

# *Smooth Emulator & Simplex Sampler*

## **User Manual**

A BAND Collaboration Project  
<https://bandframework.github.io>

Scott Pratt, Eren Erdogan, Ekaksh Kataria  
*Department of Physics and Facility for Rare Isotope Beams*  
*Michigan State University, East Lansing Michigan, 48824*  
September 26, 2023



**MICHIGAN STATE**  
UNIVERSITY



# Contents

<b>1</b>	<b>Overview</b>	<b>3</b>
<b>2</b>	<b>Installation and Getting Started</b>	<b>4</b>
2.1	Prerequisites . . . . .	4
2.2	Making Home Directory and Setting Home Environment Variable . . . . .	4
2.3	Downloading . . . . .	5
2.4	Directory and File Structure . . . . .	5
2.5	Compiling Libraries . . . . .	7
2.6	The Project Directory . . . . .	7
<b>3</b>	<b>Generating Training Points with <i>Simplex Sampler</i></b>	<b>8</b>
3.1	Summary . . . . .	8
3.2	Simplex Parameters (not model parameters!) . . . . .	8
3.3	Specifying Model Parameters and Priors . . . . .	9
3.4	Training Types . . . . .	9
3.4.1	Type 1 . . . . .	9
3.4.2	Type 2 . . . . .	9
3.4.3	Type 3 . . . . .	10
3.5	Running Simplex to Generate Training Points . . . . .	10
<b>4</b>	<b>Performing Full Model Runs</b>	<b>11</b>
<b>5</b>	<b>Tuning the Emulator</b>	<b>12</b>
5.1	Summary . . . . .	12
5.2	<i>Smooth Emulator</i> Parameters (not model parameters!) . . . . .	13
5.3	Editing Info Files . . . . .	15
5.4	Running the <i>Smooth Emulator</i> Program . . . . .	15
<b>6</b>	<b>Emulating Principal Components</b>	<b>18</b>
6.1	Summary . . . . .	18
6.2	PCA Parameters (not model parameters!) . . . . .	18
6.3	Running the PCA programs . . . . .	18

<b>7</b>	<b>MCMC Generation of Posterior</b>	<b>20</b>
<b>8</b>	<b>Template-Based Tutorial</b>	<b>21</b>
8.1	Overview . . . . .	21
8.2	Installation . . . . .	21
8.3	Creating Necessary Info Files . . . . .	22
8.4	Running <i>Simplex Sampler</i> . . . . .	22
8.5	Running the Fake Full Model . . . . .	23
8.6	Running <i>Smooth Emulator</i> . . . . .	24
<b>9</b>	<b>Generating Emulated Observables</b>	<b>26</b>
<b>9</b>	<b>Underlying Theory of <i>Smooth Emulator</i></b>	<b>29</b>
9.1	Mathematical Form of <i>Smooth Emulator</i> . . . . .	29
9.2	Alternate Forms . . . . .	31
9.3	Tuning the Emulator . . . . .	32
<b>10</b>	<b>Theoretical Basis of <i>Simplex Sampler</i></b>	<b>34</b>
10.1	The Pernicious Nature of Step-Function Priors in High Dimension . . . . .	35
<b>11</b>	<b>Tests of the <i>Smooth Emulator</i> using <i>Simplex Sampler</i> for Training</b>	<b>36</b>

# 1 Overview

This manual describes how to install and run *Smooth Emulator* software. The software performs three basic functions. First, the *Simplex Sampler* chooses a set of points in model parameter space, at which full model runs will be performed to then tune the emulator. The user must provide a description of the model parameters and the prior in a text file in a standard format. There are several options, the first of which is to choose the points that represent a simplex, e.g. an equilateral triangle in two dimensions or a tetrahedron in three dimensions. In a simplex, all points are equidistant from one another, and the number of training points is  $N_p + 1$ , where  $N_p$  is the number of parameters. In addition to the standard simplex, there are additional options which are motivated by the simplex form. For the standard form the  $N_p + 1$  training points in the simplex match the number of points needed to determine a linear fit. Another choice, which is based on the simplex, chooses enough points to determine a quadratic fit,  $(N_p + 1)(N_p + 2)/2$ . The software will write the information about the training points in a standard format, which is described in the manual. If the user decides to use training points from a different procedure, the user can still record the information about the points in a same format, and the emulator tuning will still work, as the emulator itself is not predicated on a specific choice of training points.

The user is then responsible for running the full model at the training points and expressing observables, and the uncertainties, for each training point in a standard format. The manual describes the output format.

The second functionality of the software is to build and tune the emulator, referred here as *Smooth Emulator*. The emulator reads the information above, along with another user-provided parameter file to choose which observables are to be emulated, which parameters will be varied, and which emulator options will be applied. After being trained, the Taylor coefficients representing the emulator are written to a file. One can always add additional training points, and retrain the emulator.

The third functionality of the software is to perform a MCMC exploration of parameter space using the emulator. This user must express the experimental observables and their uncertainties in a standard format. The MCMC software will read the emulator coefficients from file and perform the MCMC exploration. This procedure is also guided by a simple text file of parameters. The MCMC software uses python and Matplotlib to generate plots that describe the posterior.

## 2 Installation and Getting Started

### 2.1 Prerequisites

*Smooth Emulator* software should run on UNIX, Mac OS or Linux, but is not supported for Windows OS. *Smooth Emulator* is largely written in C++. In addition to a C++ compiler, the user needs the following software installed.

- git
- CMake
- Eigen3 (Linear Algebra Package)
- GSL (Gnu Scientific Library)
- Python/Matplotlib (only for generating plots in the MCMC procedure)

CMake is an open-source, cross-platform build system that helps automate the process of compiling and linking for software projects. Hopefully, CMake will perform the needed gymnastics to find the Eigen3 and GSL installations. To install CMake, either visit the CMake website (<https://cmake.org/>), or use the system's package manager for the specific system. For example, on Mac OS, if one uses *homebrew* as a package manager, the command is

```
% brew install cmake
```

Eigen is a C++ template library for vector and matrix math, i.e. linear algebra. The user can visit the Eigen website (<https://eigen.tuxfamily.org/dox/>), or use their system's package manager. For example on Mac OS with *homebrew*,

```
% brew install eigen
```

The GNU Scientific Library (GSL) is a numerical library for C and C++ programmers. The library provides a wide range of mathematical routines such as random number generators, special functions and least-squares fitting. There are over 1000 functions in total with an extensive test suite. One can either download the software from the GSL website (<https://www.gnu.org/software/gsl/>), or use a package manager. Again, for *homebrew* on Mac OS,

```
% brew install gsl
```

### 2.2 Making Home Directory and Setting Home Environment Variable

The software requires two repositories. They should be cloned into the same directory. For compilation purposes this directory needs to be accessed via an environmental variable. For example the directory might be named `/Users/CarlosSmith/git_msu`. The User needs to set an environmental variable, `GITHOME_MSU`, to the full path of the directory, e.g.

```
% export GITHOME_MSU="/Users/CarlosSmith/git_msu"
```

It is recommended to copy this command into the user's `.bashrc` (or equivalent) file to avoid re-defining it each time one needs to recompile.

## 2.3 Downloading

From within the `GITHOME_MSU/` directory,

```
GITHOME_MSU% git clone https://github.com/scottedwardpratt/smooth.git
GITHOME_MSU% git clone https://github.com/scottedwardpratt/commonutils.git
```

This creates 2 directories: `.../GITHOME_MSU/commonutils` and `.../GITHOME_MSU/smooth`.

The User needs to create a project directory from which the User would perform most projects. This is easiest accomplished by copying a template from the *Smooth* distribution,

```
% cp -r GITHOME_MSU/templates/myproject MY_PROJECT
```

Hence forth, `MY_PROJECT` will refer to the directory, including the path, from which the User will perform most of the analysis. The User may wish to have several such directories. These directories should be outside the main distribution, i.e. outside the `smooth/` or `commonutils` paths.

Although the main source code, include files and libraries are all located in the software directory, the main programs and executables are not. The motivation for this decision is to allow the User to easily modify their own versions of the main programs. These tend to be very short programs. For that reason there is a separate directory to store the main programs and their executables. The User can easily set this up by copying a template directory,

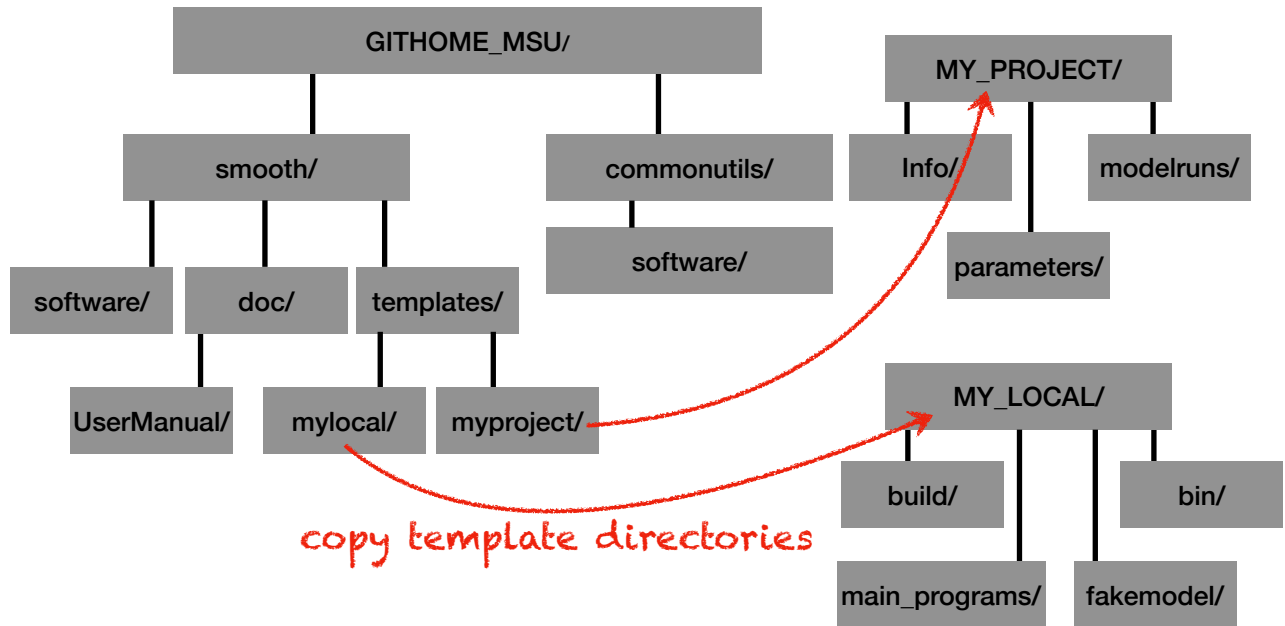
```
% cp -r GITHOME_MSU/templates/mylocal MY_LOCAL
```

Here `MY_LOCAL` will hence forth refer to the path of this directory. This directory should be outside the main distribution, i.e. outside the `smooth/` or `commonutils` paths.

For the remainder of this manual, `GITHOME_MSU`, `MY_LOCAL` and `MY_PROJECT` will be used to denote the location of these directories.

## 2.4 Directory and File Structure

Once compiled, the libraries in the `commonutils/` directory are used for a variety of tasks. These libraries are not particularly designed for *Smooth Emulator* or *Simplex Sampler*. The `smooth/` directory contains codes that are used to create libraries specific to the sampler and emulator. The executables are stored in `MY_LOCAL/bin`. The short main program source files are located in `MY_LOCAL/main_programs/`. It is not envisioned that the User would edit files in the `smooth/software/` directory, but that the User may well wish to create custom versions of the short main programs in `MY_LOCAL/main_programs/`. The main programs are compiled using the CMake files in `MY_LOCAL/build/`. The User may find it convenient to add `MY_LOCAL/bin/` to their path.



**Figure 2.1: The directory structure:** The User clones two repositories into some location, which will be referred to as **GITHOME.MSU**. The User can then copy two template directories into the User's choices of locations outside the path of the repositories. The name of those two directories will be referred to as **MY\_LOCAL**, which will contain the main programs and executables, and **MY\_PROJECT** which contains the data files related to a given project. The programs are designed to be run from within the **MY\_PROJECT/** directory.

## 2.5 Compiling Libraries

First, change into software directories, then create the makefiles with cmake, then compile them.

```
% cd GITHOME_MSU/commonutils/software
GITHOME_MSU/commonutils/software% cmake .
GITHOME_MSU/commonutils/software% make
GITHOME_MSU/commonutils/software% cd ../GITHOME_MSU/smooth/software
GITHOME_MSU/smooth/software% cmake .
GITHOME_MSU/smooth/software% make
```

At this point all the libraries are built, but this does not include the main programs. The main programs are short, and are located in a separate location, as they are meant to serve as examples which the User might copy and edit at will.

Finally, compile the main programs. Below, this illustrates how to build the programs used for generating training points with Simplex and for tuning the emulator with Smoothy:

```
% cd MY_LOCAL/build
MY_LOCAL/build% cmake .
MY_LOCAL/build% make simplex
MY_LOCAL/build% make smoothy_tune
.
.
```

Other source codes for main programs can be found in `MY_LOCAL/main_programs/`. If you build your own main programs (probably using these as examples), you can edit the `CMakeList.txt` file in `GITHOME_MSU/smooth/local/build`, using the existing entries as an example. The executables should appear in `MY_LOCAL/bin/`.

## 2.6 The Project Directory

Within `MY_PROJECT/` there are three sub-directories (assuming it was created from the template). The first is `MY_PROJECT/Info/`. Information about the model parameters, and their priors is stored in `MY_PROJECT/Info/modelpar_info.txt`, and information about the observables is stored in `MY_PROJECT/observable_info.txt`. The `MY_PROJECT/parameters` directory stores user-defined parameter files used by *Simplex Sampler*, `MY_PROJECT/parameters/simplex_parameters.txt`, and by *Smooth Emulator* `MY_PROJECT/parameters/emulator_parameters.txt`. The `MY_PROJECT/modelruns` directory will store information for each full-model run. The directories `MY_PROJECT/modelruns/run0/`, `MY_PROJECT/modelruns/run1/`,  $\dots$ , have files describing the model parameters for each run, along with the output required by the emulator for each specific full-model run. For example, the `MY_PROJECT/modelruns/run1/` directory has the files `mod_parameters.txt` and `obs.txt`. The first file stores the model parameter values for that particular training run. The User then runs their full model based on those parameters and stores the corresponding observables in `obs.txt`. The User may generate the `mod_parameters.txt` files using *Simplex Sampler*, or the user might generate them according to some other prescription. Once the User has then generated the `obs.txt` files, *Smooth Emulator* can then build and tune the emulator.



## 3 Generating Training Points with *Simplex Sampler*

### 3.1 Summary

*Simplex Sampler* produces a list of points in the  $n$ —dimensional model-parameter space to be used for training an emulator. The algorithms are based on the  $n$ —dimensional simplex. For example, in two dimensions the points are arranged in an equilateral triangle, and for three dimensions of parameters points are arranged in a tetrahedron. The program first reads a simple text file that provides the names of model parameters and the range of their prior distribution. Options for *Simplex sampler* are taken from a separate text file. This enables the User to make choices, such as which algorithm to apply when generating the training points. *Simplex Sampler* determines the number of training points based on the algorithm. The User is free to run the full model at additional points.

### 3.2 Simplex Parameters (not model parameters!)

These are parameters representing choices made by the User. Note these are NOT the model parameters, which are then generated by *Simplex Sampler*. If one visits the User’s project directory, these parameters are typically stored in the file MY\_PROJECT/parameters/simplex\_parameters.txt, where the path is either absolute or relative to the project directory. Parameters files can have any name or location. These files are text files in the format. An example of a parameter file is:

```
#Simplex_LogFileName      simplexlog.txt      # if blank, output to screen
Simplex_TrainType         1                          # Must be 1 or 2 or 3
Simplex_RTrain            0.95                      # Radius of simplex
Simplex_ModelRunDirName   modelruns                # Directory with training
                                                point information
```

For the parameter file, the first string is the parameter name and is followed by the value. Both are single strings (without spaces). The # symbol is used for comments. Each parameter has a default value, which will be used if the parameter is not mentioned in the parameter file. *Simplex Sampler* has four User-defined parameters.

1. **Simplex\_TrainType**

Possible values are “1”, “2” or “3”. The default, “1”, will position points according to a simplex, i.e. in two dimensions this is an equilateral triangle and in three dimensions, it is a tetrahedron. In  $n$  dimensions there are  $n + 1$  points separated at equal distances from one another and centered at the origin.

2. **Simplex\_RTrain**

This sets the distance from the origin that the training points will be selected throughout the parameter space. The default is 0.95.

3. **Simplex\_ModelRunDirName**

This sets the path to the directory in which the run files will be created. The default name is `modelruns`, but the user can change this to anything they want. The path is relative to the project directory, i.e. the directory from which you run the *simplex* command.

#### 4. Simplex\_LogFileName

If this is left blank, *Simplex Sampler* will write output to the string. Otherwise it will write output to a file. Given that *Simplex Sampler* runs in a few seconds, the program is usually run interactively and output is sent to the screen.

### 3.3 Specifying Model Parameters and Priors

Before proceeding, Simplex requires information about the parameters, specifically, their ranges. The User enters this information into the file `Info/modelpar_info.txt`. An example of such a file might be

NuclearCompressibility	gaussian	210	40
ScreeningMass	uniform	0.3	1.2
Viscosity	uniform	0.08	0.3

The first column is the model-parameter name, and the last three parameters describe the range of the parameters, which is usually the prior, assuming the prior is uniform or Gaussian. The second entry for each parameter defines whether the range/prior is **uniform** or **gaussian**. If the prior is **uniform**, the next two numbers specify the lower and upper ranges of the parameter. If the range/prior is **gaussian**, the third entry describes the center of the Gaussian,  $\mathbf{x}_0$ , and the fourth entry describes the Gaussian width,  $\sigma_0$ , where the prior distribution is  $\propto \exp\{-(\mathbf{x} - \mathbf{x}_0)^2 / 2\sigma_0^2\}$ . Simplex will read the information to determine the number of parameters. It will then assign the  $n$  points,  $\theta_{1...n}$  assuming each dimension of  $\theta$  varies from -1 to 1, for uniform distributions, or proportional to  $e^{-\theta^2/2}$  for Gaussian distributions. The points  $\theta_i$  are each then converted into  $\mathbf{x}_i$  by scaling and translating the values according to the ranges/priors defined in the `modelpar_info.txt` file.

### 3.4 Training Types

#### 3.4.1 Type 1

Depending on the number of parameters,  $n$ , the program creates a simplex in  $n$  dimensions. This simplex's vertices will be used to generate  $N_{\text{train}} = n + 1$  training points. These points will be scaled by different values so the training points aren't in the same radius. This results in the minimum number of required points for linear fits. Thus, if the model is perfectly linear, this option provides perfect emulation.

#### 3.4.2 Type 2

Depending on the number of parameters, the program first creates a simplex in  $n$  dimensions. This simplex's vertices will be used to generate new training points there and along the edges. These points will be scaled to be in different radii from the center. This results in the minimum number of required points for quadratic fits. The net number of training points is then  $N_{\text{train}} =$

$n + 1 + n(n + 1)/2$ . Thus, if the model is perfectly quadratic, this option provides perfect emulation.

### 3.4.3 Type 3

This training type creates two different simplexes, both centered at the origin. The second simplex is a reflection of the first. The 2-dimensional visualization of this would look like the “Star of David”. Finally, one extra training point is added at the origin. The net number of training points is  $N_{\text{train}} = 2n + 3$ .

## 3.5 Running Simplex to Generate Training Points

To run *Simplex Sampler*, first make sure the program is compiled. To compile the programs, change into the `MSU_GITHOME/smooth/local/build/` directory and enter the following command,

```
GITHOME_MSU/local/build% cmake .
GITHOME_MSU/local/build% make simplex
```

Next, change into your project directory and run the program.

```
MY_PROJECT% MY_LOCAL/bin/simplex PARAMETER_FILE_NAME
```

Here `MY_LOCAL` is the path to where the User compiles the main programs into executables, and `PARAMETER_FILE_NAME` is the name of the simplex parameter file, e.g. `parameters/simplex.parameters.txt`. Paths are either absolute, or are relative to the project directory.

Simplex will write the information about the training points in the directory defined by the `Simplex_ModelRunDirName` parameter. Within the directory, a sub-directory will be created for each training point, named `run0/`, `run1/`, `run2/...`. Within each subdirectory, Simplex creates a file `Simplex_ModelRunDirName/runI/mod_parameters.txt` for the  $I^{\text{th}}$  training point. For example, the `run0/mod_parameters.txt` file might be

NuclearCompressibility	229.08
ScreeningMass	0.453
Viscosity	0.192

At this point, it is up to the User to run their full model at each training point and create a file `runI/obs.txt`, which stores values of the observables at those training points as calculated by the full model.

## 4 Performing Full Model Runs

Once the training points are generated, the User can run the full model for each of the training points. At this point there is a directory, usually called `MY_PROJECT/modelruns/`, in which there are sub directories, `run0/`, `run1/`, `run2/...`. Within each sub-directory, `runI`, there should exist a text file `MY_PROJECT/modelruns/runI/mod_parameters.txt`. These files could have been generated by *Simplex Sampler*, but could have been generated by any other means, including by hand. The files should be of the form,

```
par1_name  par1_value
par2_name  par2_value
par3_name  par3_value
.
.
```

The parameter names must match those defined in `MY_PROJECT/Info/modelpar_info.txt`, the format of which is described in Sec. 3.

The User must then perform full model runs using the model-parameter values as defined in each sub-directory. The full model runs then need to produce results and write a list of observable values in each run directory. Each file must be named `MODEL_RUN_DIRNAME/runI/obs.txt`. The directory `MODEL_RUN_DIRNAME/`. Typically this directory is `MY_PROJECT/moderundirs`, but that can be changed by the User. The format of those text files should be

```
observable1_name  observable1_value  observable1_random_uncertainty
observable2_name  observable2_value  observable2_random_uncertainty
observable3_name  observable3_value  observable3_random_uncertainty
.
.
```

The names must match those listed in `MY_PROJECT/Info/observable_info.txt`, which will be used by *Smooth Emulator*, as described in Sec. 5. The values are the observable values as calculated by the full model for the model-parameter values listed in the corresponding `mod_parameters.txt` file in the same directory. The random uncertainties refer only to those uncertainties due to noise in the full model. Random noise is that, which if the full model would be rerun at the same model-parameter values, would represent the variation in the observable values. In most cases this would be set to zero. But, if the full model has some aspect of sampling to it, for example generating observables from event generators with a finite number of events, that variation should be listed here. This variation is required for the emulator. If there is such a variation, the emulator should not be constrained to exactly reproduce the training point observables at the training points. The principal danger being, that if two training points are very close to one another, but with a finite fluctuation, exactly producing the training points might require very high slopes to exactly reproduce the training points.

Once the observable files are produced for each of the full model runs, the User can then proceed to build and tune an emulator using *Smooth Emulator*.

## 5 Tuning the Emulator

### 5.1 Summary

Smooth emulator finds a sample set of Taylor expansion coefficients that reproduce a set of observables at a set of training points. For a given observables, a particular sample set of coefficients gives the following emulated function:

$$E(\vec{\theta}) = \sum_{\vec{n}, s.t. \sum_i n_i \leq \text{MaxRank}} \frac{d(\vec{n})}{(n_1 + n_2 + \dots)!} A_{\vec{n}} \left(\frac{\theta_1}{\Lambda}\right)^{n_1} \left(\frac{\theta_2}{\Lambda}\right)^{n_2} \dots \quad (5.1)$$

Here,  $\theta_1 \theta_2 \dots$  are the model parameters, scaled so that their priors range from -1 to +1, or if they are Gaussian, have unit variance. The degeneracy factor,  $d(\vec{n})$  is the number of different ways to sum the powers  $n_i$  to a given rank,

$$d(\vec{n}) = \frac{n_1! n_2! \dots}{(n_1 + n_2 + \dots)!} \quad (5.2)$$

The emulator finds a predetermined number of sets of coefficients, where each set of coefficients provides a function that reproduces the real model at the training points. The User sets the number of sets of coefficients, typically of order  $N_{\text{sample}} \approx 10$ , in a parameter file. Away from the training points, the uncertainty of the emulator is represented by the spread of the values amongst the  $N_{\text{sample}}$  predictions.

*Smooth Emulator* solves for the  $N_{\text{sample}}$  sets of coefficients from the training data, then stores those coefficients in files for later use. *Smooth Emulator* can emulate either the full-model observables directly, or their principal components (PCA capability coming soon). Training the emulator follows the same steps for either approach.

The executables based on *Smooth Emulator* are located in the User's `MY_LOCAL/bin` directory. Examples of such executables are `smoother_tune` or `smoother_calcobs`. These functions must be executed from within the User's project directory.

Before training the emulator, one must first run the full model at a given set of training points. In addition to a parameter file, which sets numerous options, the User must provide the following:

1. A file listing the names of observables, their uncertainties, and an estimate of the variance of each observable throughout the model-parameter space. This file is named `Info/obs.txt`, where the path is relative to the project directory.
2. If the number of full-model runs performed is  $N_{\text{train}}$ , *Smooth emulator* requires files for each run. Each file is named `MODEL_RUN_DIRNAME/runI/mod_parameters.txt`, where  $I$  varies from  $1 - N_{\text{train}}$ , describes the point in parameter space for the  $I^{\text{th}}$  full-model run. The directory `MODEL_RUN_DIRNAME/` is typically `MY_PROJECT/modelruns`, but can be defined otherwise (see below).
3. In the same directory, *Smooth Emulator* requires the observables calculated at the training points mentioned above. This information is provided in `MODEL_RUN_DIRNAME/runI/obs.txt`.

The parameter file, typically stored in `parameters/emulator_parameters.txt`, enables the User to select numerous options. For example, the User might use training data from a different directory, not `modelruns/`, or might choose to use principal components rather than the observables directly. In the following subsections, we first review the format for each of the required input files, then describe how to run *Smooth Emulator*, how its output is stored, and how to switch PCA observables for real observables.

## 5.2 *Smooth Emulator* Parameters (not model parameters!)

*Smooth Emulator* requires a parameter file. This can be located anywhere, as it will be specified on the command line when running *Smooth Emulator*, but is typically `parameters/emulator_parameters.txt`. The parameter file is simply a list, of parameter names followed by values.

```
#SmoothEmulator_LogFileName      smoothlog.txt
SmoothEmulator_LAMBDA            3.0
SmoothEmulator_MAXRANK           4
SmoothEmulator_NMC               100000 # Steps between samples
SmoothEmulator_NASample          8      # No. of coefficient samples
SmoothEmulator_TuneChooseMCMC    true   # set true if NPars>~5
SmoothEmulator_UseSigmaYRreal    false
SmoothEmulator_ConstrainAO       true
SmoothEmulator_CutoffA           false
SmoothEmulator_ModelRunDirName   modelruns
SmoothEmulator_TrainingPts       0-9
SmoothEmulator_UsePCA            false
RANDY_SEED                      1234
```

If any of these parameters are missing from the parameters file, *Smooth Emulator* will assign a default value.

### 1. **SmoothEmulator\_LAMBDA**

This is the smoothness parameter  $\Lambda$ . It sets the relative importance of terms of various rank. If  $\Lambda$  is unity or less, it suggests that the Taylor expansion converges slowly. The default is 3.

### 2. **SmoothEmulator\_LogFileName**

If this is commented out, as it is in the example above, *Smooth Emulator*'s main output will be directed to the screen. Otherwise, the output will be recorded in the specified file.

### 3. **SmoothEmulator\_MAXRANK**

As *Smooth Emulator* assumes a Taylor expansion, this the maximum power of  $\theta^n$  that is considered. Higher values require more coefficients, which in turn, slows down the tuning process. The default is 4.

### 4. **SmoothEmulator\_TuneChooseMCMC**

If set to `false`, *Smooth Emulator* will set all but  $N_{\text{train}}$  coefficients randomly, according to their Gaussian prior. Then, it will solve for the remaining coefficients in order to fit the

training data. The weight is calculated for the remaining coefficients, at which point the coefficients are kept or rejected proportional to the weight. The coefficients chosen in this manner are perfectly independent of one another, but perhaps at the cost of requiring many samplings before finding a weight to keep. This choice is efficient when the number of training points is small. If `SmoothEmulator_SmoothEmulator_TuneChooseMCMC` is set to `true`, *Smooth Emulator* will choose the coefficients as a small random step from the previous coefficients, then keep or reject the coefficients according to a Metropolis algorithm. The downside is that many steps are required to create a sampling set of coefficients that are independent of one another. Although it can be slow when there are many model parameters, this choice is more efficient for larger numbers of training points. The default is `true`.

5. **SmoothEmulator\_NMC**

When the previous parameter is set to `true`, this sets the number of steps between retained samples of coefficients. For larger numbers of parameters, this should be set at many thousands. Higher values lead to more independent sets of coefficients, but the calculation then requires more time.

6. **SmoothEmulator\_NASample**

*Smooth Emulator* finds  $N_{\text{sample}}$  sets of coefficients. Each set reproduces the training points, but differs away from the training points. Setting  $N_{\text{sample}} \sim 10$  should reasonably represent the uncertainty of the emulator. The default is set at 8.

7. **SmoothEmulator\_UseSigmaYRreal**

If the real model has noise, the emulator should not be constrained to exactly reproduce the observables at the training points. In fact, if two training points are located close to one another in parameter space, *Smooth Emulator* might be forced to find a particularly uneven function so that the points are exactly reproduced. If the User wishes to exactly reproduce the training points, this should be set to `false`, as is the default.

8. **SmoothEmulator\_ConstrainA0**

The coefficients in the Taylor expansion are assumed to have some weight,

$$W(A_i) = \frac{1}{\sqrt{2\pi\sigma_A^2}} e^{-A_i^2/2\sigma_A^2}.$$

The term  $\sigma_A$  is allowed to vary during the tuning to maximize the likelihood of the expansion. If the User wishes to exempt the lowest term, i.e. the constant term in the Taylor expansion from the weight, the User may set `SmoothEmulator_ConstrainA0` to `false`. The default is `false`.

9. **SmoothEmulator\_CutoffA**

This applies an additional multiplicative weight to the weight for  $A$  above.

$$W(A_i)_{\text{additional}} = \frac{1}{1 + \frac{1}{4} \frac{A_i^2}{\sigma_A^2}}.$$

Here  $\sigma_{A0}$  is the initial guess for the spread. This can safeguard against the width  $\sigma_A$  drifting off to arbitrarily large values. Unless necessary, it is recommended to leave this at the default, `false`.

#### 10. **SmoothEmulator\_ModelRunDirName**

This gives the directory in which the training data from the full model runs is stored. The default is `modelruns`, which is the same default `Simplex Sampler` uses for writing the coordinates of the training points.

#### 11. **SmoothEmulator\_TrainingPts**

This lists which full-model training runs `SmoothEmulator` will use to train the emulator. This provides the User with the flexibility to use some subset for training, as may be the case when testing the accuracy. The default is “1”. An example the User might enter could be

`SmoothEmulator_TrainingPts 0-4,13,15`

This would choose the training information from the directories `run0`, `run1`, `run2`, `run3`, `run4`, `run13` and `run14`, which would be found in the directory denoted by the `SmoothEmulator_ModelRunDirName` parameter.

#### 12. **RANDY\_SEED**

This sets the seed for the random number generator. If the line is commented out, it will be set to `std::time(NULL)`.

#### 13. **SmoothEmulator\_UsePCA**

By default, this is set to false. If one wishes to emulate the PCA observables, i.e. those that are linear combinations of the real observables, this should be set to true. One must then be sure to have run the pca decomposition programs first. For more, see Sec. 6.

### 5.3 Editing Info Files

#### 1. **Info/observable\_info.txt**

*Smooth Emulator* requires knowledge of the observables. An example of such a file is

```
meanpt_pion 40
meanpt_proton 60
meanv2_pion 0.05
```

The second line provides an initial estimate for the parameter  $\sigma_A$ . The User needn't worry if this is off by a few factors of two from the final value, but if it is off by orders of magnitude, it might take *Smooth Emulator* a long time to find the appropriate range.

#### 2. **Info/modelpar\_info.txt**

This file provides the names and ranges of the model parameters, i.e. the prior. *Smooth Emulator* translates the scales the parameters so that they have uniform ranges, in the case of uniform distributions, or uniform widths, and zero mean for the Gaussian distributions. This same file was used for running *Simplex Sampler* and is described in Sec. 3.3.

### 5.4 Running the *Smooth Emulator* Program

The source code for the *Smooth Emulator* main program can be found in the `MY_LOCAL/main_programs/` directory. This directory contains source code for several main programs. They are separated from



the bulk of the software, which is in the `GITHOME_MSU/smooth/software/` directory. The main programs are designed so that the User can easily copy and edit them to create versions that might be more appropriate to the User's specific needs. When compiled, from the `MY_LOCAL/build/` directory, the executables appear in the `MY_LOCAL/bin/` directory. Two of the source codes that come with the distributions are `MY_LOCAL` and `MY_LOCAL/main_programs/smoothy_calc_main.cc`. Once compiled the corresponding executables are `MY_LOCAL/bin/smoothy_tune` and `MY_LOCAL/bin/smoothy_calcobs`.  
A

```
using namespace std;
int main(int argc, char *argv[]){
    if(argc!=2){
        printf("Usage smoothy emulator parameter filename");
        exit(1);
    }
    CparameterMap *parmap=new CparameterMap();
    parmap->ReadParsFromFile(string(argv[1]));
    CSmoothMaster master(parmap);
    master.ReadTrainingInfo();
    master.GenerateCoefficientSamples();
    master.WriteCoefficientsAlly();
    return 0;
}
```

Similarly, there is a code `MY_LOCAL/main_programs/smoothy_calcobs_main.cc`, which provides an example of how one might read the coefficients and generate predictions for the emulator at specified points in parameter space.

From within the `MY_LOCAL/build/` directory, one can compile the two programs with the commands:

```
MY_LOCAL/build % cmake .
MY_LOCAL/build % make smoothy_tune
MY_LOCAL/build % make smoothy_calcobs
```

The executables `smoothy_tune` and `smoothy_calcobs` should now appear in the `MY_LOCAL/bin/` directory. Assuming the `bin/` directory has been added to the User's path, the User may switch to the User's project directory, and enter the command

```
~/MY_PROJECT % smoothy_tune PARAMETERS/MY_PARAMETERS.TXT
```

Here `PARAMETERS/MY_PARAMETERS.TXT` can be replaced by the User, but is typically `parameters/emulator_par`

The program will write the Taylor coefficients for the  $N_{\text{sample}}$  samples to files in the `coefficients` directory. The coefficients for each observable are given in separate subdirectories, named by the observables, i.e. `coefficients/OBS_NAME/sampleI.txt`. Here,  $I$ , where `OBS_NAME` is the name for each observable, and if there are  $N_{\text{sample}}$  sets of coefficients,  $0 \leq I < N_{\text{sample}}$ . Along with the coefficients, in the same directory *Smooth Emulator* writes a file for each observable. These files are

named `coefficients/OBS_NAME/meta.txt`. This file provides information, such as the maximum rank and net number of model parameters, to make it possible to read the coefficients later on.

*Smooth Emulator* will output lines describing its progress, either to the screen or to a file, as specified by the `SmoothEmulator_LogFile` parameter described above. This output includes a report on the percentage of steps in the MCMC program that were successful. The line `BestLogP/Ndof` describes the weight used to evaluate the likelihood of a coefficients sample. This value should roughly plateau once the Metropolis procedure has settled on the most likely region.

For later us, e.g. when performing the MCMC to sampler the posterior, the User would need to generate predictions for specified values of the parameters. The executable `MY_LOCAL/bin/smoothy_calcobs`, is such an example. It is compiled from the main program, `MY_LOCAL/main_programs/smoothy_calcobs.cc`:

```
int main(int argc,char *argv[]){
    if(argc!=2){
        CLog::Info("Usage smoothy_calcobs emulator parameter filename");
        exit(1);
    }
    CparameterMap *parmap=new CparameterMap();
    parmap->ReadParsFromFile(string(argv[1]));
    CSmoothMaster master(parmap);
    // Reads Emulator Coefficients for all observables
    master.ReadCoefficientsAlly();
    master.priorinfo->PrintInfo();
    //modpars carries info about single point
    CModelParameters *modpars=new CModelParameters(master.priorinfo);
    // Prompt user for model parameter values
    vector<double> X(modpars->NModelPars);
    for(int ipar=0;ipar<modpars->NModelPars;ipar++){
        cout << "Enter value for " << master.priorinfo->GetName(ipar) << ":\n";
        cin >> X[ipar];
    }
    modpars->SetX(X);
    // Calc Observables Y[iy] for X
    CObservableInfo *obsinfo=new CObservableInfo("Info/Observable_Info.txt");
    vector<double> Y(obsinfo->NObservables);
    vector<double> SigmaY(obsinfo->NObservables);
    master.CalcAlly(modpars,Y,SigmaY);
    for(int iY=0;iY<obsinfo->NObservables;iY++){
        cout << obsinfo->GetName(iY) << " = " << Y[iY] << " +/- " << SigmaY[iY] << endl;
    }
    return 0;
}
```

The User can hopefully use this template programs as a base for calling *Smooth Emulator* to calculate the emulator values for a specified point in space. Note that `SigmaY` is the emulator uncertainty, not that from experiment or from the theoretical model.

## 6 Emulating Principal Components

This section is incomplete, and the PCA software is not yet fully tested.

### 6.1 Summary

Rather than emulating all observables, it can be more efficient to emulate a handful of principal components. After generating the training-point data, one can run the `pca` program included with the distribution. This will create files that shadow those used to emulate the observables. This will create a file `Info/pca_info.txt` alongside `Info/observable_info.txt`. The difference is that the observables will be named `z1, z2, ...`. In each run directory, alongside the `obs.txt` files, there will be a `obs_pca.txt` file. Finally, there will be a file `PCA_Info/tranformation_info.txt` file that contains all the information and matrices required to perform the basis transformation. If the parameter `Use_PCA` is set to `true`, the emulator will use the PCA files above instead of the observable files. The emulator will then store the Taylor coefficients in the directory `coefficients_pca/` rather than in `coefficients/`.

To get an idea of the capabilities and functionality of the PCA elements of the *Smooth Emulator* Distribution, one can view the sample main program, `GITHOME_MSU/smooth/local/main_programs/pca_main.cc`.

### 6.2 PCA Parameters (not model parameters!)

The PCA programs uses parameter that are prefixed with **SmoothEmulator**. One would typically use the same parameter file as used for running *Smooth Emulator*. The relevant parameters are:

1. **SmoothEmulator\_UsePCA**

If one wishes to emulate the PCA observables, i.e. those that are linear combinations of the real observables, this should be set to `true`. One must then be sure to have run the PCA decomposition programs first.

2. **SmoothEmulator\_ModelRunDirName** and **SmoothEmulator\_TrainingPts** should be set the same as used by *Smooth Emulator*.

### 6.3 Running the PCA programs

The first sample program is `pca_calctransformation`, which reads the training information from the full model runs and calculates the principal component information. Quantities such as the PCA eigenvalues and eigenvectors are stored. This provides the User knowledge of which linear combination of observables carry significant resolving power. By storing the eigenvectors, the information may be retrieved later. This allows the User to easily transform from the observable,  $y_a$ , to the PCA components  $z_a$ .

One should first to `GITHOME_MSU/local/build` and compile/install the program `GITHOME_MSU/local/bin/pca_calctransformation`.

```
.../local/build % cmake .  
.../local/build % make pca_calctransformation
```

Next, from the project directory (assuming the training point information has already been collected) one can enter the command (assuming the path includes GITHOME\_MSU/local/bin)

```
../my_project/ % pca_calctransformation PARAMETER_FILENAME
```

Here, `PARAMETER_FILENAME` is likely `parameters/emulator_parameters.txt`. At this point, all the information about observables from the training has equivalent representation for the PCA components.

In order to emulate the PCA components, one must set the parameter `SmoothEmulator_UsePCA` to true. Then, running the program `smoother_tune` as described above will build and tune an emulator for the PCA components. It will store the Taylor coefficients in the directory `coefficients_pca/`.

The second sample program reads the transformation information written by `pca_calctransformation`. This program gives an example of transforming a vector of principal components,  $\mathbf{z}_a$ , to a vector of observables  $\mathbf{y}_a$ . To compile the programs, change into the build directory as above and enter:

```
.../local/build % cmake .  
.../local/build % make pca_readtransformation
```

## 7 MCMC Generation of Posterior

An MCMC module will be added soon. Hopefully, it is not too difficult for the User to incorporate the emulator into a third-party MCMC program. An MCMC code specifically written to incorporate *Smooth Emulator* should be ready by Spring of 2024.

## 8 Template-Based Tutorial

### 8.1 Overview

A template project directory is provided that the User may copy to their own space, then use this as a foundation from which to embark on their own analysis. This directory includes information files, describing the parameter priors and the observables, that correspond to an artificial model that is also provided as a template. Working through the steps in this section constitutes a tutorial, both for running *Simplex Sampler* and for running *Smooth Emulator*.

This section describes the steps of how the User would

1. Copy the required files from the template directory to the User's space, and compile the required programs.
2. Set up the information files describing the priors and observable names.
3. Run *Simplex Sampler* to generate the model-parameter values at which the full model will be trained.
4. Run an artificial full model to generate the observables for each of the full-model runs.
5. Tune *Smooth Emulator* and write the coefficients to file.
6. Run a program that prompts the User for the coordinates of a point in parameter space, then returns the emulator prediction with its uncertainty.

### 8.2 Installation

After completing the necessary prerequisites listed in section 2[Installation] and following the steps outlined in section 2[Prerequisites] to install the required cmake, eigen, and gsl libraries, and setting the Home Environment Variable by creating the Home Directory as described in section 2[Making Home Directory and Setting Home Environment Variable], the user must proceed to clone the smooth and commonutils directories and compile the libraries, as explained in sections 2[Downloading] and [Compiling Libraries].

Then, the user can establish a personalized project directory by duplicating the project\_template directory onto their computer. The User needs to copy two directories, GITHOME\_MSU/smooth/templates/mylocal and GITHOME\_MSU/smooth/templates/myproject to locations in their personal space. We will refer to the User's two new directories as MY\_LOCAL/ and MY\_PROJECT/. For the purpose of this tutorial, the User must compile three programs. This requires first changing into the MY\_LOCAL/build/ directory and entering:

```
MY_LOCAL% cmake .
MY_LOCAL% make simplex
MY_LOCAL% make smoothy_tune
MY_LOCAL% make smoothy_calcobs
```

## 8.3 Creating Necessary Info Files

Before a User can run *Simplex Sampler* they must create information files that describe the model-parameter priors and list the observable names. Both files are in the MY\_PROJECT directory. The first file is MY\_PROJECT/Info/prior\_info.txt. For the purposes of this tutorial, a file already exists,

```
par1 uniform      0.000      100.0
par2 uniform      0.000      100.0
par3 uniform      0.000      100.0
par4 uniform      0.000      100.0
```

This implies that the model has four parameters. The names, without much inspiration, are `par1`, `par2`, `par3` and `par4`. These names would normally be more descriptive, e.g. `NuclearCompressibility`. The second entry in each line is either `uniform` or `gaussian`. If the parameter is `uniform`, the last two numbers represent the range of the uniform prior,  $x_{\min}$  and  $x_{\max}$ . If the second entry is `gaussian` the third entry represents the center of the Gaussian distribution and the fourth represents the width. For a real model, the User would replace this model with one appropriate for their own model.

The second file is MY\_PROJECT/Info/observable\_info.txt. This describes output values from the model. In the template the file is

```
length  meters    10.0
mass     kg        30.0
time     s         25.0
```

The first entry in each line simply provides the names of the observable which will be processed in the Bayesian analysis. The second lists the units, and is never used, but is included for reference. The third entry is only used by `Smooth Emulator` during tuning. The spread,  $\sigma_A$ , describes the variance of the output due to a single term in the Taylor expansion. Because this is treated as a random variable, and varied from one sampling of the coefficients to another, one must have an initial choice for it before the MCMC procedures in the tuning can proceed. One simple needs to choose this value within a few factors of two from the optimized value. Otherwise, the initial burn-in stage of the MCMC might require more time to coverge in the right neighborhood. Once the burn-in of the MCMC procedure is completed, this parameter is not used.

## 8.4 Running Simplex Sampler

Both *Simplex Sampler* and `Smooth Emulator` have options. These are provided in parameter files. For this tutorial, the provided parameter file is MY\_PROJECT/parameters/simplex\_parameters.txt. The provided file is

```
#Simplex_LogFileName      simplexlog.txt # comment out to direct output to screen
Simplex_TrainType          1                # Must be 1 or 2 or 3
Simplex_RTrain             0.95              # Radius of simplex
Simplex_ModelRunDirName    modelruns         # Directory with training pt. info
```

Because the first line is commented, the output of *Simplex Sampler* will be to the screen. Otherwise it would go to the specified file. By setting `Simplex_TrainType=1`, the sampler will choose  $n + 1$  training points, where  $n = 4$  is the number of model parameters. Each point corresponds to the vertices of an  $n + 1$  dimensional simplex. The parameter `Simplex_RTrain` sets the radius of the simplex, and needs to be less than one for a uniform distribution. The value of 0.95 ensures that the 5 training points are far from one another. Finally, the parameter `Simplex_ModelRunDirName` is set to “modelruns”. This informs *Simplex Sampler* to write the coordinates of each training point and the corresponding observables in the directory `MY_PROJECT/modelruns/`.

Now the user can run *Simplex Sampler*, which must be run from the project directory,

```
~/MY_PROJECT% MY_LOCAL/bin/simplex parameters/simplex_parameters.txt
```

If all goes well there is no screen output. The program writes information about the training points in the `modelruns/` directory. Changing into that directory, there should now be five directories, corresponding to the training points at 5 places: `modelruns/run0`, `modelruns/run1`, `modelruns/run2`, `modelruns/run3`, and `modelruns/run4`. Each directory has one text file describing the training points. For example, the `modelruns/run0/mod_parameters.txt` file is

```
par1 18.3772
par2 31.7426
par3 37.0901
par4 40
```

This describes the four model parameters, which will serve as the input for the first full model run. If one had set the parameter `Simplex_TrainType` to 2, there would be 15 sub-directories in `modelruns/`. The next step will be to run the full model for the parameters in each directory. Thus for `Simplex_TrainType=1`, one would need 5 full-model runs, and for `Simplex_TrainType=2`, one would need to do 15 full-model runs. The corresponding observables will be written in the files `modelruns/runI/obs.txt`

## 8.5 Running the Fake Full Model

Once the training points have been generated, the user will input a Real full model based on the given structure, tailored to address their specific problem. For the tutorial, a fake model is provided. It reads the model-parameter values in each `modelruns/runI/mod_parameters.txt` file and writes the corresponding observables in `modelruns/runI/obs.txt`.

The fake model is a simple python script, which is enacted by entering the command:

```
MY_PROJECT% python3 MY_LOCAL/fakemodel/templatemod.py
```

The output should appear as

```
MY_PROJECT% python3 MY_LOCAL/fakemodel/templatemod.py
```



```

MY_PROJECT% Writing: length 10.871583624279312 1.0
MY_PROJECT% Writing: mass -48.906793378963606 1.0
MY_PROJECT% Writing: time -50.76565279455246 1.0

MY_PROJECT% Writing: length -27.67701646408149 1.0
MY_PROJECT% Writing: mass -50.35567713223607 1.0
MY_PROJECT% Writing: time -34.04264570717603 1.0
.
.

```

Inspecting the modelruns/run0/obs.txt file,

```

length 10.871583624279312 1.0
mass -48.906793378963606 1.0
time -50.76565279455246 1.0

```

The second entry of each line is the value of the specified observable for that specific training point. The last entry is the random uncertainty associated with the full model. This is only relevant if the model has random fluctuations, meaning the re-running the model at the same point might result in different output. For this tutorial, the emulator will not consider such fluctuations (there is an emulator parameter that can be set to either consider the randomness or ignore it), so the third entry on each line is superfluous.

## 8.6 Running *Smooth Emulator*

To tune the emulator, the User will run `MY_LOCAL/bin/smoothy_tune` which should have been compiled in the directions above. The User needs to edit one additional file at this point, the parameter file that sets numerous options for *Smooth Emulator*. For the template used in this tutorial, that file is

```

#SmoothEmulator_LogFileName smoothlog.txt
SmoothEmulator_LAMBDA 3.0
SmoothEmulator_MAXRANK 4
SmoothEmulator_NASample 8 # No. of coefficient samples
SmoothEmulator_TuneChooseMCMC true # set false if NPars<5
SmoothEmulator_UseSigmaYRreal false #
SmoothEmulator_ConstrainA0 false
SmoothEmulator_CutoffA false
SmoothEmulator_ModelRunDirName modelruns
SmoothEmulator_CoefficientsDirName coefficients
SmoothEmulator_ModelParInfoDir Info
SmoothEmulator_ObservableInfoDir Info
SmoothEmulator_TrainingPts 0-4
SmoothEmulator_MCStepSize 0.01

```

```

SmoothEmulator_MCSigmaAStepSize 50.0
SmoothEmulator_NMC 10000 # Steps between samples
SmoothEmulator_UsePCA false

```

The parameters are described in detail in Sec. 5. For the purpose of this tutorial we review a few of them. The parameter `SmoothEmulator_TrainingPts` is set above to include all five of the training points. However, if *Simplex Sampler* were run with a different choice, or if the User had added more points to be used by some other means, this parameter should be adjusted. The three parameters, `SmoothEmulator_MCStepSize`, `SmoothEmulator_MCSigmaAStepSize` and `SmoothEmulator_NMC` might be adjusted depending on the rate of Monte Carlo convergence of the tuning.

Now, running `smoothy_tune`, produces the following output,

```
MY_PROJECT% smoothy_tune parameters/emulator_parameters.txt
```

```
Tuning Emulator for length
```

```

success percentage=30.900000, SigmaA=170.951454, logP/Ndof=-0.604220,BestLogP/Ndof=-0.4826
success percentage=29.980000, SigmaA=842.098167, logP/Ndof=-0.731110,BestLogP/Ndof=-0.5440
success percentage=30.840000, SigmaA=184.488302, logP/Ndof=-0.651830,BestLogP/Ndof=-0.5288
success percentage=31.330000, SigmaA=151.392117, logP/Ndof=-0.784629,BestLogP/Ndof=-0.5839
success percentage=30.120000, SigmaA=120.719517, logP/Ndof=-0.588851,BestLogP/Ndof=-0.4810
success percentage=29.950000, SigmaA=178.510175, logP/Ndof=-0.671494,BestLogP/Ndof=-0.4950
success percentage=29.710000, SigmaA=281.210451, logP/Ndof=-0.750559,BestLogP/Ndof=-0.6051
success percentage=29.930000, SigmaA=171.487852, logP/Ndof=-0.599591,BestLogP/Ndof=-0.5212

```

```
Tuning Emulator for mass
```

```

success percentage=19.490000, SigmaA=53.649126, logP/Ndof=-0.625247,BestLogP/Ndof=-0.49247
success percentage=19.910000, SigmaA=84.204171, logP/Ndof=-0.656591,BestLogP/Ndof=-0.52713
success percentage=17.380000, SigmaA=46.138790, logP/Ndof=-0.780594,BestLogP/Ndof=-0.57353
success percentage=18.750000, SigmaA=92.444543, logP/Ndof=-0.531562,BestLogP/Ndof=-0.50588
success percentage=19.190000, SigmaA=60.575128, logP/Ndof=-0.540636,BestLogP/Ndof=-0.45743
success percentage=17.920000, SigmaA=75.219279, logP/Ndof=-0.674120,BestLogP/Ndof=-0.53911
success percentage=18.360000, SigmaA=40.116381, logP/Ndof=-0.656642,BestLogP/Ndof=-0.55767
success percentage=19.470000, SigmaA=74.261850, logP/Ndof=-0.675975,BestLogP/Ndof=-0.56443

```

```
Tuning Emulator for time
```

```

success percentage=19.760000, SigmaA=82.486706, logP/Ndof=-0.669089,BestLogP/Ndof=-0.48988
success percentage=19.970000, SigmaA=44.492860, logP/Ndof=-0.632434,BestLogP/Ndof=-0.49769
success percentage=19.520000, SigmaA=80.956164, logP/Ndof=-0.682640,BestLogP/Ndof=-0.54988
success percentage=19.340000, SigmaA=70.625835, logP/Ndof=-0.668105,BestLogP/Ndof=-0.51880
success percentage=19.590000, SigmaA=116.988915, logP/Ndof=-0.673993,BestLogP/Ndof=-0.5827
success percentage=19.080000, SigmaA=89.891122, logP/Ndof=-0.631130,BestLogP/Ndof=-0.50822
success percentage=19.880000, SigmaA=102.710384, logP/Ndof=-0.707850,BestLogP/Ndof=-0.5426
success percentage=19.450000, SigmaA=206.622539, logP/Ndof=-0.800493,BestLogP/Ndof=-0.4865

```

If the success percentage above is small, much less than 50%, then one may wish to decrease the Monte Carlo stepsizes, `SmoothEmulator_MCStepSize` and `SmoothEmulator_MCSigmaAStepSize` 0.01, to increase the success rate of the Metropolis steps. If the percentage is large, much greater

than 50%, conversely, one may wish to increase the step size to increase the coverage of the Monte Carlo trace. The quantity **SigmaA** represents the characteristic width of the distributions of coefficients. It is allowed to vary, and from experience, the range can be quite large. If the value of **SigmaA** does not seem to have randomized, one may wish to increase the number of Monte Carlo steps, **SmoothEmulator\_NMC**. The quantity **logP/Ndof** is the logarithm of the weight for a coefficient set divided by the number of degrees of freedom (number of coefficients). It should fluctuate throughout the trace. Higher values (usually less negative) indicate better fits. If the value is well below -1, and rising throughout the trace, it suggests the burn-in was insufficient and **SmoothEmulator\_NMC** should be increased. The last entry in the output lines, **BestLogP/Ndof**, simply reports the highest weight encountered in the trace. Again, if that is significantly increasing throughout the trace, **SmoothEmulator\_NMC** should be increased.

To represent the uncertainty, there were  $N_{\text{sample}}$  sets of coefficients generated by *Smooth Emulator*. In this case, the parameter **SmoothEmulator\_NASample** was set to 8. As the Monte Carlo trace proceeds, after **SmoothEmulator\_NMC** steps, the coefficients are stored, and an output line is generated as seen above. Thus there are 8 lines printed above for each observable.

The program generates Taylor coefficients for  $N_{\text{sample}}$  samples and saves them in the **coefficients** directory. Each observable has its own sub-directory with its name. Besides the coefficients, *Smooth Emulator* also produces a **meta.txt** file in the same directory. It includes essential information like the observable's maximum rank and net number of model parameters. This file is used to access the coefficients later. In this case, **smoother\_tune** created the directories, **MY\_PROJECT/coefficients/length**, **MY\_PROJECT/coefficients/mass** and **MY\_PROJECT/coefficients/time**. Within each of these sub-directories **smoother\_tune** created files **meta.txt**, **sample0.txt**, **sample1.txt** ... **sample7.txt**.

## 9 Generating Emulated Observables

Finally, now that the emulator is tuned, one may wish to generate emulated values for the observables for specified points in model-parameter space. A sample program, **MY\_LOCAL/bin/smoother\_calcobs** is provided to illustrate how this can be accomplished. If one invokes the executable, using the same parameters as those used by **smoother\_tune**, the User is prompted to enter the coordinates of a point in model-parameter space, after which **smoother\_calcobs** prints out the observables. In this case, for the case where **par1=20**, **par2=40**, **par3=60** and **par4=80**,

```
~/git/smooth_projects/scottrun% smoother_calcobs parameters/emulator_parameters.txt
Prior Info
#      ParameterName Type      Xmin_or_Xbar  Xmax_or_SigmaX
0:      par1      uniform          0           100
1:      par2      uniform          0           100
2:      par3      uniform          0           100
3:      par4      uniform          0           100
Enter value for par1:
20
Enter value for par2:
40
```

```

Enter value for par3:
60
Enter value for par4:
80
length = 61.349 +/- 5.14857
mass = -26.4496 +/- 2.6206
time = -44.9779 +/- 2.0495

```

One can test the emulator by entering the coordinates of training point. For example, one of the training points is `par1=50`, `par2=50`, `par3=88.7298` and `par4=40`. Running `smoothy_calcobs` for that coordinate,

```

~/git/smooth_projects/scottrun% smoothy_calcobs parameters/emulator_parameters.txt
Prior Info
#      ParameterName Type      Xmin_or_Xbar  Xmax_or_SigmaX
0:      par1      uniform          0           100
1:      par2      uniform          0           100
2:      par3      uniform          0           100
3:      par4      uniform          0           100
Enter value for par1:
50
Enter value for par2:
50
Enter value for par3:
88.7298
Enter value for par4:
40
length = -45.878 +/- 6.7435e-07
mass = -7.8083 +/- 8.42937e-08
time = -52.4031 +/- 1.16801e-06

```

Note that the uncertainties for the emulation are not effectively zero, as each set of the 8 sets of coefficients provides an an emulator that exactly reproduces the training points.

Of course, it is unlikely the User will wish to enter model parameters interactively as was done above. To incorporate `Smooth Emulator` into other programs, the User should inspect the main program, `MY_LOCAL/main_programs/smoothy_calcobs_main.cc`. The User can then design their own program based on this source code, and compile and link it by editing `MY_LOCAL/build/CMakeLists.txt`. By editing the CMake file, replacing the lines unique to `smoothy_calcobs`, one can easily compile new executables based on the User's main programs. To understand what might be involved, the source code in `MY_LOCAL/main_programs/smooth_calcobs_main.cc` is

```

#include "msu_commonutils/parametermap.h"
#include "msu_smooth/master.h"
#include "msu_commonutils/log.h"
using namespace std;

```

```

int main(int argc, char *argv[]){
    if(argc!=2){
        CLog::Info("Usage smoothy_calcobs emulator parameter filename");
        exit(1);
    }
    CparameterMap *parmap=new CparameterMap();
    parmap->ReadParsFromFile(string(argv[1]));
    CSmoothMaster master(parmap);
    master.ReadCoefficientsAlly();

    CModelParameters *modpars=new CModelParameters(master.priorinfo); // contains info about

    master.priorinfo->PrintInfo();
    // Prompt user for model parameter values
    vector<double> X(modpars->NModelPars);
    for(int ipar=0; ipar<modpars->NModelPars; ipar++){
        cout << "Enter value for " << master.priorinfo->GetName(ipar) << ":\n";
        cin >> X[ipar];
    }
    modpars->SetX(X);

    // Calc Observables
    CObservableInfo *obsinfo=master.observableinfo;
    vector<double> Y(obsinfo->NObservables);
    vector<double> SigmaY(obsinfo->NObservables);
    master.CalcAlly(modpars, Y, SigmaY);
    for(int iY=0; iY<obsinfo->NObservables; iY++){
        cout << obsinfo->GetName(iY) << " = " << Y[iY] << " +/- " << SigmaY[iY] << endl;
    }

    return 0;
}

```

This illustrates how one can write a code that

- a) Reads the parameter file
- b) Creates a *master* emulator file (called master because it includes emulators for all the observables)
- c) Creates a model-parameters object, **modpars**, that stores the coordinates of the model-parameter point
- d) Calculates the observables from the emulator

## 9 Underlying Theory of *Smooth Emulator*

The choice of model emulators,  $\mathbf{E}(\vec{\theta})$ , depends on the prior understanding of the model being emulated,  $\mathbf{M}(\vec{\theta})$ . If one knows that a function is linear, then a linear fit is clearly the best choice. Whereas to reproduce lumpy features, where the lumps have a characteristic length scale, Gaussian process emulators are highly effective. The quality of an emulator can be assessed through the following criteria:

- $\mathbf{E}(\vec{\theta}_t) = \mathbf{M}(\vec{\theta}_t)$  at the training points,  $\vec{\theta}_t$ .
- The emulator should reasonably reproduce the model away from the training points. This should hold true for either interpolation or extrapolation.
- The emulator should reasonably represent its uncertainty
- A minimal number of training points should be needed
- The method should easily adjust to larger numbers of parameters,  $\theta_i$ ,  $i = 1 \dots N$
- The emulator should not be affected by unitary transformations of the parameter space
- The emulator should be able to account for noisy models
- Training and running the emulator should not be numerically intensive

Here the goal is to focus on a particular class of functions: functions that are *smooth*. Smoothness is a prior knowledge of the function. It is an expectation that the linear terms of the function are likely to provide more variance than the quadratic contributions, which are in turn likely to be more important than the cubic corrections, and so on.

### 9.1 Mathematical Form of *Smooth Emulator*

To that end the following form for  $\mathbf{E}(\vec{\theta})$  is chosen,

$$\mathbf{E}(\vec{\theta}) = \sum_{\vec{n}, \text{s.t. } K(\vec{n}) \leq K_{\max}} d_{\vec{n}} f_{K(\vec{n})}(|\vec{\theta}|) \mathbf{A}_{\vec{n}} \left(\frac{\theta_1}{\Lambda}\right)^{n_1} \left(\frac{\theta_2}{\Lambda}\right)^{n_2} \dots \left(\frac{\theta_N}{\Lambda}\right)^{n_N}. \quad (9.1)$$

Each term has a rank  $K(\vec{n}) = n_1 + n_2 + \dots + n_N$ . If  $f$  is constant, the rank of that term corresponds to the power of  $|\vec{\theta}|/\Lambda$ . All terms are included up to a given rank,  $K_{\max}$ . The coefficients  $\mathbf{A}$  are stochastically distributed. The coefficients  $d_{\vec{n}}$  will ensure that the variance is independent of the direction of  $\vec{\theta}$ , with the constraint that  $d_{K,0,0,\dots} = 1$ . The function  $f_K(|\vec{\theta}|)$  provides the freedom to alter how the behavior depends on the distance from the origin,  $|\vec{\theta}|$ , and on the rank,  $K$ . Given that the variance of  $\mathbf{A}_{\vec{n}}$  can be changed,  $f_{K=0}(|\vec{\theta}| = 0)$  is also set to unity for all  $K$  without loss of generality. For each combination  $\vec{n}$ , the prior probability for any the  $\mathbf{A}$  coefficients is given by

$$p(\mathbf{A}_{\vec{n}}) = \frac{1}{\sqrt{2\pi\sigma_{K(\vec{n})}^2}} e^{-\mathbf{A}_{\vec{n}}^2 / 2\sigma_{K(\vec{n})}^2}, \quad (9.2)$$

$$\langle \mathbf{A}_{\vec{n}}^2 \rangle = \sigma_{K(\vec{n})}^2.$$

The variance,  $\sigma_K^2$ , is allowed to vary as a function of  $K$ .

The parameter  $\Lambda$  will be referred to as the *smoothness parameter*. Here, we assume that all parameters have a similar range, of order unity, e.g.  $-1 < \theta_i < 1$ . Thus, the relative importance of each term Eq. (9.1) falls with increasing rank,  $K$ , as  $(1/\Lambda)^K$ . For now, the smoothness parameter is fixed by prior knowledge, i.e. one chooses higher values of  $\Lambda$  if one believes the function to be close to linear.

First, we consider the variance of the emulator at a given point,  $\vec{\theta}$ . Requiring that the variance is independent of the direction of  $\vec{\theta}$  will fix  $d_{\vec{n}}$ . For example, transforming  $\theta_1$  and  $\theta_2$  to parameters  $(\theta_1 \pm \theta_2)/\sqrt{2}$  should not affect the accuracy or uncertainty of the emulator.

At  $|\vec{\theta}| = 0$  the only term in Eq. (9.1) that contributes to the variance is the one  $K = 0$  term. Averaging over the  $\mathbf{A}$  coefficients, which can be either positive or negative with equal probability,

$$\langle E(\vec{\theta}) \rangle = 0, \quad (9.3)$$

where the averaging refers to an average over the  $\mathbf{A}$  coefficients. At the origin,  $|\vec{\theta}| = 0$ , the variance of  $E$  is

$$\langle E(\theta_1 = \theta_2 = \dots \theta_N = 0)^2 \rangle = d_{n_i=0}^2 \sigma_{K=0}^2 f_{K=0}^2(\vec{\theta} = 0). \quad (9.4)$$

Choosing  $f_{K=0}(0) = 1$  and  $d_{n_i=0} = 1$ , the variance of  $E$  is indeed  $\sigma_0^2$ . The variance at some point  $\vec{\theta} \neq 0$  is

$$\langle E^2(\vec{\theta}) \rangle = \sum_{\vec{n}} f_K^2(|\vec{\theta}|) \sigma_{K(\vec{n})}^2 d_{\vec{n}}^2 \left( \frac{\theta_1^{2n_1}}{\Lambda^2} \right) \left( \frac{\theta_2^{2n_2}}{\Lambda^2} \right) \dots \left( \frac{\theta_N^{2n_N}}{\Lambda^2} \right). \quad (9.5)$$

If  $\langle E^2 \rangle$  is to be independent of the direction of  $\vec{\theta}$ , the sum above must be a function of  $|\vec{\theta}|^2$  only. This requires the net contribution from each rank,  $K$  to be proportional to  $|\vec{\theta}|^{2K}$  multiplied by some function of  $K$ . Using the fact that

$$(\vec{\theta}_a \cdot \vec{\theta}_b)^K = \sum_{\vec{n}, s.t. \sum_i n_i = K} \frac{K!}{n_1! \dots n_N!} (\theta_{a1} \theta_{b1})^{n_1} \dots (\theta_{aN} \theta_{bN})^{n_N}, \quad (9.6)$$

one can see that if the sum is to depend only on the norm of  $\vec{\theta}$ ,

$$d_{\vec{n}}^2 = \frac{K(\vec{n})!}{n_1! n_2! \dots n_N!}. \quad (9.7)$$

The factor of  $K!$  in the numerator was chosen to satisfy the condition that  $d_{K,0,0,0} = 1$ .

One can now calculate the correlation between the emulator at two different points, averaged over all possible values of  $\mathbf{A}$ ,

$$\langle E(\vec{\theta}_a) E(\vec{\theta}_b) \rangle = \sum_{K=0}^{K_{\max}} \sigma_K^2 f_K^2(|\vec{\theta}|) \left( \frac{\vec{\theta}_a \cdot \vec{\theta}_b}{\Lambda^2} \right)^K. \quad (9.8)$$

Requiring  $\mathbf{f}(|\boldsymbol{\theta}| = 0) = \mathbf{1}$  gives

$$\langle \mathbf{E}^2(\vec{\boldsymbol{\theta}} = 0) \rangle = \sigma_0^2, \quad (9.9)$$

and for  $\vec{\boldsymbol{\theta}}_a = \vec{\boldsymbol{\theta}}_b$ ,

$$\langle \mathbf{E}^2(\vec{\boldsymbol{\theta}} = 0) \rangle = \sum_{K=0}^{K_{\max}} \sigma_K^2 \mathbf{f}_K^2(|\vec{\boldsymbol{\theta}}|) \left( \frac{|\vec{\boldsymbol{\theta}}|^2}{\Lambda^2} \right)^K. \quad (9.10)$$

To this point, the form is completely general once one requires that the variance above is independent of the direction of  $\vec{\boldsymbol{\theta}}$ . I.e. the function  $\mathbf{f}_K(\vec{\boldsymbol{\theta}})$  could be any function satisfying the constraint,  $\mathbf{f}_K(0) = \mathbf{1}$ , and  $\sigma_K^2$  could have any function of  $\mathbf{K}$ . Below, we illustrate how different choices for  $\mathbf{f}$  or for  $\sigma_K$  affect the emulator by comparing several variations. First, we define the default form.

## 9.2 Alternate Forms

As stated above, once the form is to provide variances that are independent of the direction of  $\vec{\boldsymbol{\theta}}$ , the general form going forward is

$$\mathbf{E}(\vec{\boldsymbol{\theta}}) = \sum_{\vec{n}, \text{s.t. } K(\vec{n}) < K_{\max}} \mathbf{f}_K(|\vec{\boldsymbol{\theta}}|) \left( \frac{K(\vec{n})!}{n_1! \dots n_N!} \right)^{1/2} A_{\vec{n}} \left( \frac{\theta_1}{\Lambda} \right)^{n_1} \left( \frac{\theta_2}{\Lambda} \right)^{n_2} \dots \left( \frac{\theta_N}{\Lambda} \right)^{n_N}, \quad (9.11)$$

$$P(\vec{A}_n) = \frac{1}{(2\pi\sigma_K^2)^{1/2}} e^{-|\vec{A}_n|^2/2\sigma_K^2}.$$

Variations from the general form involve adjusting either the  $\mathbf{K}$ -dependence of the  $|\vec{\boldsymbol{\theta}}|$ -dependence of  $\mathbf{f}_K(|\vec{\boldsymbol{\theta}}|)$ , or adjusting the  $\mathbf{K}$ -dependence of  $\sigma_K$ .

### Default Form

Here, we assume  $\mathbf{f}_K(|\vec{\boldsymbol{\theta}}|)$  is independent of  $|\vec{\boldsymbol{\theta}}|$ , and that  $\sigma_K$  is independent of  $\mathbf{K}$ . Further, the  $\mathbf{K}$ -dependence of  $\mathbf{f}^2$  is assumed to be  $\mathbf{1}/\mathbf{K}!$ . With this choice

$$\begin{aligned} \mathbf{E}(\vec{\boldsymbol{\theta}}) &= \sum_{\vec{n}, K(\vec{n}) \leq K_{\max}} \frac{1}{\Lambda^K} \frac{A_{\vec{n}}}{\sqrt{n_1! n_2! \dots n_N!}} \theta_1^{n_1} \dots \theta_N^{n_N}, \\ P(A_{\vec{n}}) &\sim e^{-A_{\vec{n}}^2/2\sigma^2}. \end{aligned} \quad (9.12)$$

With this form the variance increases with  $|\vec{\boldsymbol{\theta}}|$ ,

$$\langle \mathbf{E}^2(\vec{\boldsymbol{\theta}}) \rangle = \sigma^2 e^{|\vec{\boldsymbol{\theta}}|^2/\sigma^2}. \quad (9.13)$$

If the function is trained in a region where the function is linear, the emulator's extrapolation outside the region will continue to follow the linear behavior, albeit with variation from the higher order coefficients.

The choice of  $\mathbf{f}_K^2 = \mathbf{1}/\mathbf{K}!$  ensures that the sum defining  $\mathbf{E}(\vec{\boldsymbol{\theta}})$  converges as a function of  $\mathbf{K}$  as long as  $K_{\max}$  is rather large compared to  $|\vec{\boldsymbol{\theta}}|/\Lambda$ .



### Variant A: Letting $\sigma_K$ have a $K$ dependence

One reasonable alteration to the default choice might be to allow the  $K = 0$  term to take any value, i.e.  $\sigma_{K=0} = \infty$ , while setting all the other  $\sigma_K$  terms equal to one another. This would make sense if our prior expectation of smoothness meant that we expect the  $K = 2$  terms to be less important than the  $K = 1$  terms, by some factor  $|\vec{\theta}|/\Lambda$ , but that the variation of the  $K = 1$  term is unrelated to the size of the  $K = 0$  term. This would make the emulator independent of redefinition of the emulated function by adding a constant. This may well be a reasonable choice for many circumstances.

### Variant B: Suppressing correlations for large $\Delta\vec{\theta}$

This form for  $f$  causes correlations to fall for points far removed from one another.

$$f_K(|\vec{\theta}|) = \frac{1}{\sqrt{K!}} \left\{ \sum_{K=0}^{K_{\max}} \frac{1}{\sqrt{K!}} \left( \frac{|\vec{\theta}|^2}{2\Lambda^2} \right)^K \right\}^{-1/2}.$$

In the limit that  $K_{\max} \rightarrow \infty$  the form is a simple exponential,

$$f_K(|\vec{\theta}|) \Big|_{K_{\max} \rightarrow \infty} = \frac{1}{\sqrt{K!}} e^{-|\vec{\theta}|^2/2\Lambda^2}. \quad (9.14)$$

With this form the same-point correlations remain constant over all  $\vec{\theta}$ ,

$$\langle E(\vec{\theta}) E(\vec{\theta}) \rangle = \sigma^2, \quad (9.15)$$

while the correlation between separate positions fall with increasing separation. This is especially transparent for the  $K_{\max} \rightarrow \infty$  limit,

$$\langle E(\vec{\theta}_a) E(\vec{\theta}_b) \rangle_{K_{\max} \rightarrow \infty} = \sigma^2 e^{-|\vec{\theta}_a - \vec{\theta}_b|^2/2\Lambda^2}.$$

In this limit one can also see that

$$\langle [E(\vec{\theta}) - E(\vec{\theta}')]^2 \rangle_{K_{\max} \rightarrow \infty} = 2\sigma^2 \left( 1 - e^{-|\vec{\theta} - \vec{\theta}'|^2/2\Lambda^2} \right). \quad (9.16)$$

If one trains such an emulator in one region, then extrapolates to a region separated by  $|\vec{\theta} - \vec{\theta}'| \gg \Lambda$ , the average predictions will return to zero. Thus, if the behavior would appear linear in some region the emulator's distribution of predictions far away (extrapolations) would center at zero. This behavior is similar to a Gaussian-process emulator.

### Variant C: Eliminating the $1/K!$ weight

Clearly, eliminating the  $1/K!$  weights in  $f_K$  would more emphasize the contributions from larger  $K$ . But, for and  $|\vec{\theta}| > \Lambda$  the contribution to the variance would increase as  $K$  increases and the sum would not converge if  $K_{\max}$  were allowed to approach infinity. An example of a function that expands without factorial suppression is  $1/(1-x) = 1 + x + x^2 + x^3 \dots$ , which diverges as  $x \rightarrow 1$ . If such behavior is not expected, then this choice would be unreasonable.

## 9.3 Tuning the Emulator

Here, we illustrate how an emulator of the form above can be constrained given training points. We consider  $M$  full model runs at positions  $\vec{\theta}_{m=1,M}$ , with values  $F_{m=1,M}$ . The functional form

has a large number of coefficients,  $\mathbf{A}_{\vec{n}}$ , which is much larger than the number of training points,  $M$ . However, the dependence of  $\mathbf{A}$  is purely linear. One can enumerate the coefficients as  $\mathbf{A}_{\mathbf{c}}$  with  $\mathbf{c} = 1 \cdots C$ , where  $C > M$ . One can randomly set the coefficients  $\mathbf{A}_M$  through  $\mathbf{A}_C$ , then solve for the first  $M$  coefficients by solving a linear equation. One can then apply a weight based on the values of  $\mathbf{A}$  consistent with the prior likelihood of  $\mathbf{A}$  and the constraints. One can then generate a representative set of  $\mathbf{A}$ , perhaps a dozen samples. Each sample function will go through all the training points, but will vary further from the training points. Averaging over the  $N_{\text{sample}}$  sets of coefficients can be used to make a prediction for the emulator at some point  $\vec{\theta}$ , and the variance of those  $N_{\text{sample}}$  points would represent the uncertainty of the emulator.

Here, we present two different methods for generating samples of points that are consistent with the training constraints and with the prior range of parameters. We do this for the default method, but this can be easily extended to the other model variants listed in the previous section. In addition to varying the coefficients, we will additionally vary the width parameters,  $\sigma$ .

First, we need to describe how the weights are generated. The probability of a set of coefficients is initially

$$P(\mathbf{A}_c) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\mathbf{A}_c^2/2\sigma^2}. \quad (9.17)$$

If we also vary  $\sigma$ , we can state that

$$Q(\sigma) = \frac{1}{\pi} \frac{\Gamma/2}{\sigma^2 + (\Gamma/2)^2}. \quad (9.18)$$

Here, the half width of the distribution should be set to some large value to encompass the degree to which the emulated function might vary. The Lorentzian form accommodates a large uncertainty. The result should not depend strongly on  $\Gamma$ . The joint probability is  $Q(\sigma) \prod_c P(\mathbf{A}_c)$ .

The constrained probability is

$$\begin{aligned} dP &= d\sigma Q(\sigma) \prod_c [d\mathbf{A}_c P(\mathbf{A}_c)] \prod_{m=1}^M \delta(E(\mathbf{A}, \sigma, \vec{\theta}_m) - F(\vec{\theta}_m)) \\ &= d\sigma Q(\sigma) \prod_{c=1}^M P(\mathbf{A}_c) \frac{1}{|J|} \prod_{c=M}^C [d\mathbf{A}_c P(\mathbf{A}_c)]. \end{aligned} \quad (9.19)$$

Here,  $|J|$  is the Jacobian, i.e. it is the determinant of the  $M \times M$  matrix

$$J_{mc} = \frac{\partial E(\vec{\theta}_m)}{\partial \mathbf{A}_c}, \quad (9.20)$$

For the default form in the previous section,

$$J_{mc} = \frac{1}{\sqrt{n_{c1}! \cdots n_{cN}!}} (\theta_{m1})^{n_{c1}} \cdots (\theta_{mN})^{n_{cN}}. \quad (9.21)$$

Because the Jacobian depends only on the position of the training points, it can be treated as a constant.

One can now sample  $\mathbf{A}$  and  $\sigma$  according to the weights above. Here, we list two methods.

a) **Keep and Reject Method**

One can generate the coefficients  $\mathbf{A}_{c=M} \cdots \mathbf{A}_C$  according to the Gaussian distribution with width  $\sigma$ , after generating  $\sigma$  according to the Lorentzian. The residual weight is then

$$w = \prod_{c=1}^M P(\mathbf{A}_c). \quad (9.22)$$

For each attempt, one could keep or reject the attempt with probability  $w$ . This generates perfectly independent samples, but with the caveat that in a high dimensional space the rejection level might be quite high.

b) **Metropolis Exploration of  $\mathbf{A}$  and  $\sigma$**

Beginning with any configuration  $\mathbf{A}$  and width  $\sigma$ , one can take a random step  $\delta\mathbf{A}$  and  $\delta\sigma$ . In the absence of any weight this would consider all  $\mathbf{A}$  or  $\sigma$  with values from  $-\infty$  to  $\infty$ . The residual weight would then be

$$w = Q(\sigma) \prod_{c=1}^C P(\mathbf{A}_c). \quad (9.23)$$

One then keeps the step if the new weight exceeds the previous one, or if the ratio of the new weight to the old weight is greater than a random number between zero and unity. After some number of steps,  $N_{\text{steps}}$ , one saves the configuration as a representative sample for the emulator. The disadvantage of this approach is that many steps may be needed to ensure that the samples are independent, and thus faithfully represent the variation in the emulator.

## 10 Theoretical Basis of *Simplex Sampler*

Here, we imagine a spherically symmetric prior, e.g. one that is purely Gaussian, and where the parameters are scaled so that the prior distribution is invariant to rotations. If one believe the function were close to linear, a strategy would be to find  $M = N + 1$  points in parameter space placed far apart from another. One choice is the  $N$ -dimensional simplex. Examples are an equilateral triangle in two-dimensions or a tetrahedron in three dimensions. For a simplex, one places the  $M = N + 1$  training points at a uniform distance from the origin,  $\mathbf{r}$ , with equal separation between each point. One can generate an  $N$ -dimensional simplex from an  $N - 1$  dimensional arrangement. In the  $N - 1$  dimensional arrangement, the points are arranged equidistant from one another using the coordinates  $\mathbf{x}_1 \cdots \mathbf{x}_N$ . The points would all be placed at a radius  $r_N$  from the origin in this system, and the separation would be  $d$ . In the  $N$ -dimensional system all these  $N - 1$  points had coordinate  $\mathbf{x}_N = -\mathbf{a}$ . The  $(N + 1)^{\text{th}}$  point is then placed at position  $\mathbf{x}_1 \cdots \mathbf{x}_N = \mathbf{0}$  and  $\mathbf{x}_{N+1} = N\mathbf{a}$ . This keeps the center of mass at zero. One then chooses  $\mathbf{a}$  such that the new point is equally separated from all the other points by the same distance  $d$ ,

$$\begin{aligned} d^2 &= r_{N-1}^2 + N^2 a^2, \\ a &= \sqrt{\frac{d^2}{N^2} - r_{N-1}^2}. \end{aligned} \quad (10.1)$$

Now, each of the  $N$  points is located a distance  $N\mathbf{a}$  away from the center of the  $N$ –dimensional origin. This procedure can be applied iteratively to generate the vertices of the simplex.

One might also wish to use enough training points to uniquely determine the emulator in the case that the function is quadratic. There are  $N(N+1)/2$  additional points, which is exactly the number of segments connecting the  $N+1$  equally-spaced vertices of the  $N$ –dimensional simplex. If placed at the midpoints of the segments, these points would be closer to the origin than the vertices. One of the simplex options is to place these points at the midpoints, then double their radii while maintaining their direction. This would result in arrays of points at two different radii, with  $N+1$  points positioned at the lower radius and  $N(N+1)/2$  points being placed at the larger radius.

Choosing the training points depends on prior expectation of the emulated function. The simplex choice for the first  $M = N+1$  points seems logical. Even if another method, such as a Gaussian process emulator is to be implemented, such methods are often based on first understanding the linear behavior. The simplex strategy would seem a good way to pick the first  $N+1$  training points.

One issue with the simplex is that the first set of  $N+1$  training points would all be placed at the same radius. If the prior parameter distribution is uniform within an  $N$ –dimensional hyper cube, the training points could be rather far from the corners in that space. Issues with such priors are discussed in the next section.

## 10.1 The Pernicious Nature of Step-Function Priors in High Dimension

For purely Gaussian priors, one can scale the prior parameter space to be spherically symmetric. Unfortunately, that is not true for step function priors (uniform within some range). In that case the best one can do (if the priors for each parameter are independent) is to scale the parameter space such that each parameter has the constraint,  $-1 < \theta_i < 1$ . If the number of parameters is  $N$ , the hyper-cube has  $2N$  faces and  $2^N$  corners, a face being defined as one parameter being  $\pm 1$  while the others are zero, while a corner has each parameter either  $\pm 1$ . For 10 parameters, there are 1024 corners, and for 15 parameters there are 32678 corners. Thus, it might be untenable to place a training point in each corner.

One can also see the problem with placing the training points in a spherically symmetric fashion as is done with the *Simplex Sampler*. The hyper-volume of the parameter-space hyper-cube is  $2^N$ , whereas the volume of an  $N$ –dimensional hyper-sphere of radius  $R = 1$  is

$$V_{\text{sphere}} = \Omega_N \int_0^R dr r^{N-1} = \Omega_N \frac{R^N}{N}. \quad (10.2)$$

The solid angle,  $\Omega_N$  in  $N$  dimensions is

$$\Omega_N = \frac{2\pi^{N/2}}{\Gamma(N/2)}, \quad (10.3)$$

and after putting this together, the fraction of the hyper-cube’s volume that is within the hyper-sphere is

$$\frac{V_{\text{sphere}}}{V_{\text{cube}}} = \begin{cases} \frac{(\pi/2)^{N/2}}{(\pi/2)^{(N-1)/2}}, & N = 2, 4, 6, 8, \dots \\ \frac{(\pi/2)^{N/2}}{N!!}, & N = 3, 5, 7, \dots \end{cases} \quad (10.4)$$

In two dimensions, the ratio is  $\pi/4$ , and in three dimensions it is  $\pi/6$ . In 10 dimensions it is  $2.5 \times 10^{-3}$ . For high dimensions only a small fraction of the parameter space can ever lie inside inside a sphere used to place points. And, if the model is expensive, it may not be tenable to run the full model inside every corner.

One can also appreciate the scope of the problem by considering the radius of the corners vs. the radius of the sphere. The maximum value of  $|\vec{\theta}|$  is  $\sqrt{N}$ . So, for 9 parameters, if the training points were all located at positions  $|\vec{\theta}| < 1$ , one would have to extrapolate all the way to  $|\vec{\theta}| = 3$ . Thus, unless the model is exceptionally smooth, one needs to devise a strategy to isolate the portion of likely parameter space using some original set of full-model runs, then augment those runs in the likely region.

A third handle for viewing the issue in  $N$  dimensions is to compare the r.m.s. radii of the hypersphere to that of the hyper-cube. For the cube where each side has length  $2a$ ,

$$(R_{\text{r.m.s.}}^{(\text{cube})})^2 = a^2 \frac{N}{3}. \quad (10.5)$$

whereas for a sphere of radius  $a$ ,

$$(R_{\text{r.m.s.}}^{(\text{sphere})})^2 = a^2 \frac{N+2}{N}. \quad (10.6)$$

The ratio of the radii is then

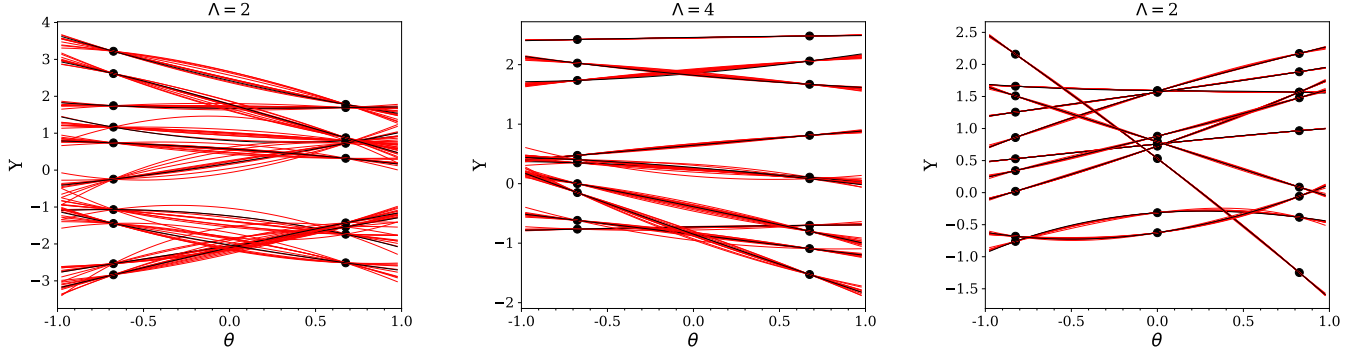
$$\frac{R_{\text{r.m.s.}}^{(\text{sphere})}}{R_{\text{r.m.s.}}^{(\text{cube})}} = \sqrt{\frac{3}{N+2}}. \quad (10.7)$$

Thus, in 10 dimensions, if the training points are placed at a distance  $a$  from the origin, the r.m.s. radii of the interior space would be half that of the entire space. Further, the r.m.s. radii of the points in the cube,  $a\sqrt{N/3}$ , would be about 83% higher than the radius of the training points.

Of course, these problems would be largely avoided if the number of parameters was a half dozen or fewer, or if one was confident that the function was extremely smooth. In the first two sections of this paper, the smoothness parameter,  $\Lambda$ , was set to a constant. There might be prior knowledge that certain parameters affect the observables only weakly. In that case, the response to these parameters can be considered as linear. This could be done by scaling those parameters so that they vary over a smaller range. If a parameter varies only between  $-0.1$  and  $0.1$ , that effectively applies a smoothness parameter in those directions that is ten times higher. Unfortunately, the choice of which parameters to rescale in this fashion would likely vary depending on which observable is being emulated. Because all the observables might be calculated in a full model run, one needs to identify parameters that would likely have weak response on all observables.

## 11 Tests of the *Smooth Emulator using Simplex Sampler* for Training

First, we show some results of one-parameter emulators. For a linear fit, two parameters (slope and intercept) would suffice. But because higher-rank contributions exist, there is a spread of functions

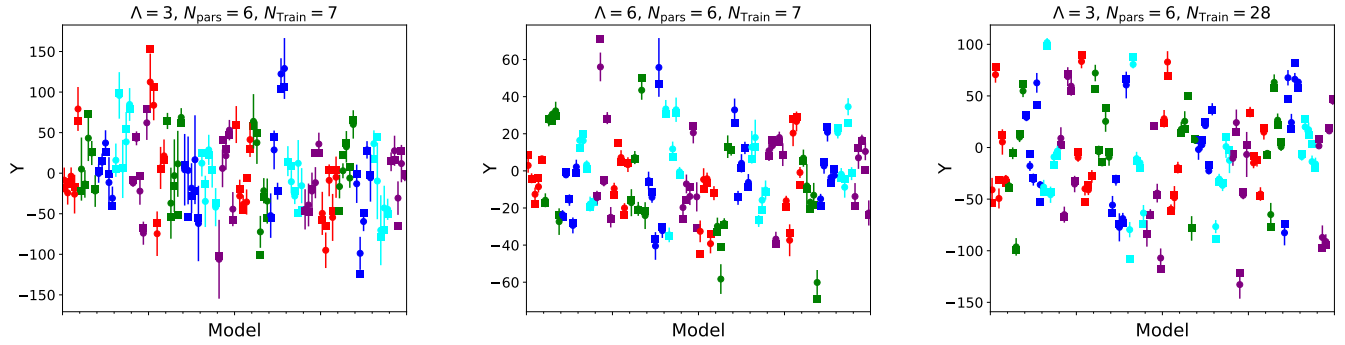


**Figure 11.1:** Real functions (black lines) are emulated with the *Smooth Emulator* (default settings). Ten different functions are emulated. Each emulation is sampled ten times (red lines), with the spread of the lines representing the uncertainty.

is that narrows with increasing smoothness parameter as shown in Fig. 11.1. Figure 11.1 also shows how increasing the number of training points, from two to three, also narrows the range of possible values.

Next, we repeat the same test with six parameters. The prior distribution of parameters was uniform in the region  $-1 < \theta_i < 1$ . In this case, the *Simplex Sampler* was used to choose the training points. The “real” model,  $\mathbf{F}$ , was constructed from sample smooth functions, with the coefficients generated randomly and a smoothness parameter set to three or six. The  $\mathbf{A}_{\vec{n}=0}$  coefficient for the real model was set to zero to better accommodate viewing results. The emulator was not given that information. Twenty instances of real models were emulated. For each real model five random point in parameter space were chosen. The emulator value and its uncertainty are plotted alongside the real-model values in Fig. 11.2. The 100 comparisons show that the emulator accurately predicts the uncertainty. This is not surprising, because the emulator used the same smoothness parameter as was used to construct the real models. The width parameters,  $\sigma$ , were explored stochastically, and remarkably they seem to fluctuate around the same value used to construct the real model, albeit with fluctuations of the order of 30%.

Figure 11.2 first shows the 100 comparisons for the case with the smoothness parameter,  $\Lambda = 3$ . The simplest simplex form was applied, which has seven parameters, the same number one would use for a linear fit. Uncertainties were provided by finding 16 independent samplings of  $\mathbf{A}$  coefficients and of  $\sigma$ , then using the variance of the 16 samples to define the uncertainty. Variant *a* of the default emulator was used, i.e. the coefficient  $\mathbf{A}_0$  was not given a constraining prior distribution, whereas all other coefficients were weight by a Gaussian of width  $\sigma$ . In the other two panels of Fig. 11.2, the procedure was repeated but once with a higher smoothness parameter, and once with a the 28 training points, placed according to the procedure mentioned in Sec. 3, where an additional set of points in parameter space was generated by choosing points that bisect all the lines connecting points in the original simplex, then doubling their radius. As expected, smoother functions are more easily emulated with a given number of training points, and using more training points also improves the accuracy.



**Figure 11.2:** Using 20 instances of real models using six parameters, choosing 5 random points in parameter space for each model, the emulator (circles) and its uncertainty were compared to the real values (squares). Neighboring points of the same color emulated the same real model. The accuracy improves if a smoother function is considered (middle panel), or if more training points are used (right panel). Estimates of the uncertainty seem reasonable given that the uncertainties illustrated in the figure represent one standard deviation. It should be emphasized that this consistency depends critically on the fact that the emulator chose the same smoothness parameter as was used to generate the real models.