
*Default: False* Ignored if seeds argument is not None.
- **<min\_bin\_freq>**, *integer, optional field*, To speed up the algorithm, accept only those bins with at least min\_bin\_freq points as seeds.  
*Default: 1.*
- **<cluster\_all>**, *boolean, optional field*, If true, then all points are clustered, even those orphans that are not within any kernel. Orphans are assigned to the nearest kernel. If false,

then orphans are given cluster label -1.  
*Default: True*

#### 15.5.8.2.5 Spectral clustering

SpectralClustering does a low-dimension embedding of the affinity matrix between samples, followed by a *KMeans* in the low dimensional space. It is especially efficient if the affinity matrix is sparse and the *pyamg* module is installed.

In order to use the *Spectral Clustering*, the user needs to set the sub-node:

```
<SKLtype>cluster|Spectral</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_clusters>**, *integer, optional field*, The dimension of the projection subspace.  
*Default: 8*
- **<affinity>**, *string, array-like or callable, optional field*, If a string, this may be one of:
  - ‘nearest\_neighbors’,
  - ‘precomputed’,
  - ‘rbf’ or
  - one of the kernels supported by *sklearn.metrics.pairwise\_kernels*.

Only kernels that produce similarity scores (non-negative values that increase with similarity) should be used. This property is not checked by the clustering algorithm.

*Default: ‘rbf’*

- **<gamma>**, *float, optional field*, Scaling factor of RBF, polynomial, exponential  $\chi^2$  and sigmoid affinity kernel. Ignored for *affinity = ‘nearest\_neighbors’*.  
*Default: 1.0*
- **<degree>**, *float, optional field*, Degree of the polynomial kernel. Ignored by other kernels.  
*Default: 3*
- **<coef0>**, *float, optional field*, Zero coefficient for polynomial and sigmoid kernels. Ignored by other kernels.  
*Default: 1*
- **<n\_neighbors>**, *integer, optional field*, Number of neighbors to use when constructing the affinity matrix using the nearest neighbors method. Ignored for *affinity=‘rbf’*.  
*Default: 10*

- **<eigen\_solver>** *string, optional field*, The eigenvalue decomposition strategy to use:
  - None,
  - ‘arpack’,
  - ‘lobpcg’, or
  - ‘amg’

**Note:** AMG requires pyamg to be installed. It can be faster on very large, sparse problems, but may also lead to instabilities

- **<random\_state>**, *integer seed, RandomState instance, or None, optional field*, A pseudo random number generator used for the initialization of the lobpcg eigen vectors decomposition when *eigen\_solver == ‘amg’* and by the K-Means initialization.

*Default: None*

- **<n\_init>**, *integer, optional field*, Number of time the k-means algorithm will be run with different centroid seeds. The final results will be the best output of n\_init consecutive runs in terms of inertia.

*Default: 10*

- **<eigen\_tol>**, *float, optional field*, Stopping criterion for eigendecomposition of the Laplacian matrix when using arpack eigen\_solver.

*Default: 0.0*

- **<assign\_labels>**, *string, optional field*, The strategy to use to assign labels in the embedding space. There are two ways to assign labels after the laplacian embedding:

- ‘kmeans’,
- ‘discretize’

k-means can be applied and is a popular choice. But it can also be sensitive to initialization. Discretization is another approach which is less sensitive to random initialization.

*Default: ‘kmeans’*

- **<kernel\_params>**, *dictionary of string to any, optional field*, Parameters (keyword arguments) and values for kernel passed as callable object. Ignored by other kernels.

*Default: None*

## Notes

If you have an affinity matrix, such as a distance matrix, for which 0 means identical elements, and high values means very dissimilar elements, it can be transformed in a similarity matrix that is well suited for the algorithm by applying the Gaussian (RBF, heat) kernel:

$$np.exp(-X ** 2 / (2. * delta ** 2)) \tag{17}$$

Another alternative is to take a symmetric version of the k nearest neighbors connectivity matrix of the points. If the *pyamg* package is installed, it is used: this greatly speeds up computation.

### 15.5.8.2.6 DBSCAN Clustering

The Density-Based Spatial Clustering of Applications with Noise (DBSCAN) algorithm views clusters as areas of high density separated by areas of low density. Due to this rather generic view, clusters found by DBSCAN can be any shape, as opposed to k-means which assumes that clusters are convex shaped.

In order to use the *DBSCAN Clustering*, the user needs to set the sub-node:

```
<SKLtype>cluster|DBSCAN</SKLtype>.
```

In addition to this XML node, several others are available:

- **<eps>**, *float, optional field*, The maximum distance between two samples for them to be considered as in the same neighborhood.  
*Default: 0.5*
- **<min\_samples>**, *integer, optional field*, The number of samples in a neighborhood for a point to be considered as a core point.  
*Default: 5*
- **<metric>**, *string, or callable, optional field* The metric to use when calculating distance between instances in a feature array. If metric is a string or callable, it must be one of the options allowed by *metrics.pairwise.calculate\_distance* for its metric parameter. If metric is “precomputed”, X is assumed to be a distance matrix and must be square.  
*Default: 'euclidean'*
- **<random\_state>**, *numpy.RandomState, optional field*, The generator used to initialize the centers.  
*Default: numpy.random.*

### 15.5.8.2.7 Clustering performance evaluation

Evaluating the performance of a clustering algorithm is not as trivial as counting the number of errors or the precision and recall of a supervised classification algorithm. In particular any evaluation metric should not take the absolute values of the cluster labels into account but rather if this clustering define separations of the data similar to some ground truth set of classes or satisfying some assumption such that members belong to the same class are more similar than members of different classes according to some similarity metric.

If the ground truth labels are not known, evaluation must be performed using the model itself. The **Silhouette Coefficient** is an example of such an evaluation, where a higher Silhouette Coefficient score relates to a model with better defined clusters. The Silhouette Coefficient is defined for each sample and is composed of two scores:

1. The mean distance between a sample and all other points in the same class.
2. The mean distance between a sample and all other points in the next nearest cluster.

The Silhouette Coefficient  $s$  for a single sample is then given as:

$$s = \frac{b - a}{\max(a, b)} \quad (18)$$

The Silhouette Coefficient for a set of samples is given as the mean of the Silhouette Coefficient for each sample. In normal usage, the Silhouette Coefficient is applied to the results of a cluster analysis.

### Advantages

- The score is bounded between -1 for incorrect clustering and +1 for highly dense clustering. Scores around zero indicate overlapping clusters.
- The score is higher when clusters are dense and well separated, which relates to a standard concept of a cluster.

### Drawbacks

The Silhouette Coefficient is generally higher for convex clusters than other concepts of clusters, such as density based clusters like those obtained through DBSCAN.

## 15.5.8.3 Decomposing signals in components (matrix factorization problems)

### 15.5.8.3.1 Principal component analysis (PCA)

- **Exact PCA and probabilistic interpretation**

Linear Dimensionality reduction using Singular Value Decomposition of the data and keeping only the most significant singular vectors to project the data to a lower dimensional space. In order to use the *Exact PCA*, the user needs to set the sub-node:

`<SKLtype>decomposition|PCA</SKLtype>`.

In addition to this XML node, several others are available:

- `<n_components>`, *integer, None or String, optional field*, Number of components to keep. if
- `<n_components>` is not set all components are kept, *Default: all components*

- **<copy>**, *boolean, optional field*, If False, data passed to fit are overwritten and running fit(X).transform(X) will not yield the expected results, use fit\_transform(X) instead.  
*Default: True*
- **<whiten>**, *boolean, optional field*, When True the components\_ vectors are divided by n\_samples times singular values to ensure uncorrelated outputs with unit component-wise variances. Whitening will remove some information from the transformed signal (the relative variance scales of the components) but can sometime improve the predictive accuracy of the downstream estimators by making there data respect some hard-wired assumptions.  
*Default: False*

**Example:**

```

<Simulation>
...
<Models>
...
  <PostProcessor name='PostProcessorName'
    subtype='DataMiningPostProcessor' >
    <KDD lib='SciKitLearn' >
      <Features>variable1,variable2,variable3,
        variable4,variable5</Features>
      <SKLtype>decomposition|PCA</SKLtype>
      <n_components>2</n_components>
    </KDD>
  </PostProcessor>
...
<Models>
...
</Simulation>

```

- **Randomized (Approximate) PCA**

Linear Dimensionality reduction using Singular Value Decomposition of the data and keeping only the most significant singular vectors to project the data to a lower dimensional space. In order to use the *Randomized PCA*, the user needs to set the sub-node:

```

<SKLtype>decomposition|RandomizedPCA</SKLtype>.

```

In addition to this XML node, several others are available:

- **<n\_components>**, *interger, None or String, optional field*, Number of components to keep. if n\_components is not set all components are kept.  
*Default: all components*

- **<copy>**, *boolean, optional field*, If False, data passed to fit are overwritten and running `fit(X).transform(X)` will not yield the expected results, use `fit_transform(X)` instead.  
*Default: True*
- **<iterated\_power>**, *integer, optional field*, Number of iterations for the power method.  
*Default: 3*
- **<whiten>**, *boolean, optional field*, When True the components\_ vectors are divided by `n_samples` times singular values to ensure uncorrelated outputs with unit component-wise variances. Whitening will remove some information from the transformed signal (the relative variance scales of the components) but can sometime improve the predictive accuracy of the downstream estimators by making there data respect some hard-wired assumptions.  
*Default: False*
- **<random\_state>**, *int, or Random State instance or None, optional field*, Pseudo Random Number generator seed control. If None, use the `numpy.random` singleton.  
*Default: None*

- **Kernel PCA**

Non-linear dimensionality reduction through the use of kernels. In order to use the *Kernel PCA*, the user needs to set the sub-node:

**<SKLtype>**`decomposition|KernelPCA`**</SKLtype>**.

In addition to this XML node, several others are available:

- **<n\_components>**, *interger, None or String, optional field*, Number of components to keep. if `n_components` is not set all components are kept.  
*Default: all components*
- **<kernel>**, *string, optional field*, name of the kernel to be used, options are:
  - \* linear
  - \* poly
  - \* rbf
  - \* sigmoid
  - \* cosine
  - \* precomputed

*Default: linear* **<degree>**, *integer, optional field*, Degree for poly kernels, ignored by other kernels.

*Default: 3* **<gamma>**, *float, optional field*, Kernel coefficient for rbf and poly kernels, ignored by other kernels.

*Default: 1/n\_features*



- **<coef0>**, *float, optional field*, independent term in poly and sigmoig kernels, ignored by other kernels.
- **<kernel\_params>**, *mapping of string to any, optional field*, Parameters (keyword arguments) and values for kernel passed as callable object. Ignored by other kernels.  
*Default: 3*
- **alpha**, *int, optional field*, Hyperparameter of the ridge regression that learns the inverse transform (when `fit_inverse_transform=True`).  
*Default: 1.0*
- **<fit\_inverse\_transform>**, *bool, optional field*, Learn the inverse transform for non-precomputed kernels. (i.e. learn to find the pre-image of a point)  
*Default: False*
- **<eigen\_solver>**, *string, optional field*, Select eigensolver to use. If `n_components` is much less than the number of training samples, `arpack` may be more efficient than the dense eigensolver. Options are:
  - \* auto
  - \* dense
  - \* arpack

*Default: False*

- **tol**, *float, optional field*, convergence tolerance for `arpack`.  
*Default: 0 (optimal value will be chosen by arpack)*
- **max\_iter**, *int, optional field*, maximum number of iterations for `arpack`.  
*Default: None (optimal value will be chosen by arpack)*
- **<remove\_zero\_eig>**, *boolean, optional field*, If `True`, then all components with zero eigenvalues are removed, so that the number of components in the output may be  $\leq n\_components$  (and sometimes even zero due to numerical instability). When `n_components` is `None`, this parameter is ignored and components with zero eigenvalues are removed regardless.  
*Default: True*

- **Sparse PCA**

Finds the set of sparse components that can optimally reconstruct the data. The amount of sparseness is controllable by the coefficient of the L1 penalty, given by the parameter `alpha`. In order to use the *Sparse PCA*, the user needs to set the sub-node:

**<SKLtype>**`decomposition|SparsePCA`**</SKLtype>**.

In addition to this XML node, several others are available:

- **<n\_components>**, *integer, optional field*, Number of sparse atoms to extract.  
*Default: None*

- **<alpha>**, *float, optional field*, Sparsity controlling parameter. Higher values lead to sparser components.  
*Default: 1.0*
- **<ridge\_alpha>**, *float, optional field*, Amount of ridge shrinkage to apply in order to improve conditioning when calling the transform method.  
*Default: 0.01*
- **<max\_iter>**, *float, optional field*, maximum number of iterations to perform.  
*Default: 1000*
- **<tol>**, *float, optional field*, convergence tolerance.  
*Default: 1E-08*
- **<method>**, *string, optional field*, method to use, options are:
  - \* lars: uses the least angle regression method to solve the lasso problem (linear\_model.lars\_path)
  - \* cd: uses the coordinate descent method to compute the Lasso solution (linear\_model.Lasso)

Lars will be faster if the estimated components are sparse.  
*Default: lars*
- **<n\_jobs>**, *int, optional field*, number of parallel runs to run.  
*Default: 1*
- **<U\_init>**, *array of shape (n\_samples, n\_components) , optional field*, Initial values for the loadings for warm restart scenarios  
*Default: None*
- **<V\_init>**, *array of shape (n\_components, n\_features), optional field*, Initial values for the components for warm restart scenarios  
*Default: None*
- **verbose**, *boolean, optional field*, Degree of verbosity of the printed output.  
*Default: False*
- **random\_state**, *int or Random State, optional field*, Pseudo number generator state used for random sampling.  
*Default: None*

- **Mini Batch Sparse PCA**

Finds the set of sparse components that can optimally reconstruct the data. The amount of sparseness is controllable by the coefficient of the L1 penalty, given by the parameter alpha. In order to use the *Mini Batch Sparse PCA*, the user needs to set the sub-node:

**<SKLtype>**decomposition|MiniBatchSparsePCA**</SKLtype>**.

In addition to this XML node, several others are available:

- **<n\_components>**, *integer, optional field*, Number of sparse atoms to extract.  
*Default: None*

- **<alpha>**, *float, optional field*, Sparsity controlling parameter. Higher values lead to sparser components.  
*Default: 1.0*
- **<ridge\_alpha>**, *float, optional field*, Amount of ridge shrinkage to apply in order to improve conditioning when calling the transform method.  
*Default: 0.01*
- **<n\_iter>**, *float, optional field*, number of iterations to perform per mini batch.  
*Default: 100*
- **<callback>**, *callable, optional field*, callable that gets invoked every five iterations.  
  
*Default: None*
- **<batch\_size>**, *int, optional field*, the number of features to take in each mini batch.  
  
*Default: 3*
- **<verbose>**, *boolean, optional field*, Degree of verbosity of the printed output.  
*Default: False*
- **<shuffle>**, *boolean, optional field*, whether to shuffle the data before splitting it in batches.  
*Default: True*
- **<n\_jobs>**, *integer, optional field*, Parameters (keyword arguments) and values for kernel passed as callable object. Ignored by other kernels.  
*Default: 3*
- **<metho>**, *string, optional field*, method to use, options are:
  - \* lars: uses the least angle regression method to solve the lasso problem (`linear_model.lars_path`),
  - \* cd: uses the coordinate descent method to compute the Lasso solution (`linear_model.Lasso`)
 Lars will be faster if the estimated components are sparse.  
*Default: lars*
- **<random\_state>**, *integer or Random State, optional field*, Pseudo number generator state used for random sampling.  
*Default: None*

### 15.5.8.3.2 Truncated singular value decomposition

Dimensionality reduction using truncated SVD (aka LSA). In order to use the *Truncated SVD*, the user needs to set the sub-node:

```
<SKLtype>decomposition|TruncatedSVD</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_components>**, *integer, optional field*, Desired dimensionality of output data. Must be strictly less than the number of features. The default value is useful for visualisation. For LSA, a value of 100 is recommended.

*Default: 2*

- **<algorithm>**, *string, optional field*, SVD solver to use:
  - Randomized: randomized algorithm
  - Arpack: ARPACK wrapper in.

*Default: Randomized*

- **<n\_iter>**, *float, optional field*, number of iterations randomized SVD solver. Not used by ARPACK.

*Default: 5*

- **<random\_state>**, *int or Random State, optional field*, Pseudo number generator state used for random sampling. If not given, the numpy.random singleton is used.

*Default: None*

- **<tol>**, *float, optional field*, Tolerance for ARPACK. 0 means machine precision. Ignored by randomized SVD solver.

*Default: 0.0*

### 15.5.8.3.3 Fast ICA

A fast algorithm for Independent Component Analysis. In order to use the *Fast ICA*, the user needs to set the sub-node:

**<SKLtype>**decomposition|FastICA**</SKLtype>**.

In addition to this XML node, several others are available:

- **<n\_components>**, *integer, optional field*, Number of components to use. If none is passed, all are used.

*Default: None*

- **<algorithm>**, *string, optional field*, algorithm used in FastICA:
  - parallel,
  - deflation.

*Default: parallel*

- **<fun>**, *string or function, optional field*, The functional form of the G function used in the approximation to neg-entropy. Could be either:
  - logcosh,
  - exp, or
  - cube.

One can also provide own function. It should return a tuple containing the value of the function, and of its derivative, in the point.  
*Default: logcosh*
- **<fun\_args>**, *dictionary, optional field*, Arguments to send to the functional form. If empty and if fun='logcosh', fun\_args will take value 'alpha' : 1.0.  
*Default: None*
- **<max\_iter>**, *float, optional field*, maximum number of iterations during fit.  
*Default: 200*
- **<tol>**, *float, optional field*, Tolerance on update at each iteration.  
*Default: 0.0001*
- **<w\_init>**, *None or an (n\_components, n\_components) ndarray, optional field*, The mixing matrix to be used to initialize the algorithm.  
*Default: None*
- **<randome\_state>**, *int or Random State, optional field*, Pseudo number generator state used for random sampling.  
*Default: None*

#### 15.5.8.4 Manifold learning

A manifold is a topological space that resembles a Euclidean space locally at each point. Manifold learning is an approach to non-linear dimensionality reduction. It assumes that the data of interest lie on an embedded non-linear manifold within the higher-dimensional space. If this manifold is of low dimension, data can be visualized in the low-dimensional space. Algorithms for this task are based on the idea that the dimensionality of many data sets is only artificially high.

##### 15.5.8.4.1 Isomap

Non-linear dimensionality reduction through Isometric Mapping (Isomap). In order to use the *Isometric Mapping*, the user needs to set the sub-node:

```
<SKLtype>manifold|Isomap</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_neighbors>**, *integer, optional field*, Number of neighbors to consider for each point.  
*Default: 5*
- **<n\_components>**, *integer, optional field*, Number of coordinates to manifold.  
*Default: 2*
- **<eigen\_solver>**, *string, optional field*, eigen solver to use:
  - auto: Attempt to choose the most efficient solver for the given problem,
  - arpack: Use Arnoldi decomposition to find the eigenvalues and eigenvectors
  - dense: Use a direct solver (i.e. LAPACK) for the eigenvalue decomposition

*Default: auto*

- **<tol>**, *float, optional field*, Convergence tolerance passed to arpack or lobpcg. not used if eigen\_solver is 'dense'.  
*Default: 0.0*
- **<max\_iter>**, *float, optional field*, Maximum number of iterations for the arpack solver. not used if eigen\_solver == 'dense'.  
*Default: None*
- **<path\_method>**, *string, optional field*, Method to use in finding shortest path. Could be either:
  - Auto: attempt to choose the best algorithm
  - FW: Floyd-Warshall algorithm
  - D: Dijkstra algorithm with Fibonacci Heaps

*Default: auto*

- **<neighbors\_algorithm>**, *string, optional field*, Algorithm to use for nearest neighbors search, passed to neighbors.NearestNeighbors instance.
  - auto,
  - brute
  - kd\_tree
  - ball\_tree

Default: auto

**Example:**

```
<Simulation>
...
<Models>
...
  <PostProcessor name='PostProcessorName'
    subType='DataMiningPostProcessor'>
    <KDD lib='SciKitLearn'>
      <Features>input</Features>
      <SKLtype>manifold|Isomap</SKLtype>
      <n_neighbors>5</n_neighbors>
      <n_components>3</n_components>
      <eigen_solver>arpack</eigen_solver>
      <neighbors_algorithm>kd_tree</neighbors_algorithm>
    </KDD>
  </PostProcessor>
...
<Models>
...
</Simulation>
```

**15.5.8.4.2 Locally Linear Embedding**

In order to use the *Locally Linear Embedding*, the user needs to set the sub-node:

```
<SKLtype>manifold|LocallyLinearEmbedding</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_neighbors>**, *integer, optional field*, Number of neighbors to consider for each point.  
Default: 5
- **<n\_components>**, *integer, optional field*, Number of coordinates to manifold.  
Default: 2
- **<reg>**, *float, optional field*, regularization constant, multiplies the trace of the local covariance matrix of the distances.  
Default: 0.01
- **<eigen\_solver>**, *string, optional field*, eigen solver to use:

- auto: Attempt to choose the most efficient solver for the given problem,
- arpack: use arnoldi iteration in shift-invert mode.
- dense: use standard dense matrix operations for the eigenvalue

*Default: auto*

- **<tol>**, *float, optional field*, Convergence tolerance passed to arpack. not used if eigen\_solver is 'dense'.

*Default: 1E-06*

- **<max\_iter>**, *int, optional field*, Maximum number of iterations for the arpack solver. not used if eigen\_solver == 'dense'.

*Default: 100*

- **<method>**, *string, optional field*, Method to use. Could be either:

- Standard: use the standard locally linear embedding algorithm
- hessian: use the Hessian eigenmap method
- itsa: use local tangent space alignment algorithm

*Default: standard*

- **<hessian\_tol>**, *float, optional field*, Tolerance for Hessian eigenmapping method. Only used if method == 'hessian'

*Default: 0.0001*

- **<modified\_tol>**, *float, optional field*, Tolerance for modified LLE method. Only used if method == 'modified'

*Default: 0.0001*

- **<neighbors\_algorithm>**, *string, optional field*, Algorithm to use for nearest neighbors search, passed to neighbors.NearestNeighbors instance.

- auto,
- brute
- kd\_tree
- ball\_tree

*Default: auto*

- **<random\_state>**, *int or numpy random state, optional field*, the generator or seed used to determine the starting vector for arpack iterations.

*Default: None*



#### 15.5.8.4.3 Spectral Embedding

Spectral embedding for non-linear dimensionality reduction, it forms an affinity matrix given by the specified function and applies spectral decomposition to the corresponding graph laplacian. The resulting transformation is given by the value of the eigenvectors for each data point. In order to use the *Spectral Embedding*, the user needs to set the sub-node:

```
<SKLtype>manifold|SpectralEmbedding</SKLtype>.
```

In addition to this XML node, several others are available:

- **<n\_components>**, *integer, optional field*, the dimension of projected sub-space.  
*Default: 2*
- **<eigen\_solver>**, *string, optional field*, the eigen value decomposition strategy to use:
  - none,
  - arpack.
  - lobpcg,
  - amg

*Default: none*

- **<random\_state>**, *integer or numpy random state, optional field*, A pseudo random number generator used for the initialization of the lobpcg eigen vectors decomposition when `eigen_solver == 'amg'`.  
*Default: None*
- **<affinity>**, *string or callable, optional field*, How to construct the affinity matrix:
  - *nearest\_neighbors* : construct affinity matrix by knn graph
  - *rbf* : construct affinity matrix by rbf kernel
  - *precomputed* : interpret X as precomputed affinity matrix
  - *callable* : use passed in function as affinity the function takes in data matrix (n\_samples, n\_features) and return affinity matrix (n\_samples, n\_samples).

*Default: nearest\_neighbor*

- **<gamma>**, *float, optional field*, Kernel coefficient for rbf kernel.  
*Default: None*
- **<n\_neighbors>**, *int, optional field*, Number of nearest neighbors for nearest\_neighbors graph building.  
*Default: None*

#### 15.5.8.4.4 Multi-dimensional Scaling (MDS)

In order to use the *Multi Dimensional Scaling*, the user needs to set the sub-node:

```
<SKLtype>manifold|MDS</SKLtype>.
```

In addition to this XML node, several others are available:

- **<metric>**, *boolean, optional field*, compute metric or nonmetric SMACOF (Scaling by Majorizing a Complicated Function) algorithm  
*Default: True*
- **<n\_components>**, *integer, optional field*, number of dimension in which to immerse the similarities overridden if initial array is provided.  
*Default: 2*
- **<n\_init>**, *integer, optional field*, Number of time the smacof algorithm will be run with different initialisation. The final results will be the best output of the n\_init consecutive runs in terms of stress.  
*Default: 4*
- **<max\_iter>**, *integer, optional field*, Maximum number of iterations of the SMACOF algorithm for a single run  
*Default: 300*
- **<verbose>**, *integer, optional field*, level of verbosity  
*Default: 0*
- **<eps>**, *float, optional field*, relative tolerance with respect to stress to declare converge  
*Default: 1E-06*
- **<n\_jobs>**, *integer, optional field*, The number of jobs to use for the computation. This works by breaking down the pairwise matrix into n\_jobs even slices and computing them in parallel. If -1 all CPUs are used. If 1 is given, no parallel computing code is used at all, which is useful for debugging. For n\_jobs below -1, (n\_cpus + 1 + n\_jobs) are used. Thus for n\_jobs = -2, all CPUs but one are used.  
*Default: 1*
- **<random\_state>**, *<integer or numpy random state, optional field>*, The generator used to initialize the centers. If an integer is given, it fixes the seed. Defaults to the global numpy random number generator.  
*Default: None*
- **<dissimilarity>**, *string, optional field*, Which dissimilarity measure to use. Supported are 'euclidean' and 'precomputed'.  
*Default: euclidean*

## 15.5.9 Interfaced

The **Interfaced** post-processor is a Post-Processor that allows the user to create its own Post-Processor. While the External Post-Processor (see Section 15.5.6) allows the user to create case-dependent Post-Processors, with this new class the user can create new general purpose Post-Processors.

In order to use the *Interfaced* PP, the user needs to set the `subType` of a `<PostProcessor>` node:

```
<PostProcessor subType=' Interfaced' />.
```

Several sub-nodes are available:

- `<method>`, *comma separated string, required field*, lists the method names of a method that will be computed (each returning a post-processing value). All available methods need to be included in the “/raven/framework/PostProcessorFunctions/” folder

### Example:

```
<Simulation>
...
<Models>
...
  <PostProcessor name="example"
    subType=' InterfacedPostProcessor ' verbosity=' debug' >
    <method>testInterfacedPP</method>
    <!--Here, the xml nodes required by the chosen method
      have to be
      included.
    -->
  </PostProcessor>
...
</Models>
...
</Simulation>
```

All the **Interfaced** post-processors need to be contained in the “/raven/framework/PostProcessorFunctions/” folder. In fact, once the **Interfaced** post-processor is defined in the RAVEN input file, RAVEN search that the method of the post-processor is located in such folder.

The class specified in the **Interfaced** post-processor has to inherit the `PostProcessorInterface`-Base class and the user must specify this set of methods:

- **initialize:** in this method, the internal parameters of the post-processor are initialized. Mandatory variables that needs to be specified are the following:
  - `self.inputFormat`: type of `dataObject` expected in input
  - `self.outputFormat`: type of `dataObject` generated in output
- **readMoreXML:** this method is in charge of reading the `PostProcessor` xml node, parse it and fill the `PostProcessor` internal variables.
- **run:** this method performs the desired computation of the `dataObject`.

```

from PostProcessorInterfaceBaseClass import PostProcessorInterfaceBase
class testInterfacedPP (PostProcessorInterfaceBase) :
    def initialize (self)
    def readMoreXML (self, xmlNode)
    def run (self, inputDic)

```

### 15.5.9.1 Data Format

The user is not allowed to modify directly the `DataObjects`, however the content of the `DataObjects` is available in the form of a python dictionary. Both the dictionary give in input and the one generated in the output of the `PostProcessor` are structured as follows:

```
inputDict = {'data': {}, 'metadata': {}}
```

where:

```
inputDict['data'] = {'input': {}, 'output': {}}
```

In the input dictionary, each input variable is listed as a dictionary that contains a numpy array with its own values as shown below for a simplified example

```
inputDict['input'] = {'inputVar1': array([ 1., 2., 3.]),
                    'inputVar2': array([4., 5., 6.])}
```

Similarly, if the `dataObject` is a `PointSet` then the output dictionary is structured as follows:

```
inputDict['output'] = {'outputVar1': array([ .1, .2, .3]),
                    'outputVar2': array([.4, .5, .6])}
```

Howevers, if the `dataObject` is a `HistorySet` then the output dictionary is structured as follows:

```
inputDict['output'] = {'hist1': {}, 'hist2': {}}
```

where

```
inputDict['output'][hist1] = {'time': array([ .1,.2,.3]),  
                              'outputVar1':array([ .4,.5,.6])}  
inputDict['output'][hist2] = {'time': array([ .1,.2,.3]),  
                              'outputVar1':array([ .14,.15,.16])}
```

## 16 Functions

The RAVEN code provides support for the usage of user-defined external functions. These functions are python modules, with a format that is automatically interpretable by the RAVEN framework. For example, users can define their own method to perform a particular post-processing activity and the code will be embedded and use the function as though it were an active part of the code itself. In this section, the XML input syntax and the format of the accepted functions are fully specified.

The specifications of an external function must be defined within the XML block **<External>**. This XML node requires the following attributes:

- **name**, *required string attribute*, user-defined name of this function. **Note:** As with other objects, this name can be used to refer to this specific entity from other input blocks in the XML.
- **file**, *required string attribute*, absolute or relative path specifying the code associated to this function. **Note:** If a relative path is specified, it must be relative with respect to where the user is running the instance of RAVEN.

In order to make the RAVEN code aware of the variables the user is going to manipulate/use in her/his own python function, the variables need to be specified in the **<External>** input block. The user needs to input, within this block, only the variables directly used by the external code and not the local variables that the user does not want, for example, those stored in a RAVEN internal object. These variables are named within consecutive **<variable>** XML nodes:

- **<variable>**, *string, required parameter*, in the body of this XML node, the user needs to specify the name of the variable. This variable needs to match a variable used/defined in the external python function.

When the external function variables are defined, at runtime, RAVEN initializes them and keeps track of their values during the simulation. Each variable defined in the **<External>** block is available in the function as a python `self.` member. In the following, an example of a user-defined external function is reported (a python module and its related XML input specifications).

Example Python Function:

```
import numpy as np
def residuumSign(self):
    if self.var1 < self.var2 :
        return 1
    else:
        return -1
```

Example XML Input:

```
...
<Functions>
  ...
  <External name='whatever' file='path_to_python_file'>
    ...
    <variable>var1</variable>
    <variable>var2</variable>
    ...
  </External>
  ...
</Functions>
...
</Simulation>
```

## 17 Steps

The core of the RAVEN calculation flow is the **Step** system. The **Step** is in charge of assembling different entities in RAVEN (e.g. Samplers, Models, Databases, etc.) in order to perform a task defined by the kind of step being used. A sequence of different **Steps** represents the calculation flow.

Before analyzing each **Step** type, it is worth to 1) explain how a general **Step** entity is organized, and 2) introduce the concept of step “role”. In the following example, a general example of a **Step** is shown below:

```
<Simulation>
...
<Steps>
...
  <WhateverStepType name='aName' >
    <Role1 class='aMainClassType'
      type='aSubType' >userDefinedName1</Role1>
    <Role2 class='aMainClassType'
      type='aSubType' >userDefinedName2</Role2>
    <Role3 class='aMainClassType'
      type='aSubType' >userDefinedName3</Role3>
    <Role4 class='aMainClassType'
      type='aSubType' >userDefinedName4</Role4>
  </WhateverStepType>
...
</Steps>
...
</Simulation>
```

As shown above each **Step** consists of a list of entities organized into “Roles.” Each role represents a behavior the entity (object) will assume during the evaluation of the **Step**. In RAVEN, several different roles are available:

- **Input** represents the input of the **Step**. The allowable input objects depend on the type of **Model** in the **Step**.
- **Output** defines where to collect the results of an action performed by the **Model**. It is generally one of the following types: **DataObjects**, **Databases**, or **OutputStreamManager**.
- **Model** represents a physical or mathematical system or behavior. The object used in this role defines the allowable types of **Inputs** and **Outputs** usable in this step.



- **Sampler** defines the sampling strategy to be used to probe the model. It is worth to mention that, when a sampling strategy is employed, the “variables” defined in the `<variable>` blocks are going to be directly placed in the **Output** objects of type **DataObjects** and **Databases**).
- **Function** is an extremely important role. It introduces the capability to perform pre or post processing of Model **Inputs** and **Outputs**. Its specific behavior depends on the **Step** is using it.
- **ROM** defines an acceleration Reduced Order Model to use for a **Step**.
- **SolutionExport** represents the container of the eventual output of a Sampler. For the moment, it is only used when a **Step** is employing the search of the Limit Surface (LS), through the class of Adaptive **Samplers**). In this case, it contains the coordinates of the LS in the input space.

Depending on the **Step** type, different combinations of these roles can be used. For this reason, it is important to analyze each **Step** type in details.

The available steps are the following

- SingleRun (see Section 17.1)
- MultiRun(see Section 17.2)
- IOStep(see Section 17.3)
- RomTrainer(see Section 17.4)
- PostProcess(see Section 17.5)

## 17.1 SingleRun

The **SingleRun** is the simplest step the user can use to assemble a calculation flow: perform a single action of a **Model**. For example, it can be used to run a single job (Code Model) and collect the outcome(s) in a “**DataObjects**” object of type **Point** or **History** (see Section 12 for more details on available data representations).

The specifications of this Step must be defined within a `<SingleRun>` XML block. This XML node has the following definable attributes:

- **name**, *required string attribute*, user-defined name of this **Step**. **Note:** This name is used to reference this specific entity in the `<RunInfo>` block, under the `<Sequence>` node. If the name of this **Step** is not listed in the `<Sequence>` block, its action is not going to be performed.

- **pauseAtEnd**, *optional boolean/string attribute (case insensitive)*, if True (True values = True, yes, y, t), the code will pause at the end of the step, waiting for a user signal to continue. This is used in case one or more of the **Outputs** are of type **OutputStreamManager**. For example, it can be used when an **OutputStreamManager** of type **Plot** is output to the screen. Thus, allowing the user to interact with the **Plot** (e.g. rotate the figure, change the scale, etc.).

*Default: False.*

In the **<SingleRun>** input block, the user needs to specify the objects needed for the different allowable roles. This step accepts the following roles:

- **<Input>**, *string, required parameter*, names an entity (defined elsewhere in the RAVEN input) that will be used as input for the model specified in this step. This XML node accepts the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main object's type used in the input. For example, 'Files', 'DataObjects', 'Databases', etc.
  - **type**, *required string attribute*, the actual entity type. This attribute needs to specify the object type within the main object class. For example, if the **class** attribute is 'DataObjects', the **type** attribute might be 'PointSet'. **Note:** The class 'Files' has no type (i.e. **type**='').

**Note:** The **class** and, consequently, the **type** usable for this role depends on the particular **<Model>** being used. In addition, the user can specify as many **<Input>** nodes as needed.

- **<Model>**, *string, required parameter*, names an entity defined elsewhere in the input file to be used as a model for this step. This XML node accepts the following attributes:
  - **class**, *required string attribute*, main object class type. For this role, only 'Models' can be used.
  - **type**, *required string attribute*, the actual entity type. This attribute needs to specify the object type within the Models object class. For example, the **type** attribute might be 'Code', 'ROM', etc.
- **<Output>**, *string, required parameter* names an entity defined elsewhere in the input to use as the output for the **Model**. This XML node recognizes the following attributes:
  - **class**, *required string attribute*, main object class type. For this role, only 'DataObjects', 'Databases', and 'OutputStreamManager' can be used.
  - **type**, *required string attribute*, the actual entity type. This attribute needs to specify the object type within the main object class. For example, if the **class** attribute is 'DataObjects', the **type** attribute might be 'PointSet'.

**Note:** The number of `<Output>` nodes is unlimited.

Example:

```
<Steps>
...
<SingleRun name='StepName' pauseAtEnd='false'>
  <Input class='Files'
    type=''>anInputFile.i</Input>
  <Input class='Files' type=''>aFile</Input>
  <Model class='Models' type='Code'>aCode</Model>
  <Output class='Databases' type='HDF5'>aDatabase</Output>
  <Output class='DataObjects'
    type='History'>aData</Output>
</SingleRun>
...
</Steps>
```

## 17.2 MultiRun

The **MultiRun** step allows the user to assemble the calculation flow of an analysis that requires multiple “runs” of the same model. This step is used, for example, when the input (space) of the model needs to be perturbed by a particular sampling strategy.

The specifications of this type of step must be defined within a `<MultiRun>` XML block. This XML node recognizes the following list of attributes:

- **name**, *required string attribute*, user-defined name of this Step. **Note:** As with other objects, this name is used to reference this specific entity in the `<RunInfo>` block, under the `<Sequence>` node. If the name of this **Step** is not listed in the `<Sequence>` block, its action is not going to be performed.
- **pauseAtEnd**, *optional boolean/string attribute*, if True (True values = True, yes, y, t), the code will pause at the end of the step, waiting for a user signal to continue. This is used in case one or more of the **Outputs** are of type **OutputStreamManager**. For example, it can be used when an **OutputStreamManager** of type **Plot** is output to the screen. Thus, allowing the user to interact with the **Plot** (e.g. rotate the figure, change the scale, etc.).
- **sleepTime**, *optional float attribute*, in this attribute the user can specify the waiting time (seconds) between two subsequent inquiries of the status of the submitted job (i.e. check if a run has finished).

*Default: 0.05.*

In the `<MultiRun>` input block, the user needs to specify the objects that need to be used for the different allowable roles. This step accepts the following roles:

- **<Input>**, *string, required parameter*, names an entity to be used as input for the model specified in this step. This XML node accepts the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main object's type used in the input. For example, 'Files', 'DataObjects', 'Databases', etc.
  - **type**, *required string attribute*, the actual entity type. This attribute specifies the object type within the main object class. For example, if the **class** attribute is 'DataObjects', the **type** attribute might be 'PointSet'. **Note:** The class 'Files' has no type (i.e. **type**='').

**Note:** The **class** and, consequently, the **type** usable for this role depend on the particular **<Model>** being used. The user can specify as many **<Input>** nodes as needed.

- **<Model>**, *string, required parameter* names an entity defined elsewhere in the input that will be used as the model for this step. This XML node recognizes the following attributes:
  - **class**, *required string attribute*, main object class type. For this role, only 'Models' can be used.
  - **type**, *required string attribute*, the actual entity type. This attribute needs to specify the object type within the Models object class. For example, the **type** attribute might be 'Code', 'ROM', etc.
- **<Sampler>**, *string, required parameter* names an entity defined elsewhere in the input file to be used as a sampler. As mentioned in Section 11, the **Sampler** is in charge of defining the strategy to characterize the input space. This XML node recognizes the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main object's type used. Only 'Samplers' can be used for this role.
  - **type**, *required string attribute*, the actual entity type. This attribute needs to specify the object type within the Samplers object class. For example, the **type** attribute might be 'MonteCarlo', 'Adaptive', 'AdaptiveDET', etc. See Section 11 for all the different types currently supported.
- **<SolutionExport>**, *string, optional parameter* identifies an entity to be used for exporting key information coming from the **Sampler** object during the simulation. This XML node accepts the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main object's type used in the input. For this role, only 'DataObjects' can be used.
  - **type**, *required string attribute*, the actual entity type. This attribute needs to specify the object type within the DataObjects object class. For example, the **type** attribute might be 'Code', 'ROM', etc.

**Note:** Whether or not it is possible to export the **Sampler** solution depends on the **type**. Currently, only the Samplers in the '**Adaptive**' category can export their solution into a **<SolutionExport>** entity. The **<Outputs>** node in the `DataObjects` needs to contain the goal **<Function>** name. For example, if **<Sampler>** is of type '**Adaptive**', the **<SolutionExport>** needs to be of type '**PointSet**' and it will contain the coordinates, in the input space, that belong to the "Limit Surface."

- **<Output>**, *string, required parameter* identifies an entity to be used as output for this step. This XML node recognizes the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main object's type used in the input. For this role, only '**DataObjects**', '**Databases**', and '**OutputStreamManager**' may be used.
  - **type**, *required string attribute*, the actual entity type. This attribute specifies the object type within the main object class. For example, if the **class** attribute is '**DataObjects**', the **type** attribute might be '**PointSet**'.

**Note:** The number of **<Output>** nodes is unlimited.

Example:

```
<Steps>
...
<MultiRun name='StepName1' pauseAtEnd='False' sleepTime='0.01'>
  <Input class='Files' type=''>anInputFile.i</Input>
  <Input class='Files' type=''>aFile</Input>
  <Sampler class='Samplers' type='Grid'>aGridName</Sampler>
  <Model class='Models' type='Code'>aCode</Model>
  <Output class='Databases' type='HDF5'>aDatabase</Output>
  <Output class='DataObjects' type='History'>aData</Output>
</MultiRun >
<MultiRun name='StepName2' pauseAtEnd='True' sleepTime='0.02'>
  <Input class='Files' type=''>anInputFile.i</Input>
  <Input class='Files' type=''>aFile</Input>
  <Sampler class='Samplers' type='Adaptive'>anAS</Sampler>
  <Model class='Models' type='Code'>aCode</Model>
  <Output class='Databases' type='HDF5'>aDatabase</Output>
  <Output class='DataObjects' type='History'>aData</Output>
  <SolutionExport class='DataObjects' type='PointSet'>
    aTPS
  </SolutionExport>
</MultiRun>
...
</Steps>
```

## 17.3 IOStep

As the name suggests, the **IOStep** is the step where the user can perform input/output operations among the different I/O entities available in RAVEN. This step type is used to:

- construct/update a *Database* from a *DataObjects* object, and vice versa;
- construct/update a *DataObject* from a *CSV* file contained in a directory;
- construct/update a *Database* or a *DataObjects* object from *CSV* files contained in a directory;
- stream the content of a *Database* or a *DataObjects* out through an **OutputStream** object (see section 14);
- store/retrieve a *ROM* to/from an external *File* using Pickle module of Python.

This last function can be used to create and store mathematical model of fast solution trained to predict a response of interest of a physical system. This model can be recovered in other simulations or used to evaluate the response of a physical system in a Python program by the implementing of the Pickle module. The specifications of this type of step must be defined within an **<IOStep>** XML block. This XML node can accept the following attributes:

- **name**, *required string attribute*, user-defined name of this Step. **Note:** As for the other objects, this is the name that can be used to refer to this specific entity in the **<RunInfo>** block, under the **<Sequence>** node.
- **pauseAtEnd**, *optional boolean/string attribute (case insensitive)*, if True (True values = True, yes, y, t), the code will pause at the end of the step, waiting for a user signal to continue. This is used in case one or more of the **Outputs** are of type **OutputStreamManager**. For example, it can be used when an **OutputStreamManager** of type **Plot** is output to the screen. Thus, allowing the user to interact with the **Plot** (e.g. rotate the figure, change the scale, etc.).

*Default: False.*

In the **<IOStep>** input block, the user specifies the objects that need to be used for the different allowable roles. This step accepts the following roles:

- **<Input>**, *string, required parameter*, names an entity that is going to be used as a source (input) from which the information needs to be extracted. This XML node recognizes the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main object's type used in the input. As already mentioned, the allowable main classes are '**DataObjects**', '**Databases**', '**Models**' and '**Files**'.



- **type**, *required string attribute*, the actual entity type. This attribute needs to specify the object type within the main object class. For example, if the **class** attribute is 'DataObjects', the **type** attribute might be 'PointSet'. If the **class** attribute is 'Models', the **type** attribute must be 'ROM' and if the **class** attribute is 'Files', the **type** attribute must be ' '.
- **<Output>**, *string, required parameter* names an entity to be used as the target (output) where the information extracted in the input will be stored. This XML node needs to contain the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main object's type used in the input. The allowable main classes are 'DataObjects', 'Databases', 'OutputStreamManager', 'Models' and 'Files'.
  - **type**, *required string attribute*, the actual entity type. This attribute specifies the object type within the main object class. For example, if the **class** attribute is 'OutputStreamManager', the **type** attribute might be 'Plot'.

This step acts as a “transfer network” among the different RAVEN storing (or streaming) objects. The number of **<Input>** and **<Output>** nodes is unlimited, but should match. This step assumes a 1-to-1 mapping (e.g. first **<Input>** is going to be used for the first **<Output>**, etc.).

**Note:** This 1-to-1 mapping is not present when **<Output>** nodes are of **class** 'OutputStreamManager', since **OutputStreamManager** objects are already linked to a Data object in the relative RAVEN input block.

In this case, the user needs to provide in the **<Input>** nodes, all of the “DataObjects” objects linked to the **OutputStreamManager** objects (see the example below):

```

<Steps>
...
<IOStep name='OutputStreamStep'>
  <Input class='DataObjects'
    type='HistorySet'>aHistorySet</Input>
  <Input class='DataObjects' type='PointSet'>aTPS</Input>
  <Output class='OutputStreamManager' type='Plot'>plot_hist
  </Output>
  <Output class='OutputStreamManager' type='Print'>print_hist
  </Output>
  <Output class='OutputStreamManager' type='Print'>print_tps
  </Output>
  <Output class='OutputStreamManager' type='Print'>print_tp
  </Output>
</IOStep>
<IOStep name='PushDataObjectsIntoDatabase'>

```

```

    <Input class='DataObjects'
        type='HistorySet'>aHistorySet</Input>
    <Input class='DataObjects' type='PointSet'>aTPS</Input>
    <Output class='Databases' type='HDF5'>aDatabase</Output>
    <Output class='Databases' type='HDF5'>aDatabase</Output>
</IOStep>
<IOStep name='ConstructDataObjectsFromCSV'>
    <Input class='Files' type=''>aCSVFile</Input>
    <Output class='DataObjects' type='PointSet'>aPS</Output>
</IOStep>
<IOStep name='ConstructDataObjectsFromDatabase'>
    <Input class='Databases' type='HDF5'>aDatabase</Input>
    <Input class='Databases' type='HDF5'>aDatabase</Input>
    <Output class='DataObjects'
        type='HistorySet'>aHistorySet</Output>
    <Output class='DataObjects' type='PointSet'>aTPS</Output>
</IOStep>
<IOStep name='PushROMIntoFile'>
    <Input class='Models' type='ROM'>aROM</Input>
    <Output class='Files' type=''>aFile</Output>
</IOStep>
<IOStep name='ImportROMFromFile'>
    <Input class='Files' type=''>zFile</Input>
    <Output class='Models' type='ROM'>zROM</Output>
</IOStep>
    ...
</Steps>

```

Example of how to use a ROM exported by RAVEN using the IOStep, where (x1,x2) are the input variable names for which the ROM has been trained: Example Python Function:

```

import os
from glob import glob
import inspect
import utils
import numpy as np
for dirr,_,_ in os.walk("path_to_framework"):
    utils.add_path(dirr)
import pickle
fileobj = open('zFile','rb')
unpickledObj = pickle.load(fileobj)
dictionary={"x1":np.atleast_1d(Value1),"x2":np.atleast_1d(Value2)}
eval=unpickledObj.evaluate(dictionary)

```



```
print str(eval)
fileobj.close()
```

## 17.4 RomTrainer

The **RomTrainer** step type performs the training of a Reduced Order Model. The specifications of this step must be defined within a **<RomTrainer>** block. This XML node accepts the attributes:

- **name**, *required string attribute*, user-defined name of this step. **Note:** As for the other objects, this is the name that can be used to refer to this specific entity in the **<RunInfo>** block under **<Sequence>**.

In the **<RomTrainer>** input block, the user will specify the objects needed for the different allowable roles. This step accepts the following roles:

- **<Input>**, *string, required parameter* names an entity to be used as a source (input) from which the information needs to be extracted. This XML node accepts the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main object's type used in the input. The only allowable main class is **'DataObjects'**.
  - **type**, *required string attribute*, the actual entity type. This attribute specifies the object type within the main object class. For example, the **type** attribute might be **'PointSet'**. **Note:** Since the ROMs currently present in RAVEN are not time-dependent, the only allowable types are **'Point'** and **'PointSet'**.
- **<Output>**, *string, required parameter*, names a ROM entity that is going to be trained. This XML node recognizes the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main objects type used in the input. The only allowable main class is **'Models'**.
  - **type**, *required string attribute*, the actual entity type. This attribute needs to specify the object type within the main object class. The only type accepted here is, currently, **'ROM'**.

Example:

```

<Steps>
  ...
  <RomTrainer name='aStepName' >
    <Input class='DataObjects' type='PointSet' >aTPS</Input>
    <Output class='Models' type='ROM' >aROM</Output>
  </RomTrainer>
  ...
</Steps>

```

## 17.5 PostProcess

The **PostProcess** step is used to post-process data or manipulate RAVEN entities. It is aimed at performing a single action that is employed by a **Model** of type **PostProcessor**.

The specifications of this type of step is defined within a **<PostProcess>** XML block. This XML node specifies the following attributes:

- **name**, *required string attribute*, user-defined name of this Step. **Note:** As for the other objects, this is the name that is used to refer to this specific entity in the **<RunInfo>** block under the **<Sequence>** node.
- **pauseAtEnd**, *optional boolean/string attribute (case insensitive)*, if True (True values = True, yes, y, t), the code will pause at the end of the step, waiting for a user signal to continue. This is used in case one or more of the **Outputs** are of type **OutStreamManager**. For example, it can be used when an **OutStreamManager** of type **Plot** is output to the screen. Thus, allowing the user to interact with the **Plot** (e.g. rotate the figure, change the scale, etc.).

*Default: False.*

In the **<PostProcess>** input block, the user needs to specify the objects needed for the different allowable roles. This step accepts the following roles:

- **<Input>**, *string, required parameter*, names an entity to be used as input for the model specified in this step. This XML node accepts the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main object's type used in the input. For example, '**Files**', '**DataObjects**', '**Databases**', etc.
  - **type**, *required string attribute*, the actual entity type. This attribute specifies the object type within the main object class. For example, if the **class** attribute is '**DataObjects**', the **type** attribute might be '**PointSet**'. **Note:** The class '**Files**' has no type (i.e. **type**='').

**Note:** The **class** and, consequently, the **type** usable for this role depends on the particular type of **PostProcessor** being used. In addition, the user can specify as many **<Input>** nodes as needed by the model.

- **<Model>**, *string, required parameter*, names an entity to be used as a model for this step. This XML node recognizes the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main object's type used in the input. For this role, only '**Models**' can be used.
  - **type**, *required string attribute*, the actual entity type. This attribute needs to specify the object type within the '**Models**' object class. The only type accepted here is '**PostProcessor**'.
- **<Output>**, *string, required/optional parameter*, names an entity to be used as output for the PostProcessor. The necessity of this XML block and the types of entities that can be used as output depend on the type of **PostProcessor** that has been used as a **Model** (see section 15.5). This XML node specifies the following attributes:
  - **class**, *required string attribute*, main object class type. This string corresponds to the tag of the main object's type used in the input.
  - **type**, *required string attribute*, the actual entity type. This attribute specifies the object type within the main object class. For example, if the **class** attribute is '**DataObjects**', the **type** attribute might be '**PointSet**'.

**Note:** The number of **<Output>** nodes is unlimited.

Example:

```
<Steps>
...
<PostProcess name='PP1'>
  <Input class='DataObjects' type='PointSet' >aData</Input>
  <Model class='Models' type='PostProcessor'>aPP</Model>
  <Output class='Files' type=''>anOutputFile</Output>
</PostProcess>
...
</Steps>
```

## 18 Existing Interfaces

### 18.1 Generic Interface

The GenericCode interface is meant to handle a wide variety of generic codes that take straightforward input files and produce output CSV files. There are some limitations for this interface. If a code:

- accepts a keyword-based input file with no cross-dependent inputs,
- has no more than one filetype extension per command line flag,
- and returns a CSV with the input parameters and output parameters,

the GenericCode interface should cover the code for RAVEN.

If a code contains cross-dependent data, the generic interface is not able to edit the correct values. For example, if a geometry-building script specifies `inner_radius`, `outer_radius`, and `thickness`, the generic interface cannot calculate the thickness given the outer and inner radius, or vice versa.

An example of the code interface is shown here. The input parameters are read from the input files `gen.one` and `gen.two` respectively. The code is run using `python`, so that is part of the `<prepend>` node. The command line entry to normally run the code is

```
python poly_inp.py -i gen.one -a gen.two -o myOut
```

and produces the output `myOut.csv`.

Example:

```
<Code name="poly" subType="GenericCode">
  <executable>GenericInterface/poly_inp.py</executable>
  <inputExtensions>.one,.two</inputExtensions>
  <clargs type='prepend' arg='python' />
  <clargs type='input' arg='-i' extension='.one' />
  <clargs type='input' arg='-a' extension='.two' />
  <clargs type='output' arg='-o' />
  <prepend>python</prepend>
</Code>
```

If a code doesn't accept necessary Raven-editable auxiliary input files or output filenames through the command line, the GenericCode interface can also edit the input files and insert the filenames there. For example, in the previous example, say instead of `-a gen.two` and `-o myOut` in the command line, `gen.one` has the following lines:

```
...
auxfile = gen.two
case = myOut
...
```

Then, our example XML for the code would be

Example:

```
<Code name="poly" subType="GenericCode">
  <executable>GenericInterface/poly_inp.py</executable>
  <inputExtentions>.one, .two</inputExtentions>
  <clargs type='prepend' arg='python' />
  <clargs type='input' arg='-i' extension='.one' />
  <fileargs type='input' arg='two' extension='.two' />
  <fileargs type='output' arg='out' />
  <prepend>python</prepend>
</Code>
```

and the corresponding template input file lines would be changed to read

```
...
auxfile = $RAVEN-two$
case = $RAVEN-out$
...
```

In addition, the “wild-cards” above can contain two special and optional symbols:

- `:`, that defines an eventual default value;
- `|`, that defines the format of the value. The Generic Interface currently supports the following formatting options (\* in the examples means blank space):
  - **plain integer**, in this case the value that is going to be replaced by the Generic Interface, will be left-justified with a string length equal to the integer value specified here (e.g. “| 6”, the value is left-justified with a string length of 6);
  - **d**, signed integer decimal, the value is going to be formatted as an integer (e.g. if the value is 9 and the format “| 10d”, the replaced value will be formatted as follows: “\*\*\*\*\*9”);
  - **e**, floating point exponential format (lowercase), the value is going to be formatted as a float in scientific notation (e.g. if the value is 9.1234 and the format “| 10.3e”, the replaced value will be formatted as follows: “\*9.123e+00” );

- **E**, floating point exponential format (uppercase), the value is going to be formatted as a float in scientific notation (e.g. if the value is 9.1234 and the format “| 10.3E”, the replaced value will be formatted as follows: “\*9.123E+00” );
- **f or F**, floating point decimal format, the value is going to be formatted as a float in decimal notation (e.g. if the value is 9.1234 and the format “| 10.3f”, the replaced value will be formatted as follows: “\*\*\*\*\*9.123” );
- **g**, floating point format. Uses lowercase exponential format if exponent is less than -4 or not less than precision, decimal format otherwise (e.g. if the value is 9.1234 and the format “| 10.3g”, the replaced value will be formatted as follows: “\*\*\*\*\*9.12” );
- **G**, floating point format. Uses uppercase exponential format if exponent is less than -4 or not less than precision, decimal format otherwise (e.g. if the value is 0.000009 and the format “| 10.3G”, the replaced value will be formatted as follows: “\*\*\*\*\*9E-06” ).

For example:

```
...
auxfile = $RAVEN-two:3$
case = $RAVEN-out:5|10$
...
```

Where,

- :, in case the variable “two” is not defined in the RAVEN XML input file, the Parser, will replace it with the value “3”.;
- |, the value that is going to be replaced by the Generic Interface, will be left- justified with a string length of “10”;

## 18.2 RELAP5 Interface

### 18.2.1 Sequence

In the **<Sequence>** section, the names of the steps declared in the **<Steps>** block should be specified. As an example, if we called the first multirun “Grid\_Sampler” and the second multirun “MC\_Sampler” in the sequence section we should see this:

```
<Sequence>Grid_Sampler,MC_Sampler</Sequence>
```

## 18.2.2 batchSize and mode

For the `<batchSize>` and `<mode>` sections please refer to the `<RunInfo>` block in the previous chapters.

## 18.2.3 RunInfo

After all of these blocks are filled out, a standard example RunInfo block may look like the example below:

```
<RunInfo>
  <WorkingDir>~/workingDir</WorkingDir>
  <Sequence>Grid_Sampler,MC_Sampler</Sequence>
  <batchSize>1</batchSize>
  <mode>mpi</mode>
  <expectedTime>1:00:00</expectedTime>
  <ParallelProcNumb>1</ParallelProcNumb>
</RunInfo>
```

## 18.2.4 Files

In the `<Files>` section, as specified before, all of the files needed for the code to run should be specified. In the case of RELAP5, the files typically needed are:

- RELAP5 Input file
- Table file or files that RELAP needs to run
- The Relap5 executable

Example:

```
<Files>
  <Input name='inputrelap.i' type=''>inputfilerelap.i</Input>
  <Input name='tpfh2o' type=''>tpfh2o</Input>
  <Input name='r5executable.x' type=''>r5executabler.x</Input>
  <Input name='X10.i' type=''>X10.i</Input>
</Files>
```

It is a good practice to put inside the working directory all of these files and also:

- the RAVEN input file
- the license for the executable of RELAP5

### 18.2.5 Models

For the `<Models>` block here is a standard example of how it would look when using RELAP5 as the external model:

```
<Models>
  <Code name='MyRELAP' subType='Relap5'>
    <executable>~/path_to_the_executable</executable>
  </Code>
</Models>
```

### 18.2.6 Distributions

The `<Distribution>` block defines the distributions that are going to be used for the sampling of the variables defined in the `<Samplers>` block. For all the possible distributions and all their possible inputs please see the chapter about Distributions (see 10). Here we give a general example of three different distributions:

```
<Distributions verbosity='debug'>
  <Triangular name='BPfailtime'>
    <apex>5.0</apex>
    <min>4.0</min>
    <max>6.0</max>
  </Triangular>
  <LogNormal name='BPrepairtime'>
    <mean>0.75</mean>
    <sigma>0.25</sigma>
  </LogNormal>
  <Uniform name='ScalFactPower'>
    <lowerBound>1.0</lowerBound>
    <upperBound>1.2</upperBound>
  </Uniform>
</Distributions>
```

It is good practice to name the distribution something similar to what kind of variable is going to be sampled, since there might be many variables with the same kind of distributions but different input parameters.



## 18.2.7 Samplers

In the `<Samplers>` block we want to define the variables that are going to be sampled. **Example:** We want to do the sampling of 3 variables:

- Battery Fail Time
- Battery Repair Time
- Scaling Factor Power Rate

We are going to sample these 3 variables using two different sampling methods: grid and MonteCarlo.

In RELAP5, the sampler reads the variable as, given the name, the first number is the card number and the second number is the word number. In this example we are sampling:

- For card 0000588 (trip) the word 6 (battery failure time)
- For card 0000575 (trip) the word 6 (battery repair time)
- For card 20210000 (reactor power) the word 4 (reactor scaling factor)

We proceed to do so for both the Grid sampling and the MonteCarlo sampling.

```
<Samplers verbosity='debug'>
  <Grid name='Grid_Sampler' >
    <variable name='0000588:6'>
      <distribution>BPfailtime</distribution>
      <grid type='value' construction='equal' lowerBound='0.0'
        steps='10'>2880</grid>
    </variable>
    <variable name='0000575:6'>
      <distribution>BPrepairtime</distribution>
      <grid type='value' construction='equal' lowerBound='0.0'
        steps='10'>2880</grid>
    </variable>
    <variable name='20210000:4'>
      <distribution>ScalFactPower</distribution>
      <grid type='value' construction='equal' lowerBound='1.0'
        steps='10'>0.02</grid>
    </variable>
  </Grid>
```

```

<MonteCarlo name='MC_Sampler' limit='1000'>
  <variable name='0000588:6'>
    <distribution>BPfailtime</distribution>
  </variable>
  <variable name='0000575:6'>
    <distribution>BPrepairtime</distribution>
  </variable>
  <variable name='20210000:4'>
    <distribution>ScalFactPower</distribution>
  </variable>
</MonteCarlo>
</Samplers>

```

It can be seen that each variable is connected with a proper distribution defined in the **<Distributions>** block (from the previous example). The following demonstrates how the input for the first variable is read.

We are sampling a variable situated in word 6 of the card 0000588 using a Grid sampling method. The distribution that this variable is following is a Triangular distribution (see section above). We are sampling this variable beginning from 0.0 in 10 *equal* steps of 2880.

### 18.2.8 Steps

For a RELAP interface, the **<MultiRun>** step type will most likely be used. First, the step needs to be named: this name will be one of the names used in the **<Sequence>** block. In our example, `Grid_Sampler` and `MC_Sampler`.

```

<MultiRun name='Grid_Sampler' verbosity='debug'>

```

With this step, we need to import all the files needed for the simulation:

- RELAP input file
- element tables – `tpfh2o`

```

<Input class='Files' type=''>inputrelap.i</Input>
<Input class='Files' type=''>tpfh2o</Input>

```

We then need to define which model will be used:

```

<Model class='Models' type='Code'>MyRELAP</Model>

```

We then need to specify which Sampler is used, and this can be done as follows:

```
<Sampler class='Samplers' type='Grid'>Grid_Sampler</Sampler>
```

And lastly, we need to specify what kind of output the user wants. For example the user might want to make a database (in RAVEN the database created is an HDF5 file). Here is a classical example:

```
<Output class='Databases' type='HDF5'>MC_out</Output>
```

Following is the example of two MultiRun steps which use different sampling methods (grid and Monte Carlo), and creating two different databases for each one:

```
<Steps verbosity='debug'>
  <MultiRun name='Grid_Sampler' verbosity='debug'>
    <Input class='Files' type='' >X10.i</Input>
    <Input class='Files' type='' >r5executable.x</Input>
    <Input class='Files' type='' >tpfh2o</Input>
    <Model class='Models' type='Code'>MyRELAP</Model>
    <Sampler class='Samplers' type='Grid'>Grid_Sampler</Sampler>
    <Output class='Databases' type='HDF5'>Grid_out</Output>
  </MultiRun>
  <MultiRun name='MC_Sampler' verbosity='debug'
    re-seeding='210491'>
    <Input class='Files' type='' >X10.i</Input>
    <Input class='Files' type=''
      >r5executable.x</Input>
    <Input class='Files' type='' >tpfh2o</Input>
    <Model class='Models' type='Code' >MyRELAP</Model>
    <Sampler class='Samplers'
      type='MonteCarlo'>MC_Sampler</Sampler>
    <Output class='Databases' type='HDF5' >MC_out</Output>
  </MultiRun>
</Steps>
```

### 18.2.9 Databases

As shown in the `<Steps>` block, the code is creating two database objects called `Grid_out` and `MC_out`. So the user needs to input the following:

```
<Databases>
  <HDF5 name="Grid_out"/>
  <HDF5 name="MC_out"/>
</Databases>
```

As listed before, this will create two databases. The files will have names corresponding to their **name** appended with the .h5 extension (i.e. `Grid_out.h5` and `MC_out.h5`).

## 18.3 RELAP7 Interface

This section covers the input specifications for running RELAP7 through RAVEN. It is important to notice that this short explanation assumes that the reader already knows how to use the control logic system in RELAP7. Since the presence of the control logic system in RELAP7, this code interface is different with respect to the others and uses some special keyword available in RAVEN (see the following).

### 18.3.1 Files

In the **<Files>** section, as specified before, all of the files needed for the code to run should be specified. In the case of RELAP7, the files typically needed are the following:

- RELAP7 Input file
- Control Logic file

Example:

```
<Files>
  <Input name='nat_circ.i' type=''>nat_circ.i</Input>
  <Input name='control_logic.py' type=''>control_logic.py</Input>
</Files>
```

The RAVEN/RELAP7 interface recognizes as RELAP7 inputs the files with the extensions “\*.i”, “\*.inp” and “\*.in”.

### 18.3.2 Models

For the **<Models>** block RELAP7 uses the RAVEN executable, since through this executable the stochastic environment gets activated (possibility to sample parameters directly in the control logic system) Here is a standard example of what can be used to use RELAP7 as the model:

```
<Models>
  <Code name='MyRAVEN'
    subType='RAVEN'><executable>~path/to/RAVEN-opt</executable></Code>
</Models>
```

### 18.3.3 Distributions

As for all the other codes interfaces the **<Distributions>** block needs to be specified in order to employ as sampling strategy (e.g. MonteCarlo, Stratified, etc.). In this block, the user specifies the distributions that need to be used. Once the user defines the distributions in this block, RAVEN activates the Distribution environment in the RAVEN/RELAP7 control logic system. The sampling of the parameters is then performed directly in the control logic input file.

For example, let's consider the sampling of a normal distribution for the primary pressure in RELAP7:

```
<Distributions>  
<Normal name="Prim_Pres">  
<mean>1000000</mean>  
<sigma>100<sigma/>  
</Normal>  
</Distributions>
```

In order to change a parameter (independently on the sampling strategy), the control logic input file should be modified as follows:

```
def initial_function(monitored, controlled, auxiliary)  
    print ("monitored",monitored,"controlled",  
          controlled,"auxiliary",auxiliary)  
  
    controlled.pressureInPressurizer =  
        distributions.Prim_Pres.getDistributionsRandom()  
    return
```

### 18.3.4 Samplers

In the **<Samplers>** block, all the variables that needs to be sampled must be specified. In case some of these variables are directly sampled in the Control Logic system, the **<variable>** needs to be replaced with **<Distribution>**. In this way, RAVEN is able to understand which variables needs to be directly modified through input file (i.e. modifying the original input file \*.i) and which variables are going to be “sampled” through the control logic system. For the example, we are performing Grid Sampling. The global initial pressure wasn't specified in the control logic so it is going to be specified using the node **<variable>**. The “pressureInPressurizer “ variable is instead sampled in the control logic system; for this reason, it is going to be specified using the node **<Distribution>**. For example,

```
<Samplers>
```

```

<Grid name="MC_samp">
  <samplerInit> <limit>500</limit> </samplerInit>
  <variable name="GlobalParams|global_init_P">
    <distribution>Prim_Pres</distribution>
    <grid construction="equal" steps="10" type="CDF">0.0
      1.0</grid>
  </variable>
  <Distribution name="pressureInPressurizer">
    <distribution>Prim_Pres</distribution>
    <grid construction="equal" steps="10" type="CDF">0.0
      1.0</grid>
  </Distribution>
</Grid>
</Samplers>

```

## 18.4 MooseBasedApp Interface

### 18.4.1 Files

In the **<Files>** section, as specified before, all of the files needed for the code to run should be specified. In the case of any MooseBasedApp, the files typically needed are the following:

- MooseBasedApp YAML input file
- Restart Files (if the calculation is instantiated from a restart point)

Example:

```

<Files>
  <Input name='mooseBasedApp.i' type=''>mooseBasedApp.i</Input>
  <Input name='0020_mesh.cpr' type=''>0020_mesh.cpr</Input>
  <Input name='0020.xdr.0000'>0020.xdr.0000</Input>
  <Input name='0020.rd-0'>0020.rd-0</Input>
</Files>

```

### 18.4.2 Models

In the **<Models>** block particular MooseBasedApp executable needs to be specified. Here is a standard example of what can be used to use with a typical MooseBasedApp (Bison) as the model:

```

<Models>
  <Code name='MyMooseBasedApp'
    subType='MooseBasedApp' ><executable>~path/to/Bison-opt</executable></Code>
</Models>

```

### 18.4.3 Distributions

The `<Distributions>` block defines the distributions that are going to be used for the sampling of the variables defined in the `<Samplers>` block. For all the possible distributions and all their possible inputs please see the chapter about Distributions (see 10). Here we give a general example of three different distributions:

```

<Distributions>
  <Normal name='ThermalConductivity1'>
    <mean>1</mean>
    <sigma>0.001</sigma>
    <lowerBound>0.5</lowerBound>
    <upperBound>1.5</upperBound>
  </Normal>
  <Normal name='SpecificHeat'>
    <mean>1</mean>
    <sigma>0.4</sigma>
    <lowerBound>0.5</lowerBound>
    <upperBound>1.5</upperBound>
  </Normal>
  <Triangular name='ThermalConductivity2'>
    <apex>1</apex>
    <min>0.1</min>
    <max>4</max>
  </Triangular>
</Distributions>

```

It is good practice to name the distribution something similar to what kind of variable is going to be sampled, since there might be many variables with the same kind of distributions but different input parameters.

## 18.4.4 Samplers

In the `<Samplers>` block we want to define the variables that are going to be sampled. **Example:** We want to do the sampling of 3 variables:

- Thermal Conductivity of the Fuel;
- Specific Heat Transfer Ratio of the Cladding;
- Thermal Conductivity of the Cladding.

We are going to sample these 3 variables using two different sampling methods: Grid and Monte-Carlo.

In order to perturb any MooseBasedApp, the user needs to specify the variables to be sampled indicating the path to the value separated with the symbol “|”. For example, if the variable that we want to perturb is specified in the input as follows:

```
[Materials]
...
[./heatStructure]
...
  thermal_conductivity = 1.0
...
[../]
...
[]
```

the variable name in the Sampler input block needs to be named as follows:

```
...
<Samplers>
  <aSampler name='aUserDefinedName' >
    <variable
      name='Materials|heatStructure|thermal_conductivity'>
      ...
    </variable>
  </aSampler>
</Samplers>
...
```

In this example, we proceed to do so for both the Grid sampling and the Monte-Carlo sampling.



```

<Samplers verbosity='debug'>
  <Grid name='myGrid'>
    <variable
      name='Materials|heatStructure1|thermal_conductivity' >
      <distribution>ThermalConductivity1</distribution>
      <grid type='value' construction='custom' >0.6
        0.7 0.8</grid>
    </variable>
    <variable name='Materials|heatStructure1|specific_heat' >
      <distribution >SpecificHeat</distribution>
      <grid type='CDF' construction='custom'>0.5
        1.0 0.0</grid>
    </variable>
    <variable
      name='Materials|heatStructure2|thermal_conductivity'>
      <distribution >ThermalConductivity2</distribution>
      <grid type='value' upperBound='4' construction='equal'
        steps='1'>0.5</grid>
    </variable>
  </Grid>
  <MonteCarlo name='MC_Sampler' limit='1000'>
    <variable
      name='Materials|heatStructure1|thermal_conductivity' >
      <distribution>ThermalConductivity1</distribution>
    </variable>
    <variable name='Materials|heatStructure1|specific_heat' >
      <distribution >SpecificHeat</distribution>
    </variable>
    <variable
      name='Materials|heatStructure2|thermal_conductivity'>
      <distribution >ThermalConductivity2</distribution>
    </variable>
  </MonteCarlo>
</Samplers>

```

### 18.4.5 Steps

For a MooseBasedApp, the **<MultiRun>** step type will most likely be used, as first step. First, the step needs to be named: this name will be one of the names used in the **<Sequence>** block. In our example, Grid\_Sampler and MC\_Sampler.

```
<MultiRun name='Grid_Sampler' >
```

With this step, we need to import all the files needed for the simulation:

- MooseBasedApp YAML input file;
- eventual restart files (optional);
- other auxiliary files (e.g., powerHistory tables, etc.).

```
<Input class='Files' type=''>mooseBasedApp.i</Input>  
<Input class='Files' type=''>0020_mesh.cpr</Input>  
<Input class='Files' type=''>0020.xdr.0000</Input>  
<Input class='Files' type=''>0020.rd-0</Input>
```

We then need to define which model will be used:

```
<Model class='Models' type='Code'>MyMooseBasedApp</Model>
```

We then need to specify which Sampler is used, and this can be done as follows:

```
<Sampler class='Samplers' type='Grid'>Grid_Sampler</Sampler>
```

And lastly, we need to specify what kind of output the user wants. For example the user might want to make a database (in RAVEN the database created is an HDF5 file) and a DataObject of type PointSet, to use in sub-sequential post-processing. Here is a classical example:

```
<Output class='Databases' type='HDF5'>MC_out</Output>  
<Output class='DataObjects'  
  type='PointSet'>MCOutData</Output>
```

Following is the example of two MultiRun steps which use different sampling methods (grid and Monte Carlo), and creating two different databases for each one:

```
<Steps verbosity='debug'>  
  <MultiRun name='Grid_Sampler' verbosity='debug'>  
    <Input class='Files' type=''>mooseBasedApp.i</Input>  
    <Input class='Files' type=''>0020_mesh.cpr</Input>  
    <Input class='Files' type=''>0020.xdr.0000</Input>  
    <Input class='Files' type=''>0020.rd-0</Input>  
    <Model class='Models' type='Code'>MyMooseBasedApp</Model>  
    <Sampler class='Samplers' type='Grid'>Grid_Sampler</Sampler>
```

```

    <Output class='Databases' type='HDF5'>Grid_out</Output>
    <Output class='DataObjects'
      type='PointSet'>gridOutData</Output>
  </MultiRun>
  <MultiRun name='MC_Sampler' verbosity='debug'
    re-seeding='210491'>
    <Input class='Files' type=''>mooseBasedApp.i</Input>
    <Input class='Files' type=''>0020_mesh.cpr</Input>
    <Input class='Files' type=''>0020.xdr.0000</Input>
    <Input class='Files' type=''>0020.rd-0</Input>
    <Model class='Models' type='Code'>MyMooseBasedApp</Model>
    <Sampler class='Samplers' type='MonteCarlo'
      >MC_Sampler</Sampler>
    <Output class='Databases' type='HDF5' >MC_out</Output>
    <Output class='DataObjects'
      type='PointSet'>MCOutData</Output>
  </MultiRun>
</Steps>

```

#### 18.4.6 Databases

As shown in the `<Steps>` block, the code is creating two database objects called `Grid_out` and `MC_out`. So the user needs to input the following:

```

<Databases>
  <HDF5 name="Grid_out"/>
  <HDF5 name="MC_out"/>
</Databases>

```

As listed before, this will create two databases. The files will have names corresponding to their `name` appended with the `.h5` extension (i.e. `Grid_out.h5` and `MC_out.h5`).

#### 18.4.7 DataObjects

As shown in the `<Steps>` block, the code is creating two `DataObjects` of type `PointSet` called `gridOutData` and `MCOutData`. So the user needs to input the following:

```

<DataObjects>
  <PointSet name='gridOutData'>
    <Input>

```

```

        Materials|heatStructure2|thermal_conductivity,
        Materials|heatStructure1|specific_heat,
        Materials|heatStructure2|thermal_conductivity
    </Input>
    <Output>aveTempLeft</Output>
</PointSet>
<PointSet name='MCOutData'>
    <Input>
        Materials|heatStructure2|thermal_conductivity,
        Materials|heatStructure1|specific_heat,
        Materials|heatStructure2|thermal_conductivity
    </Input>
    <Output>aveTempLeft</Output>
</PointSet>
</DataObjects>

```

As listed before, this will create two DataObjects that can be used in sub-sequential post-processing.

#### 18.4.8 OutStreamManager

As fully explained in section 14, if the user want to print out or plot the content of a **DataObjects**, he needs to create an **OutStream** in the **<OutStreamManager>** XML block.

As it shown in the example below, for MooseBasedApp (and any other Code interface that might use the symbol | for the Sampler's variable syntax), in the Plot **<x>** and **<y>** specification, the user needs to utilize curly brackets.

```

<OutStreamManager>
  <Print name='gridOutDataDumpCSV'>
    <type>csv</type>
    <source>gridOutData</source>
  </Print>
  <Plot verbosity='debug' name='test' dim='2' overwrite='False'>
    <plotSettings>
      <plot>
        <type>line</type>
        <x>MCOutData|Input|{Materials|heatStructure2|thermal_conductivity}</x>
        <y>MCOutData|Output|aveTempLeft</y>
        <kwargs><color>blue</color></kwargs>
      </plot>
    </plotSettings>
    <actions><how>screen,png</how></actions>
  </Plot>

```

```
</OutputStreamManager>
```

## 18.5 MooseVPP Interface

The Moose Vector Post Processor is used mainly in the solid mechanics analysis. This interface loads the values of the vector output processor to a `<DataObjects>` object.

To use this interface the [DomainIntegral] needs to be present in the MooseBasedApp's input file and the subnode `<fileargs>` should be defined in the subnode `<Code>` in the `<Models>` block of the RAVEN input file. The `<fileargs>` is required to have attributes with the below specified values:

- `type`, *string, required field*, must be "MooseVPP"
- `arg`, *string, required field*, the string value attached to the vector post processor action while creating the output files.

This interface is actually identical to the MooseBasedApp interface, however there is few constraints on defining the output values of the post processor. The definition of these outputs in the `<DataObjects>` depends on the definition of the [DomainIntegral].

The location of the value outputted is defined as *ID#* and the value is as *value#*. The "#" defines the number of the location. The example below contains 3 locations in the [DomainIntegral] where the values are outputted.

Example:

```
...
<Models>
  <Code name="MOOSETestApp" subType="MooseBasedApp">
    <executable>%FRAMEWORK_DIR%/../../moose/
      modules/combined/modules-%METHOD%</executable>
    <fileargs type = "MooseVPP" arg = "_J_1_" />
    <alias variable = "poissonsRatio" >
      Materials|stiffStuff|poissons_ratio</alias>
    <alias variable = "youngModulus" >
      Materials|stiffStuff|youngs_modulus</alias>
  </Code>
</Models>
...
<DataObjects>
```

```

<PointSet name="collset">
  <Input>youngModulus,poissonsRatio</Input>
  <Output>ID1, ID2, ID3, value1, value2, value3</Output>
</PointSet>
</DataObjects>
...

```

## 18.6 OpenModelica Interface

OpenModelica (<http://www.openmodelica.org>) is an open source implementation of the Modelica simulation language. Modelica is "a non-proprietary, object-oriented, equation based language to conveniently model complex physical systems containing, e.g., mechanical, electrical, electronic, hydraulic, thermal, control, electric power or process-oriented subcomponents."<sup>1</sup> Modelica models are specified in text files with a file extension of .mo. A standard Modelica example called BouncingBall which simulates the trajectory of an object falling in one dimension from a height is shown as an example:

```

model BouncingBall
  parameter Real e=0.7 "coefficient_of_restitution";
  parameter Real g=9.81 "gravity_acceleration";
  Real h(start=1) "height_of_ball";
  Real v "velocity_of_ball";
  Boolean flying(start=true) "true,_if_ball_is_flying";
  Boolean impact;
  Real v_new;
  Integer foo;

equation
  impact = h <= 0.0;
  foo = if impact then 1 else 2;
  der(v) = if flying then -g else 0;
  der(h) = v;

  when {h <= 0.0 and v <= 0.0, impact} then
    v_new = if edge(impact) then -e*pre(v) else 0;
    flying = v_new > 0;
    reinit(v, v_new);
  end when;

end BouncingBall;

```

<sup>1</sup><http://www.modelica.org>

---

### 18.6.1 Files

An OpenModelica installation specific to the operating system is used to create a stand-alone executable program that performs the model calculations. A separate XML file containing model parameters and initial conditions is also generated as part of the build process. The RAVEN OpenModelica interface modifies input parameters by changing copies of this file. Both the executable and XML parameter file names must be provided to RAVEN. In the case of the BouncingBall model previously mentioned on the Windows operating system, the <Files>specification would look like:

```
<Files>
  <Input name='BouncingBall_init.xml'
    type=''>BouncingBall_init.xml</Input>
  <Input name='BouncingBall.exe' type=''>BouncingBall.exe</Input>
</Files>
```

### 18.6.2 Models

OpenModelica models may provide simulation output in a number of formats. The particular format used is specified during the model generation process. RAVEN works best with Comma-Separated Value (CSV) files, which is one of the possible output format options. Models are generated using the OpenModelica Shell (OMS) command-line interface, which is part of the OpenModelica installation. To generate an executable that provides CSV-formatted output, use OMSI commands as follows:

1. Change to the directory containing the .mo file to generate an executable for:

```
>> cd("C:/MinGW/msys/1.0/home/bobk/projects/raven/framework/
↪ CodeInterfaces/OpenModelica")
"C:/MinGW/msys/1.0/home/bobk/projects/raven/framework/
↪ CodeInterfaces/OpenModelica"
```

2. Load the model file into memory:

```
>> loadFile("BouncingBall.mo")
true
```

3. Create the model executable, specifying CSV output format:

```
>> buildModel(BouncingBall, outputFormat="csv")
{"C:/MinGW/msys/1.0/home/bobk/projects/raven/framework/
  ↳ CodeInterfaces/OpenModelica/BouncingBall", "
  ↳ BouncingBall_init.xml"}
Warning: The initial conditions are not fully specified. Use
  ↳ +d=initialization for more information.
```

At this point the model executable and XML initialization file should have been created in the same directory as the original model file.

The model executable is specified to RAVEN using the <Models>section of the input file as follows:

```
<Simulation>
  ...
  <Models>
    <Code name="BouncingBall" subType = "OpenModelica">
      <executable>BouncingBall.exe</executable>
    </Code>
  </Models>
  ...
</Simulation>
```

### 18.6.3 CSV Output

The CSV files produced by OpenModelica model executables require adjustment before it may be read by RAVEN. The first few lines of original CSV output from the BouncingBall example is shown below:

```
"time", "h", "v", "der(h)", "der(v)", "v_new", "foo", "flying", "impact",
0, 1, 0, 0, -9.810000000000001, 0, 2, 1, 0,
  ...
```

RAVEN will not properly read this file as-generated for two reasons:

- The variable names in the first line are each enclosed in double-quotes.
- Each line has a trailing comma.

The OpenModelica interface will automatically remove the double-quotes and trailing commas through its implementation of the finalizeCodeOutput function.



## 18.7 Dymola Interface

Modelica is "a non-proprietary, object-oriented, equation-based language to conveniently model complex physical systems containing, e.g., mechanical, electrical, electronic, hydraulic, thermal, control, electric power or process-oriented subcomponents."<sup>2</sup> Modelica models (with a file extension of .mo) are built, translated (compiled), and simulated in Dymola (<http://www.modelon.com/products/dymola/>), which is a commercial modeling and simulation environment based on the Modelica modeling language. A standard Modelica example called BouncingBall, which simulates the trajectory of an object falling in one dimension from a height, is shown as an example:

```
model BouncingBall
  parameter Real e=0.7 "coefficient_of_restitution";
  parameter Real g=9.81 "gravity_acceleration";
  parameter Real hstart = 10 "height_of_ball_at_time_zero";
  parameter Real vstart = 0 "velocity_of_ball_at_time_zero";
  Real h(start=hstart,fixed=true) "height_of_ball";
  Real v(start=vstart,fixed=true) "velocity_of_ball";
  Boolean flying(start=true) "true,_if_ball_is_flying";
  Boolean impact;
  Real v_new;
  Integer foo;

equation
  impact = h <= 0.0;
  foo = if impact then 1 else 2;
  der(v) = if flying then -g else 0;
  der(h) = v;

  when {h <= 0.0 and v <= 0.0,impact} then
    v_new = if edge(impact) then -e*pre(v) else 0;
    flying = v_new > 0;
    reinit(v, v_new);
  end when;

  annotation (uses(Modelica(version="3.2.1")),
    experiment(StopTime=10, Interval=0.1),
    __Dymola_experimentSetupOutput);

end BouncingBall;
```

---

<sup>2</sup><http://www.modelica.org>

## 18.7.1 Files

When a modelica model, e.g., BouncingBall model, is implemented in Dymola, the platform dependent C-code from a Modelica model and the corresponding executable code (i.e., by default dymosim.exe on the Windows operating system) are generated for simulation. After the executable is generated, it may be run multiple times (with Dymola license). A separate TEXT file (by default dsin.txt) containing model parameters and initial conditions are also generated as part of the build process. The RAVEN Dymola interface modifies input parameters by changing copies of this file. Both the executable and TEXT parameter file (or simulation initialization file) names must be provided to RAVEN. In the case of the BouncingBall model previously mentioned on the Windows operating system, the <Files> specification would look like:

```
<Files>
  <Input name='dsin.txt' type=''>dsin.txt</Input>
</Files>
```

## 18.7.2 Models

An executable (dymosim.exe) and a simulation initialization file (dsin.txt) can be generated after either translating or simulating the Modelica model (BouncingBall.mo) using the Dymola Graphical User Interface (GUI) or Dymola Application Programming Interface (API)-routines. To generate an executable and a simulation initialization file, use the Dymola API-routines (or Dymola GUI) to translate the model as follows:

1. Change to the directory containing the .mo file to generate an executable. In Dymola GUI, this corresponds to File/Change Directory in menus:

```
>> cd("C:/msys64/home/KIMJ/projects/raven/framework/
    ↪ CodeInterfaces/Dymola");
C:/msys64/home/KIMJ/projects/raven/framework/CodeInterfaces/
    ↪ Dymola
= true
```

2. Reads the specified file and displays its window. In Dymola GUI, this corresponds to File/Open in the menus:

```
>> openModel("BouncingBall.mo")
= true
```

3. Compile the model (with current settings), and create the model executable and the corresponding simulation initialization file. In Dymola GUI, this corresponds to Translate Model in the menus:

```
>> translateModel("BouncingBall");  
= true
```

At this point the model executable and the simulation initialization file should have been created in the same directory as the original model file. Additionally, they could be created by simulating the model. The following command corresponds to Simulate in the menus in Dymola GUI:

```
>> simulateModel("BouncingBall", stopTime=10,  
↳ numberOfIntervals=0, outputInterval=0.1, method="dassl"  
↳ , resultFile="BouncingBall");  
= true
```

The file extension (.mat) is automatically added to a output file (resultFile), e.g., BouncingBall.mat. If the generated executable code is triggered directly from a command prompt, the output file is always named as "dsres.mat".

The model executable is specified to RAVEN using the <Models>section of the input file as follows:

```
<Simulation>  
...  
<Models>  
  <Code name="BouncingBall" subType = "Dymola">  
    <executable>dymosim.exe</executable>  
  </Code>  
</Models>  
...  
</Simulation>
```

RAVEN works best with Comma-Separated Value (CSV) files. Therefore, the default .mat output type needs to be converted to .csv output. The Dymola interface will automatically convert the .mat output to human-readable forms, i.e., .csv output, through its implementation of the finalizeCodeOutput function.

## 18.8 Mesh Generation Coupled Interfaces

Some software requires a provided mesh that requires a separate code run to generate. In these cases, we use sampled geometric variables to generate a new mesh for each perturbation of the original problem, then run the input with the remainder of the perturbed parameters and the perturbed mesh. RAVEN currently provides two interfaces for this type of calculation, listed below.

## 18.8.1 MooseBasedApp and Cubit Interface

Many MOOSE-based applications use Cubit (<https://cubit.sandia.gov>) to generate Exodus II files as geometry and meshing for calculations. To use the developed interface, Cubit's bin directory must be added to the user's PYTHONPATH. Input parameters for Cubit can be listed in a journal (.jou) file. Parameter values are typically hardcoded into the Cubit command syntax, but variables may be predefined in a journal file through Aprepro syntax. This is an example of a journal file that generates a rectangle of given height and width, meshes it, defines its volume and sidesets, lists its element type, and writes it as an Exodus file:

```
#{x = 3}
#{y = 3}
#{out_name = "'out_mesh.e'"}
create surface rectangle width {x} height {y} zplane
mesh surface 1
set duplicate block elements off
block 1 surface 1
Sideset 1 curve 3
Sideset 2 curve 4
Sideset 3 curve 1
Sideset 4 curve 2
Block all element type QUAD4
export genesis {out_name} overwrite
```

The first three lines are the Aprepro variable definitions that RAVEN requires to insert sampled variables. All variables that RAVEN samples need to be defined as Aprepro variables in the journal file. One essential caveat to running this interface is that an Aprepro variable **MUST** be defined with the name "out\_name". In order to run this script without RAVEN inserting the correct syntax for the output file name and properly generate the Exodus file for a mesh, the output file name is **REQUIRED** to be in both single and double quotation marks with the file extension appended to the end of the file base name (e.g. "'output\_file.e'").

### 18.8.1.1 Files

**<Files>** works the same as in other interfaces with name and type attributes for each node entry. The **name** attribute is a user-chosen internal name for the file contained in the node, and **type** identifies which base-level interface the file is used within. **<type>** should only be specified for inputs that RAVEN will perturb. For Moose input files, **<type>** should be 'MooseInput' and for Cubit journal files, the **<type>** should be 'CubitInput'. The node should contain the path to the file from the working directory. The following is an example of a typical **<Files>** block.

```
<Files>
```

```

<Input name='moose_test'
  ↪ type='MooseInput'>simple_diffusion.i</Input>
<Input name='mesh_in'
  ↪ type='CubitInput'>rectangle.jou</Input>
<Input name='other_file' type=''
  ↪ >some_file_moose_input_needs.ext</Input>
</Files>

```

### 18.8.1.2 Models

A user provides paths to executables and aliases for sampled variables within the **<Models>** block. The **<Code>** block will contain attributes name and subType. Name identifies that particular **<Code>** model within RAVEN, and subType specifies which code interface the model will use. The **<executable>** block should contain the absolute or relative (with respect to the current working directory) path to the MooseBasedApp that RAVEN will use to run generated input files. The absolute or relative path to the Cubit executable is specified within **<preexec>**. If the **<preexec>** block is not needed, the MooseBasedApp interface is probably preferable to the Cubit-Moose interface.

Aliases are defined by specifying the variable attribute in an **<alias>** node with the internal RAVEN variable name chosen with the node containing the model variable name. The Cubit-Moose interface uses the same syntax as the MooseBasedApp to refer to model variables, with pipes separating terms starting with the highest YAML block going down to the individual parameter that RAVEN will change. To specify variables that are going to be used in the Cubit journal file, the syntax is "Cubit—aprepro\_var". The Cubit-Moose interface will look for the Cubit tag in all variables passed to it and upon finding it, send it to the Cubit interface. If the model variable does not begin with 'Cubit', the variable MUST be specified in the MooseBasedApp input file. While the model variable names are not required to have aliases defined (the **<alias>** blocks are optional), it is highly suggested to do so not only to ensure brevity throughout the RAVEN input, but to easily identify where variables are being sent in the interface.

An example **<Models>** block follows.

```

<Models>
  <Code name="moose-modules" subType="CubitMoose">
    <executable>%FRAMEWORK_DIR%/../../../../moose/modules/combined/...
      modules-%METHOD%</executable>
    <preexec>/hpc-common/apps/local/cubit/13.2/bin/cubit</preexec>
    <alias variable="length">Cubit|y</alias>
    <alias variable="bot_BC">BCs|bottom|value</alias>
  </Code>
</Models>

```

### 18.8.1.3 Distributions

The `<Distributions>` block defines all distributions used to sample variables in the current RAVEN run.

For all the possible distributions and their possible inputs please refer to the Distributions chapter (see 10). It is good practice to name the distribution something similar to what kind of variable is going to be sampled, since there might be many variables with the same kind of distributions but different input parameters.

### 18.8.1.4 Samplers

The `<Samplers>` block defines the variables to be sampled.

After defining a sampling scheme, the variables to be sampled and their distributions are identified in the `<variable>` blocks. The name attribute in the `<variable>` block must either be the full MooseBasedApp model variable name or the alias name specified in `<Models>`. If the sampled variable is a geometric property that will be used to generate a mesh with Cubit, remember the syntax for variables being passed to journal files (Cubit—`aprepro_var`).

For listings of available samplers refer to the Samplers chapter (see 11).

See the following for an example of a grid based sampler for length and the bottom boundary condition (both of which have aliases defined in `<Models>`).

```
<Samplers>
  <Grid name="Grid_sampling">
    <variable name="length" >
      <distribution>length_dist</distribution>
      <grid type="value" construction="custom">1.0 2.0</grid>
    </variable>
    <variable name="bot_BC">
      <distribution>bot_BC_dist</distribution>
      <grid type="value" construction="custom">3.0 6.0</grid>
    </variable>
  </Grid>
</Samplers>
```

### 18.8.1.5 Steps,OutStreamManager,DataObjects

This interface's `<Steps>`, `<OutStreamManager>`, and `<DataObjects>` blocks do not deviate significantly from other interfaces' respective nodes. Please refer to previous entries for these blocks if needed.

### 18.8.1.6 File Cleanup

The Cubit-Moose interface automatically removes files that are commonly unwanted after the RAVEN run reaches completion. Cubit has been described as "talkative" due to additional journal files with execution information being generated by the program after every completed journal file run. The quantity of these files can quickly become unwieldy if the working directory is not kept clean; thus these files are removed. In addition, some users may wish to remove Exodus files after the RAVEN run is complete as the typical size of each file is quite large and it is assumed that any output quantities of interest will be collected by appropriate postprocessors and the OutStreamManager. Exodus files are not automatically removed, but by using the `<deleteOutExtension>` node in `<RunInfo>`, one may specify the Exodus extension to save a fair amount of storage space after RAVEN completes a sequence. For example:

```
<RunInfo>
...
  <deleteOutExtension>e</deleteOutExtension>
...
</RunInfo>
```

## 18.8.2 MooseBasedApp and Bison Mesh Script Interface

For BISON users, a Python mesh generation script is included in the `%BISON_DIR%/tools/UO2/` directory. This script generates 3D or 2D (RZ) meshes for nuclear fuel rods using Cubit with templated commands. The BISON Mesh Script (BMS) is capable of generating rods with discrete fuel pellets of various size in assorted configurations. To use this interface, Cubit's bin directory must be added to the user's PYTHONPATH.

### 18.8.2.1 Files

Similar to the Cubit-Moose interface, the BisonAndMesh interface requires users to specify all files required to run their input so that these file may be copied into the respective sequence's working directory. The user will give each file an internal RAVEN designation with the name attribute, and

the MooseBasedApp and BISON Mesh Script inputs must be assigned their respective types in another attribute of the `<Input>` node. An example follows.

```

<Files>
  <Input name='bison_test'
    ↪ type='MooseInput'>simple_bison_test.i</Input>
  <Input name='mesh_in'
    ↪ type='BisonMeshInput'>coarse_input.py</Input>
  <Input name='other_file'
    ↪ type=''>some_file_moose_input_needs.ext</Input>
</Files>

```

### 18.8.2.2 Models

A user provides paths to executables and aliases for sampled variables within the `<Models>` block. The `<Code>` block will contain attributes `name` and `subType`. `name` identifies that particular `<Code>` model within RAVEN, and `subType` specifies which code interface the model will use. The `<executable>` block should contain the absolute or relative (with respect to the current working directory) path to the MooseBasedApp that RAVEN will use to run generated input files. The absolute or relative path to the mesh script python file is specified within `<preexec>`. If the `<preexec>` block is not needed, use the MooseBasedApp interface.

Aliases are defined by specifying the variable attribute in an `<alias>` node with the internal RAVEN variable name chosen with the node containing the model variable name. The BisonAndMesh interface uses the same syntax as the MooseBasedApp to refer to model variables, with pipes separating terms starting with the highest YAML block going down to the individual parameter that RAVEN will change. To specify variables that are going to be used in the BISON Mesh Script python input, the syntax is "Cubit—dict\_name—var\_name". The interface will look for the Cubit tag in all variables passed to it and upon finding the tag, send it to the BISON Mesh Script interface. If the model variable does not begin with Cubit, the variable MUST be specified in the MooseBasedApp input file. While the model variable names are not required to have aliases defined (the `<alias>` blocks are optional), it is highly suggested to do so not only to ensure brevity throughout the RAVEN input, but to easily identify where variables are being sent in the interface.

An example `<Models>` block follows.

```

<Models>
  <Code name="Bison-opt" subType="BisonAndMesh">
    <executable>%FRAMEWORK_DIR%/../../../../bison/bison-%METHOD%</executable>
    <preexec>%FRAMEWORK_DIR%/../../../../bison/tools/U02/mesh_script.py</preexec>
    <alias variable="pellet_radius"
      ↪ >Cubit|Pellet1|outer_radius</alias>

```



```

    <alias
      ↪ variable="clad_thickness">Cubit|clad|clad_thickness</alias>
    <alias variable="fuel_k"
      ↪ >Materials|fuel_thermal|thermal_conductivity</alias>
    <alias variable="clad_k"
      ↪ >Materials|clad_thermal|thermal_conductivity</alias>
  </Code>
</Models>

```

### 18.8.2.3 Distributions

The **<Distributions>** block defines all distributions used to sample variables in the current RAVEN run.

For all the possible distributions and their possible inputs please refer to the Distributions chapter (see 10). It is good practice to name the distribution something similar to what kind of variable is going to be sampled, since there might be many variables with the same kind of distributions but different input parameters.

### 18.8.2.4 Samplers

The **<Samplers>** block defines the variables to be sampled.

After defining a sampling scheme, the variables to be sampled and their distributions are identified in the **<variable>** blocks. The name attribute in the **<variable>** block must either be the full MooseBasedApp model variable name or the alias name specified in **<Models>**. If the sampled variable is a geometric property that will be used to generate a mesh with Cubit, remember the syntax for variables being passed to journal files (Cubit—`aprepro_var`).

For listings of available samplers refer to the Samplers chapter (see 11).

See the following for an example of a grid based sampler for length and the bottom boundary condition (both of which have aliases defined in **<Models>**).

```

<Samplers>
  <Grid name="Grid_sampling">
    <variable name="length" >
      <distribution>length_dist</distribution>
      <grid type="value" construction="custom">1.0 2.0</grid>
    </variable>
    <variable name="bot_BC">

```

```
    <distribution>bot_BC_dist</distribution>
    <grid type="value" construction="custom">3.0 6.0</grid>
  </variable>
</Grid>
</Samplers>
```

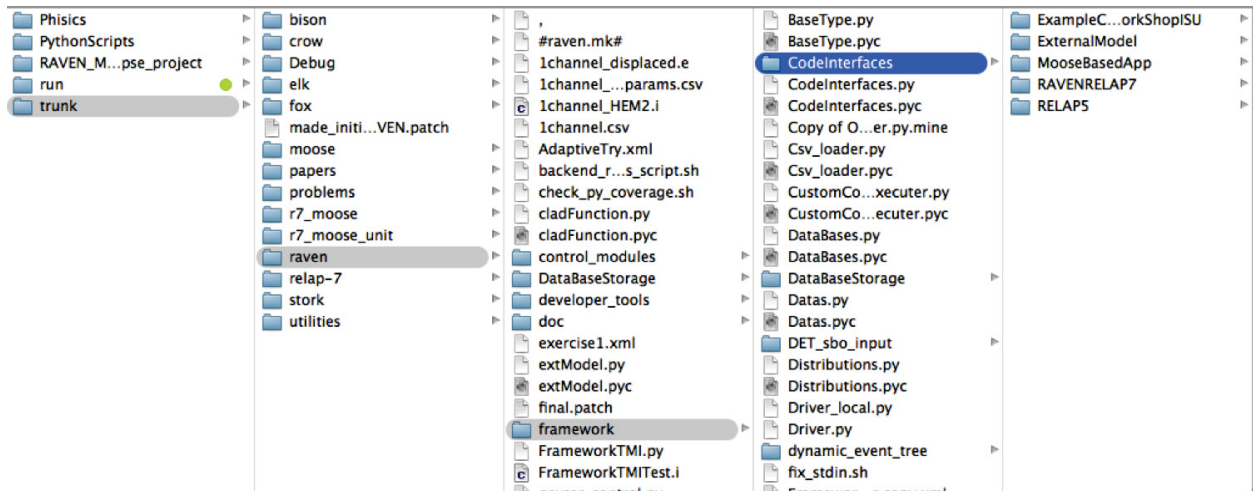
### 18.8.2.5 Steps,OutStreamManager,DataObjects

This interface's `<Steps>`, `<OutStreamManager>`, and `<DataObjects>` blocks do not deviate significantly from other interfaces' respective nodes. Please refer to previous entries for these blocks if needed.

### 18.8.2.6 File Cleanup

The BisonAndMesh interface automatically removes files that are commonly unwanted after the RAVEN run reaches completion. Cubit has been described as "talkative" due to additional journal files with execution information being generated by the program after every completed journal file run. The BISON Mesh Script creates a journal file to run with cubit after reading input parameters; so Cubit will generate its "redundant" journal files, and .pyc files will litter the working directory as artifacts of the python mesh script reading from the .py input files. The quantity of these files can quickly become unwieldy if the working directory is not kept clean, thus these files are removed. Some users may wish to remove Exodus files after the RAVEN run is complete as the typical size of each file is quite large and it is assumed that any output quantities of interest will be collected by appropriate postprocessors and the OutStreamManager. Exodus files are not automatically removed, but by using the `<deleteOutExtension>` node in `<RunInfo>`, one may specify the Exodus extension (\*.e) to save a fair amount of storage space after RAVEN completes a sequence. For example:

```
<RunInfo>
  ...
  <deleteOutExtension>e</deleteOutExtension>
  ...
</RunInfo>
```



**Figure 1.** Code Interface Location.

## 19 Advanced Users: How to couple a new code

The procedure of coupling a new code/application with RAVEN is a straightforward process. For all the systems codes currently supported by RAVEN (e.g. RELAP-7, RELAP5-3D, BISON, MOOSE), the coupling is performed through a Python interface that interprets the information coming from RAVEN and translates them to the input of the system code. The coupling procedure does not require modifying RAVEN itself. Instead, the developer creates a new Python interface that is going to be embedded in RAVEN at run-time (no need to introduce hard-coded coupling statements). This interface needs to be placed in a folder (whatever name) located in (see figure 1):

```
path/to/raven/distribution/raven/framework/CodeInterfaces/
```

At the initialization stage, RAVEN imports all the Interfaces that are contained in this directory and performs some preliminary cross-checks. In the following sub-sections, a step-by-step procedure for coupling a code to RAVEN is outlined.

### 19.1 Pre-requisites.

In order to couple a newer application to the RAVEN code, some pre-requisites need to be satisfied.

#### *Input*

The first pre-requisite is the knowledge of the input syntax of the application the developer wants

to couple. Indeed, RAVEN task “ends” at the Code Interface stage. RAVEN transfers the information needed to perturb the input space into the Code interface and expects that the newly developed Interface is able to perturb the input files based on the information passed through.

This means that the developer needs to code a Python-compatible parser of the system code input (a module that is able to read and modify the input of the code that needs to be coupled).

For example, let’s suppose the input syntax of the code the developer needs to couple is as follows:

```
keyword1 = aValue1
keyword2 = aValue2
keyword3 = aValue3
keyword4 = aValue4
```

The Python input parser would be:

```
class simpleInputParser():
    def __init__(self, filename):
        #
        # @ In, string, filename, input file name (with path)
        #
        self.keywordDictionary = {}
        # open the file
        fileobject = open(filename)
        # store all the lines into a list
        lines = fileobject.readlines()
        # parse the list to construct
        # self.keywordDictionary dictionary
        for line in lines:
            # split the line with respect
            # to the symbol "=" and store the
            # outcomes into the dictionary
            # listSplitted[0] is the keyword
            # listSplitted[1] is the value
            listSplitted = line.split("=")
            keyword = listSplitted[0]
            value = listSplitted[1]
            self.keywordDictionary[keyword] = value
        # close the file
        fileobject.close()

    def modifyInternalDictionary(self, inDictionary):
        #
        # @ In, dictionary {keyword:value},
        # inDictionary, dictionary containing
        # the keywords to perturb
```

```

#

# we just parse the dictionary and replace the
# matching keywords
for keyword,newvalue in inDictionary.items():
    self.keywordDictionary[keyword] = newvalue

def writeNewInput(self,filename):
#
# @ In, string, filename, newer input file name (with path)
#

# open the file
fileobject = open(filename)
# write line by line
for keyword,newvalue in self.keywordDictionary.items():
    fileobject.write(keyword + '=' + str(newvalue) + '\n')
# close the file
fileobject.close()

```

### Output

RAVEN is able to handle Comma Separated Value (CSV) files (as outputs of the system code). In order make RAVEN able to retrieve the information from the newly coupled code, these files need to be either generated by the system code itself or the developer needs to code a Python-compatible module to convert the whatever code output format to a CSV one. This module can be directly called in the new code interface (see following section).

Let's suppose that the output format of the code (the same of the previous input parser example) is as follows:

```

result1 = aValue1
result2 = aValue2
result3 = aValue3

```

The Python output converter would be as simple as:

```

def convertOutputFileToCSV(outputfile):
    keywordDictionary = {}
    # open the original file
    fileobject = open(outputfile)
    outputCSVfile = open (outputfile + '.csv')
    # store all the lines into a list
    lines = fileobject.readlines()

```

```

# parse the list to construct
# self.keywordDictionary dictionary
for line in lines:
    # split the line with respect
    # to the symbol "=" and store the
    # outcomes into the dictionary
    # listSplitted[0] is the keyword
    # listSplitted[1] is the value
    listSplitted = line.split("=")
    keyword = listSplitted[0]
    value    = listSplitted[1]
    keywordDictionary[keyword] = value
outputCSVfile.write(','.join(keywordDictionary.keys()))
outputCSVfile.write(','.join(keywordDictionary.values()))
outputCSVfile.close()

```

And the output CSV becomes:

```

result1, result2, result3
aValue1, aValue2, aValue3

```

## 19.2 Code Interface Creation

As already mentioned, RAVEN imports all the “Code Interfaces” at run-time, without actually knowing the syntax of the driven codes. In order to make RAVEN able to drive a newer software, the developer needs to code a Python module that will contain few methods (with strict syntax) that are called by RAVEN during the simulation.

When loading a “Code Interface”, RAVEN expects to find, in the class representing the code, the following functions:

```

from CodeInterfaceBaseClass import CodeInterfaceBase
class NewCode(CodeInterfaceBase):
    def generateCommand(self, inputFiles, executable, clargs=None,
        ↪ fargs=None)
    def createNewInput(self, currentInputFiles, oriInputFiles,
        samplerType, **Kwargs)
    def finalizeCodeOutput(self, command, output, workingDir)
    def getInputExtension(self)
    def checkForOutputFailure(self, output, workingDir)

```

In the following sub-sections all the methods are fully explained, providing examples (referring to the simple code used as example for the previous sections)

### 19.2.1 Method: `generateCommand`

```
def generateCommand(self, inputFiles, executable, clargs=None,
    ↪ fargs=None)
```

The `generateCommand` method is used to generate the commands (in `string` format) needed to launch the driven Code, as well as the root name of the output of the perturbed inputs (in `string` format). The return for this command is a two-part Python tuple. The first entry is a list of two-part tuples, each which specifies whether the corresponding command should be run exclusively in serial, or whether it can be run in parallel, as well as the command itself. For example, for a command where two successive commands are called, the first in serial and the second in parallel,

```
def generateCommand(self, inputFiles, executable, clargs=None,
    ↪ fargs=None):
    . . .
    commands = [('serial', first_command), ('parallel',
    ↪ second_command)]
    return (commands, outFileRoot)
```

For each command, the second entry in the tuple is a string containing the full command that the internal JobHandler is going to use to run the Code this interface refers to. The return data type must be a Python tuple with a list of tuples and a string: (`commands`, `outfileRoot`). Note that in most cases, only a single command needs to be run, so only a single command tuple is necessary. At run time, RAVEN will string together commands attached by double ampersands (&&), and each command labeled as parallel-compatible will be prepended with appropriate mpi arguments. For the example above, the command executed will be (with `<NumMPI>` equal to 4)

```
$ first_command && mpiexec -n 4 second_command
```

RAVEN is going to call the `generateCommand` function passing in the following arguments:

- **inputFiles**, data type = list: List of input files (length of the list depends on the number of inputs listed in the Step which is running this code);
- **executable**, data type = string, executable name with absolute path (e.g. `/home/path_to_executable/code.exe`);
- **clargs**, *optional*, data type = dictionary, a dictionary containing the command-line flags the user can specify in the input (e.g. under the node `<Code>` `< clargstype = 'input' arg = ' -i' extension = '.inp' / >` `< /Code >`).

- **fargs**, *optional*, data type = dictionary, a dictionary containing the auxiliary input file variables the user can specify in the input (e.g. under the node `< Code >< clargstype = 'input' arg = 'aux' extension = '.aux' / >< /Code >`).

For the example referred to in the previous section, this method would be implemented as follows:

```
def generateCommand(self, inputFiles, executable, clargs=None,
    ↪ fargs=None):
    found = False
    for index, inputFile in enumerate(inputFiles):
        if inputFile.endswith(self.getInputExtension()):
            found = True
            break
    if not found: raise IOError(
        `None of the input files has one of the following
        ↪ extensions: ` +
        ` ` .join(self.getInputExtension())
    )
    outputfile = 'out~'+os.path.split(inputFiles[index])[1].split
    ↪ ('.')[0]
    executeCommand = [('parallel', executable+ ` -i ` +os.path.
    ↪ split(inputFiles[index])[1])]
    return executeCommand, outputfile
```

### 19.2.2 Method: createNewInput

```
def createNewInput(self, currentInputFiles, oriInputFiles,
    ↪ samplerType, **Kwargs)
```

The **createNewInput** method is used to generate an input based on the information RAVEN passes in. In this function the developer needs to call the driven code input parser in order to modify the input file, accordingly with respect to the variables RAVEN is providing. This method needs to return a list containing the path and filenames of the modified input files. **Note:** RAVEN expects that at least one input file of the original list gets modified and re-named.

RAVEN is going to call this function passing in the following arguments:

- **currentInputFiles**, data type = list: List of current input files (input files from last this method call);
- **oriInputFiles** , data type = list, List of the original input files;



- **samplerType** , data type = string, Sampler type (e.g. MonteCarlo, Adaptive, etc.). **Note:** None if no Sampler has been used;
- **Kwargs** , data type = kwarded dictionary, dictionary of parameters. In this dictionary there is another dictionary called "SampledVars" where RAVEN stores the variables that got sampled (Kwargs['SampledVars'] = {'var1':10,'var2':40});

For the example referred in the previous section, this method would implemented as follows:

```
def createNewInput (self, currentInputFiles,
                   oriInputFiles, samplerType, **Kwargs):
    for index, inputFile in enumerate(oriInputFiles):
        if inputFile.endswith(self.getInputExtension()):
            break
    parser = simpleInputParser(currentInputFiles[index])
    parser.modifyInternalDictionary(**Kwargs['SampledVars'])
    temp = str(oriInputFiles[index][:])
    newInputFiles = copy.copy(currentInputFiles)
    newInputFiles[index] = os.path.join(os.path.split(temp)[0],
                                       Kwargs['prefix']+"~"+os.path.split(temp)[1])
    parser.writeNewInput(newInputFiles[index])
    return newInputFiles
```

### 19.2.3 Method: getInputExtension

```
def getInputExtension(self)
```

The **getInputExtension** function is an optional method. If present, it is called by RAVEN code at run time. This function can be considered an utility method, since its main goal is to return a tuple of strings, where the developer can place all the input extensions the code interface needs to support (i.e. the extensions of the input(s) the code interface is going to “perturb”). If this method is not implemented, the default extensions are (“.i”, “.inp”, “.in”). This function does not accept any input argument. For the example referred in the previous section, this method would implemented as follows:

```
def getInputExtension(self):
    return (".i", ".input")
```

### 19.2.4 Method: `finalizeCodeOutput`

```
def finalizeCodeOutput(self, command, output, workingDir)
```

The **`finalizeCodeOutput`** function is an optional method. If present, it is called by RAVEN code at the end of each run. It can be used for those codes, that do not create CSV files as output to convert the whatever output format into a CSV. RAVEN checks if a string is returned; if so, RAVEN interprets that string as the new output file name (CSV).

RAVEN is going to call this function passing in the following arguments:

- **`command`**, data type = string: the command used to run the just ended job;
- **`output`**, data type = string, the Output name root;
- **`workingDir`**, data type = string, current working directory.

For the example referred in the previous section, this method would implemented as follows:

```
def finalizeCodeOutput(self, command, output, workingDir):  
    outfile = os.path.join(workingDir, output+".o")  
    convertOutputFileToCSV(outfile)
```

### 19.2.5 Method: `checkForOutputFailure`

```
def checkForOutputFailure(self, output, workingDir)
```

The **`checkForOutputFailure`** function is an optional method. If present, it is called by RAVEN code at the end of each run. This method needs to be implemented by the codes that, if a run fails, return a “returncode” = 0. This can happen in those codes that record the failure of a run (e.g. not converged, etc.) as normal termination (returncode == 0) This method can be used, for example, to parse the outputfile looking for a special keyword that testifies that a particular job failed (e.g. in RELAP5 would be the keyword ”\*\*\*\*\*”). This method **MUST** return a boolean (True if failed, False otherwise).

RAVEN is going to call this function passing in the following arguments:

- **`output`**, data type = string, the Output name root;
- **`workingDir`**, data type = string, current working directory.

For the example referred in the previous section, this method would implemented as follows:

```
def checkForOutputFailure(self, command, output, workingDir):
    from __builtin__ import any
    errorWord = "ERROR"
    return any(errorWord in x for x in
               open(os.path.join(workingDir, output+'.o'), "r").readlines())
```

## 19.3 Tools for Developing Code Interfaces

To make generating a code interface as simple as possible, there are several tools RAVEN makes available within the Code Interface objects.

### 19.3.1 File Objects

RAVEN has created a wrapper for files within Python in order to carry along some additional information. This allows the user to tag particular files for reference in the Code Interface, using the **type** XML attribute in **<Files>** nodes. To differentiate, RAVEN file objects will use the capital Files, whereas typical files will use the lowercase files.

When the Files are passed in to `createNewInput`, they are passed in as Files objects. To access the `xmlAttrtype` of a file, use the method `getType`. For instance, instead of looking for an extension, a Code Interface might identify an input file by looking for a particular type, as shown in the example below. **Note:** RAVEN does not access a File's **type**; it is exclusively an optional tool for Code Interface developers.

```
found = False
for inFile in inputFiles:
    if inFile.getType()=='mainInput':
        found = True
        break
if not found:
    raise IOError('Desired file with type ``mainInput`` not found!'
                 ↪ )
```

Using Files **type** attributes can especially help when multiple input files have the same extension. For example, say a Code execution command normally has the following appearance on the command line:

```
/home/path/to/executable/myexec.sh -i mainInp.xml -a auxInp.xml  
↪ --mesh cube.e
```

The **<Files>** block in the RAVEN XML might appear as follows:

```
<Files>  
  <Input name='main' type='base'>mainInp.xml</Input>  
  <Input name='two' type='aux' >auxInp.xml</Input>  
  <Input name='cube' type='mesh' perturbable='False'>cube.e</  
    ↪ Input>  
</Files>
```

The search for these files in the Code Interface might then look like the example below, assuming one file per type:

```
# populate a type dictionary  
typesDict={}  
for inFile in inputFiles:  
    typesDict[inFile.getType()]=inFile  
# check all the necessary files are there  
if 'base' not in typesDict.keys():  
    raise IOError('File type ``base`` not listed in input file!')  
if 'aux' not in typesDict.keys():  
    raise IOError('File type ``aux`` not listed in input file!')  
if 'mesh' not in typesDict.keys():  
    raise IOError('File type ``mesh`` not listed in input file!')  
mainFile = typesDict['base']  
# do operations on file, etc.
```

Additionally, a Code Interface developer can access the **perturbable** through the `getPerturbable()` method of a Files object. This can be useful, for example, in preventing searching binary files for variable names when creating new input. For example,

```
for inFile in inputFiles:  
    if not inFile.getPerturbable(): continue  
    # etc
```

## A Appendix: Example Primer

In this Appendix, a set of examples are reported. In order to be as general as possible, the *Model* type “ExternalModel” has been used.

### A.1 Example 1.

This simple example is about the construction of a “Lorentz attractor”, sampling the relative input space. The parameters that are sampled represent the initial coordinate  $(x_0, y_0, z_0)$  of the attractor origin.

```
<?xml version="1.0" encoding="UTF-8"?>
<Simulation verbosity="debug">
<!-- RUNINFO -->
<RunInfo>
  <WorkingDir>externalModel</WorkingDir>
  <Sequence>FirstMRun</Sequence>
  <batchSize>3</batchSize>
</RunInfo>
<!-- Files -->
<Files>
  <Input name='lorenzAttractor.py'
    ↪ type=''>lorenzAttractor</Input>
</Files>
<!-- STEPS -->
<Steps>
  <MultiRun name='FirstMRun' re-seeding='25061978'>
    <Input class='Files' type=''
      ↪ >lorenzAttractor.py</Input>
    <Model class='Models' type='ExternalModel'
      ↪ >PythonModule</Model>
    <Sampler class='Samplers' type='MonteCarlo'
      ↪ >MC_external</Sampler>
    <Output class='DataObjects' type='HistorySet'
      ↪ >testPrintHistorySet</Output>
    <Output class='Databases' type='HDF5'
      ↪ >test_external_db</Output>
    <Output class='OutputStreamManager' type='Print'
      ↪ >testPrintHistorySet_dump</Output>
  </MultiRun >
</Steps>
```

```

<!-- MODELS -->
<Models>
  <ExternalModel name='PythonModule' subType=''
    ↪ ModuleToLoad='externalModel/lorenzAttractor'>
    <variable>sigma</variable>
    <variable>rho</variable>
    <variable>beta</variable>
    <variable>x</variable>
    <variable>y</variable>
    <variable>z</variable>
    <variable>time</variable>
    <variable>x0</variable>
    <variable>y0</variable>
    <variable>z0</variable>
  </ExternalModel>
</Models>
<!-- DISTRIBUTIONS -->
<Distributions>
  <Normal name='x0_distrib'>
    <mean>4</mean>
    <sigma>1</sigma>
  </Normal>
  <Normal name='y0_distrib'>
    <mean>4</mean>
    <sigma>1</sigma>
  </Normal>
  <Normal name='z0_distrib'>
    <mean>4</mean>
    <sigma>1</sigma>
  </Normal>
</Distributions>
<!-- SAMPLERS -->
<Samplers>
  <MonteCarlo name='MC_external'>
    <samplerInit>
      <limit>3</limit>
    </samplerInit>
    <variable name='x0' >
      <distribution >x0_distrib</distribution>
    </variable>
    <variable name='y0' >
      <distribution >y0_distrib</distribution>

```

```

    </variable>
    <variable name='z0' >
      <distribution >z0_distrib</distribution>
    </variable>
  </MonteCarlo>
</Samplers>
<!-- DATABASES -->
<Databases>
  <HDF5 name="test_external_db"/>
</Databases>
<!-- OUTSTREAMS -->
<OutputStreamManager>
  <Print name='testPrintHistorySet_dump'>
    <type>csv</type>
    <source>testPrintHistorySet</source>
  </Print>
</OutputStreamManager>
<!-- DATA OBJECTS -->
<DataObjects>
  <HistorySet name='testPrintHistorySet'>
    <Input>x0,y0,z0</Input>
    <Output>time,x,y,z</Output>
  </HistorySet>
</DataObjects>
</Simulation>

```

The Python *ExternalModel* is reported below:

```

import numpy as np

def run(self, Input):
    max_time = 0.03
    t_step = 0.01

    numberTimeSteps = int(max_time/t_step)

    self.x = np.zeros(numberTimeSteps)
    self.y = np.zeros(numberTimeSteps)
    self.z = np.zeros(numberTimeSteps)
    self.time = np.zeros(numberTimeSteps)

    self.x0 = Input['x0']
    self.y0 = Input['y0']

```

```

self.z0 = Input['z0']

self.x[0] = Input['x0']
self.y[0] = Input['y0']
self.z[0] = Input['z0']
self.time[0]= 0

for t in range (numberTimeSteps-1):
    self.time[t+1] = self.time[t] + t_step
    self.x[t+1]    = self.x[t] + self.sigma*
                    (self.y[t]-self.x[t]) * t_step
    self.y[t+1]    = self.y[t] + (self.x[t]*
                    (self.rho-self.z[t])-self.y[t]) * t_step
    self.z[t+1]    = self.z[t] + (self.x[t]*
                    self.y[t]-self.beta*self.z[t]) * t_step

```

## A.2 Example 2.

This example shows a slightly more complicated example, that employs the usage of:

- *Samplers*: Grid and Adaptive;
- *Models*: External, Reduce Order Models and Post-Processors;
- *OutStreams*: Prints and Plots;
- *Data Objects*: PointSets;
- *Functions*: ExternalFunctions.

The goal of this input is to compute the “SafestPoint”. It provides the coordinates of the farthest point from the limit surface that is given as an input. The safest point coordinates are expected values of the coordinates of the farthest points from the limit surface in the space of the “controllable” variables based on the probability distributions of the “non-controllable” variables.

The term “controllable” identifies those variables that are under control during the system operation, while the “non-controllable” variables are stochastic parameters affecting the system behaviour randomly.

The “SafestPoint” post-processor requires the set of points belonging to the limit surface, which must be given as an input.



```

<Simulation verbosity='debug'>

<!-- RUNINFO -->
<RunInfo>
  <WorkingDir>SafestPointPP</WorkingDir>
  <Sequence>pth1,pth2,pth3,pth4</Sequence>
  <batchSize>50</batchSize>
</RunInfo>

<!-- STEPS -->
<Steps>
  <MultiRun name = 'pth1' pauseAtEnd = 'False'>
    <Sampler class = 'Samplers' type = 'Grid'
      ↪ >grd_vl_ql_smp_dpt</Sampler>
    <Input class = 'DataObjects' type = 'PointSet'
      ↪ >grd_vl_ql_smp_dpt_dt</Input>
    <Model class = 'Models' type = 'ExternalModel'
      ↪ >xtr_md1</Model>
    <Output class = 'DataObjects' type = 'PointSet'
      ↪ >nt_phy_dpt_dt</Output>
  </MultiRun >

  <MultiRun name = 'pth2' pauseAtEnd = 'True'>
    <Sampler class = 'Samplers' type = 'Adaptive'
      ↪ >dpt_smp</Sampler>
    <Input class = 'DataObjects' type =
      ↪ 'PointSet' >bln_smp_dt</Input>
    <Model class = 'Models' type = 'ExternalModel'
      ↪ >xtr_md1</Model>
    <Output class = 'DataObjects' type =
      ↪ 'PointSet' >nt_phy_dpt_dt</Output>
    <SolutionExport class = 'DataObjects' type =
      ↪ 'PointSet' >lmt_srf_dt</SolutionExport>
  </MultiRun>

  <PostProcess name='pth3' pauseAtEnd = 'False'>
    <Input class = 'DataObjects' type = 'PointSet'
      ↪ >lmt_srf_dt</Input>
    <Model class = 'Models' type = 'PostProcessor'
      ↪ >SP</Model>
    <Output class = 'DataObjects' type = 'PointSet'
      ↪ >sfs_pnt_dt</Output>
  </PostProcess>
</Steps>

```

```

</PostProcess>

<OutputStreamStep name = 'pth4' pauseAtEnd = 'True'>
  <Input class = 'DataObjects' type =
    ↪ 'PointSet' >lmt_srf_dt</Input>
  <Output class = 'OutputStreamManager' type = 'Print'
    ↪ >lmt_srf_dmp</Output>
  <Input class = 'DataObjects' type = 'PointSet'
    ↪ >sfs_pnt_dt</Input>
  <Output class = 'OutputStreamManager' type = 'Print'
    ↪ >sfs_pnt_dmp</Output>
</OutputStreamStep>
</Steps>

<!-- DATA OBJECTS -->
<DataObjects>
  <PointSet name = 'grd_vl_ql_smp_dpt_dt'>
    <Input>x1,x2,gammay</Input>
    <Output>OutputPlaceholder</Output>
  </PointSet>

  <PointSet name = 'nt_phy_dpt_dt'>
    <Input>x1,x2,gammay</Input>
    <Output>g</Output>
  </PointSet>

  <PointSet name = 'bln_smp_dt'>
    <Input>x1,x2,gammay</Input>
    <Output>OutputPlaceholder</Output>
  </PointSet>

  <PointSet name = 'lmt_srf_dt'>
    <Input>x1,x2,gammay</Input>
    <Output>g_zr</Output>
  </PointSet>

  <PointSet name = 'sfs_pnt_dt'>
    <Input>x1,x2,gammay</Input>
    <Output>p</Output>
  </PointSet>
</DataObjects>

```

```

<!-- DISTRIBUTIONS -->
<Distributions>
  <Normal name = 'x1_dst'>
    <upperBound>10</upperBound>
    <lowerBound>-10</lowerBound>
    <mean>0.5</mean>
    <sigma>0.1</sigma>
  </Normal>

  <Normal name = 'x2_dst'>
    <upperBound>10</upperBound>
    <lowerBound>-10</lowerBound>
    <mean>-0.15</mean>
    <sigma>0.05</sigma>
  </Normal>

  <Normal name = 'gammay_dst'>
    <upperBound>20</upperBound>
    <lowerBound>-20</lowerBound>
    <mean>0</mean>
    <sigma>15</sigma>
  </Normal>
</Distributions>

<!-- SAMPLERS -->
<Samplers>
  <Grid name = 'grd_v1_q1_smp_dpt'>
    <variable name = 'x1' >
      <distribution>x1_dst</distribution>
      <grid type = 'value' construction = 'equal' steps = '10'
        ↪ upperBound = '10'>2</grid>
    </variable>
    <variable name='x2' >
      <distribution>x2_dst</distribution>
      <grid type = 'value' construction = 'equal' steps = '10'
        ↪ upperBound = '10'>2</grid>
    </variable>
    <variable name='gammay' >
      <distribution>gammay_dst</distribution>
      <grid type = 'value' construction = 'equal' steps = '10'
        ↪ lowerBound = '-20'>4</grid>
    </variable>
  </Grid>
</Samplers>

```

```

</Grid>

<Adaptive name = 'dpt_smp' verbosity='debug'>
  <ROM class = 'Models' type = 'ROM'
    ↪ >accelerated_ROM</ROM>
  <Function class = 'Functions' type = 'External'
    ↪ >g_zr</Function>
  <TargetEvaluation class = 'DataObjects' type =
    ↪ 'PointSet' >nt_phy_dpt_dt</TargetEvaluation>
  <Convergence limit = '3000' forceIteration = 'False' weight
    ↪ = 'none' persistence = '5'>1e-2</Convergence>
  <variable name = 'x1'>
    <distribution>x1_dst</distribution>
  </variable>
  <variable name = 'x2'>
    <distribution>x2_dst</distribution>
  </variable>
  <variable name = 'gammay'>
    <distribution>gammay_dst</distribution>
  </variable>
</Adaptive>
</Samplers>

<!-- MODELS -->
<Models>
  <ExternalModel name = 'xtr_mdl' subType = '' ModuleToLoad =
    ↪ 'SafestPointPP/safest_point_test_xtr_mdl'>
  <variable>x1</variable>
  <variable>x2</variable>
  <variable>gammay</variable>
  <variable>g</variable>
</ExternalModel>

  <ROM name = 'accelerated_ROM' subType = 'SciKitLearn'>
    <Features>x1, x2, gammay</Features>
    <Target>g_zr</Target>
    <SKLtype>svm|SVC</SKLtype>
    <kernel>rbf</kernel>
    <gamma>10</gamma>
    <tol>1e-5</tol>
    <C>50</C>
  </ROM>

```

```

<PostProcessor name='SP' subType='SafestPoint'>
  <!-- List of Objects (external with respect to this PP)
  ↪ needed by this post-processor -->
  <Distribution class = 'Distributions' type =
  ↪ 'Normal'>x1_dst</Distribution>
  <Distribution class = 'Distributions' type =
  ↪ 'Normal'>x2_dst</Distribution>
  <Distribution class = 'Distributions' type =
  ↪ 'Normal'>gammay_dst</Distribution>
  <!-- end of the list -->
  <controllable>
    <variable name = 'x1'>
      <distribution>x1_dst</distribution>
      <grid type = 'value' steps = '20'>1</grid>
    </variable>
    <variable name = 'x2'>
      <distribution>x2_dst</distribution>
      <grid type = 'value' steps = '20'>1</grid>
    </variable>
  </controllable>
  <non-controllable>
    <variable name = 'gammay'>
      <distribution>gammay_dst</distribution>
      <grid type = 'value' steps = '20'>2</grid>
    </variable>
  </non-controllable>
</PostProcessor>
</Models>

<!-- FUNCTIONS -->
<Functions>
  <External name='g_zr'
  ↪ file='SafestPointPP/safest_point_test_g_zr.py'>
    <variable>g</variable>
  </External>
</Functions>

<!-- OUT-STREAMS -->
<OutputStreamManager>
  <Print name = 'lmt_srf_dmp'>
    <type>csv</type>

```

```
        <source>lmt_srf_dt</source>
    </Print>

    <Print name = 'sfs_pnt_dmp'>
        <type>csv</type>
        <source>sfs_pnt_dt</source>
    </Print>
</OutputStreamManager>

</Simulation>
```

The Python *ExternalModel* is reported below:

```
def run(self, Input):
    self.g = self.x1+4*self.x2-self.gammay
```

The “Goal Function”, the function that defines the transitions with respect the input space coordinates, is as follows:

```
def __residuuumSign(self):
    if self.g<0 : return 1
    else       : return -1
```

## References

- [1] M. P. Forum, “Mpi: A message-passing interface standard,” tech. rep., Knoxville, TN, USA, 1994.
- [2] “Portable batch system.” <http://www.pbsworks.com>.
- [3] J. D. Hunter, “Matplotlib: A 2d graphics environment,” *Computing In Science & Engineering*, vol. 9, no. 3, pp. 90–95, 2007.
- [4] D. Cournapeau, “Scikit-learn library for machine learning.” [http://scikit-learn.org/stable/user\\_guide.html](http://scikit-learn.org/stable/user_guide.html).







