

CHAPTER 2

Getting started

2.1 Installation

PRISM can be easily installed by either cloning the [repository](#) and installing it manually:

```
$ git clone https://github.com/1313e/PRISM
$ cd PRISM
$ pip install .
```

or by installing it directly from [PyPI](#) with:

```
$ pip install prism
```

PRISM can now be imported as a package with `import prism`. For using *PRISM* in MPI, `mpi4py >= 3.0.0` is required (not installed automatically).

The *PRISM* package comes with two `ModelLink` subclasses. These `ModelLink` subclasses can be used to experiment with *PRISM* to see how it works. [Using PRISM](#) has several examples explaining the different functionalities of the package.

2.2 Running tests

If one wants to run `pytest` on *PRISM*, all [requirements_dev](#) are required. The easiest way to run the tests is by cloning the [repository](#), installing all requirements and then running `pytest` on it:

```
$ git clone https://github.com/1313e/PRISM
$ cd PRISM
$ pip install -r requirements_dev.txt
$ pytest
```

If *PRISM* and all [requirements_dev](#) are already installed, one can run the tests by running `pytest` in the installation directory:

```
$ cd <path_to_installation_directory>/prism
$ pytest
```

When using Anaconda, the installation directory path is probably of the form `<HOME>/anaconda3/envs/<environment_name>/lib/pythonX.X/site-packages`.

2.3 Example usage

See *Minimal example* for a documented explanation on this example.

```
# Imports
from prism import Pipeline
from prism.modellink import GaussianLink

# Define model data and create ModelLink object
model_data = {3: [3.0, 0.1], 5: [5.0, 0.1], 7: [3.0, 0.1]}
modellink_obj = GaussianLink(model_data=model_data)

# Create Pipeline object
pipe = Pipeline(modellink_obj)

# Construct first iteration of the emulator
pipe.construct()

# Create projections
pipe.project()
```

The PRISM pipeline

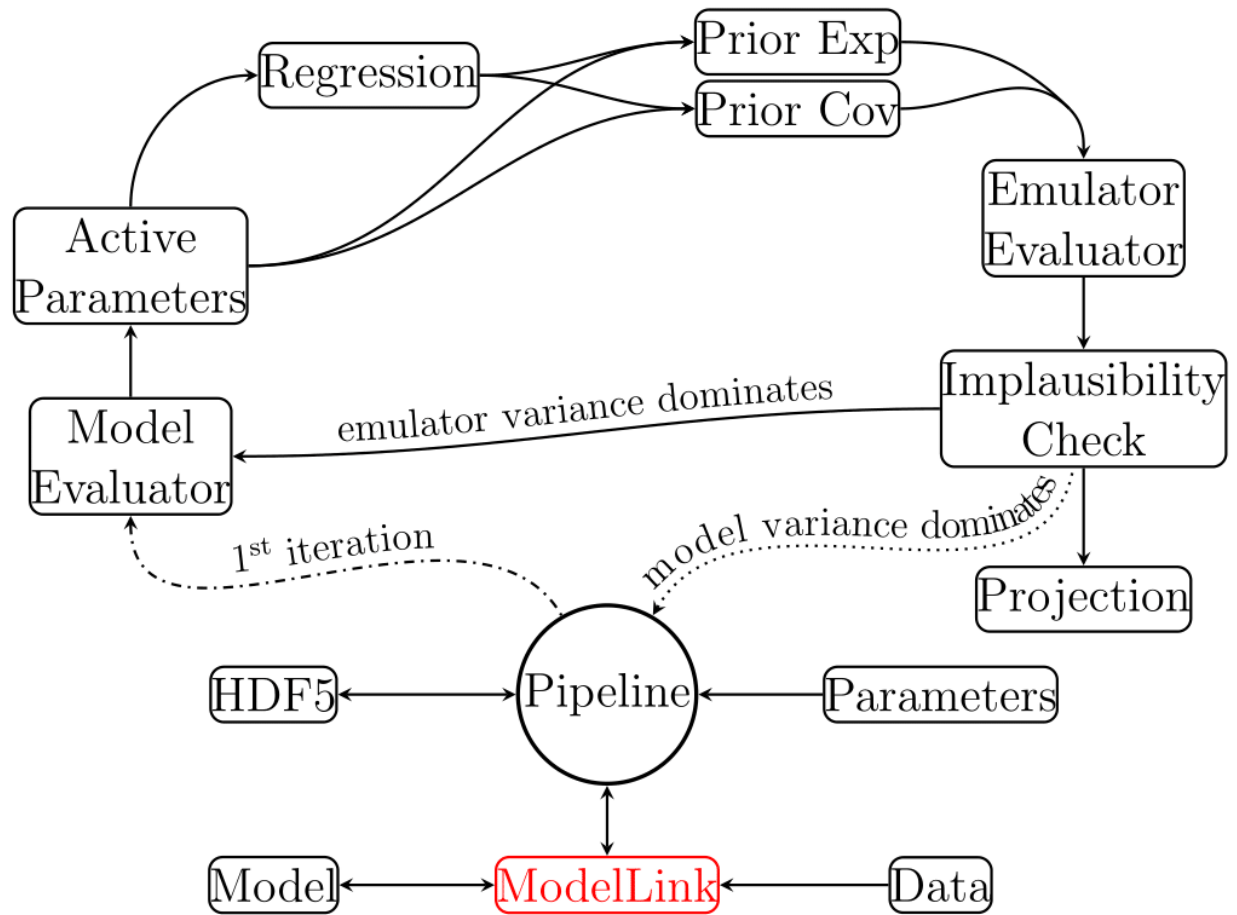
The overall structure of *PRISM* can be seen in Fig. 3.1 and will be discussed below. The *Pipeline* object plays a key-role in the *PRISM* framework as it governs all other objects and orchestrates their communications and method calls. It also performs the process of history matching and refocusing (see the [PRISM paper](#) for the methodology used in *PRISM*). It is linked to the model by a user-written *ModelLink* object (see [ModelLink: A crash course](#)), allowing the *Pipeline* object to extract all necessary model information and call the model. In order to ensure flexibility and clarity, the *PRISM* framework writes all of its data to one or several [HDF5-files](#) using [h5py](#), as well as [numpy](#).

The analysis of a provided model and the construction of the emulator systems for every output value, starts and ends with the *Pipeline* object. When a new emulator is requested, the *Pipeline* object creates a large Latin-Hypercube design (LHD) of model evaluation samples to get the construction of the first iteration of the emulator systems started. To ensure that the maximum amount of information can be obtained from evaluating these samples, a custom Latin-Hypercube sampling code was written. This produces LHDs that attempt to satisfy both the *maximin* criterion as well as the *correlation* criterion. This code is customizable through *PRISM* and publicly available in the [e13Tools](#) Python package.

This Latin-Hypercube design is then given to the *Model Evaluator*, which through the provided *ModelLink* object evaluates every sample. Using the resulting model outputs, the *Active Parameters* for every emulator system (individual data point) can now be determined. Next, depending on the user, polynomial functions will be constructed by performing an extensive *Regression* process for every emulator system, or this can be skipped in favor of a sole Gaussian analysis (faster, but less accurate). No matter the choice, the emulator systems now have all the required information to be constructed, which is done by calculating the *Prior Expectation* and *Prior Covariance* values for all evaluated model samples ($E(D_i)$ and $\text{Var}(D_i)$).

Afterward, the emulator systems are fully constructed and are ready to be evaluated and analyzed. Depending on whether the user wants to prepare for the next emulator iteration or create a projection (see [Projections](#)), the *Emulator Evaluator* creates one or several LHDs of emulator evaluation samples, and evaluates them in all emulator systems, after which an *Implausibility Check* is carried out. The samples that survive the check can then either be used to construct the new iteration of emulator systems by sending them to the *Model Evaluator*, or they can be analyzed further by performing a *Projection*. The *Pipeline* object performs a single cycle by default (to allow for user-defined analysis algorithms), but can be easily set to continuously cycle.

In addition to the above, *PRISM* also features a high-level *Message Passing Interface* (MPI) implementation using the Python package [mpi4py](#). All emulator systems in *PRISM* can be constructed independently from each other, in any order, and only require to communicate when performing the implausibility cut-off checks during history matching.

Fig. 3.1: The structure of the *PRISM* pipeline.

Additionally, since different models and/or architectures require different amounts of computational resources, *PRISM* can run on any number of MPI processes (including a single one in serial to accommodate for OpenMP codes) and the same emulator can be used on a different number of MPI processes than it was constructed on (e.g., constructing an emulator using 8 MPI processes and reloading it with 6). More details on the MPI implementation and its scaling can be found in *MPI implementation*.

In *Using PRISM* and *ModelLink: A crash course*, the various components of *PRISM* are described more extensively.

3.1 MPI implementation

Given that most scientific models are either already parallelized or could benefit from parallelization, we had to make sure that *PRISM* allows for both MPI and OpenMP coded models to be connected. Additionally, since individual emulator systems in an emulator iteration are independent of each other, the extra CPUs required for the model should also be usable by the emulator. For that reason, *PRISM* features a high-level MPI implementation for using MPI-coded models, while the Python package *NumPy* handles the OpenMP side. A mixture of both is also possible (using the *worker_mode* context manager).

Here, we discuss the MPI scaling tests that were performed on *PRISM*. For the tests, the same *GaussianLink* class was used as in *Minimal example*, but this time with 32 emulator systems (comparison data points) instead of 3. In *PRISM*, all emulator systems are spread out over the available number of MPI processes as much as possible while also trying to balance the number of calculations performed per MPI process. Since all emulator systems are stored in different HDF5-files, it is possible to reinitialize the *Pipeline* using the same *Emulator* class and *ModelLink* subclass on a different number of MPI processes. To make sure that the results are not influenced by the variation in evaluation rates, we constructed an emulator of the Gaussian model and used the exact same emulator in every test.

The tests were carried out using any number of MPI processes between 1 and 32, and using a single OpenMP thread each time for consistency. We generated a Latin-Hypercube design of $3 \cdot 10^6$ samples and measured the average evaluation rate of the emulator using the same Latin-Hypercube design each time. To take into account any variations in the evaluation rate caused by initializations, this test was performed 20 times. As a result, this Latin-Hypercube design was evaluated in the emulator a total of 640 times, giving an absolute total of $1.92 \cdot 10^9$ emulator evaluations.

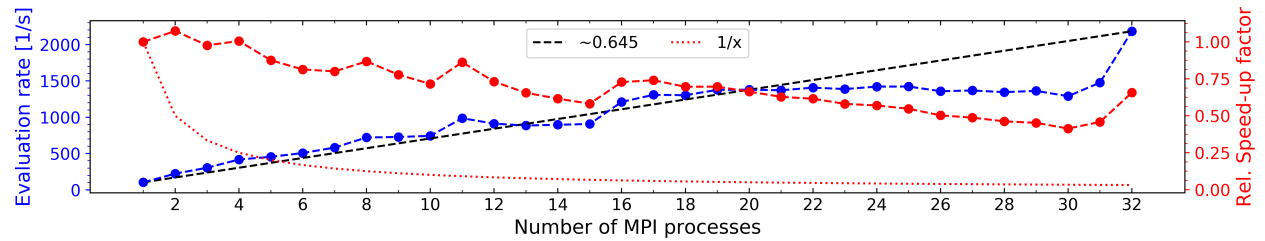


Fig. 3.2: Figure showing the MPI scaling of *PRISM* using the emulator of a simple Gaussian model with 32 emulator systems. The tests involved analyzing a Latin-Hypercube design of $3 \cdot 10^6$ samples in the emulator, determining the average evaluation rate and executing this a total of 20 times using the same sample set every time. The emulator used for this was identical in every instance. **Left axis:** The average evaluation rate of the emulator vs. the number of MPI processes it is running on. **Right axis:** The relative speed-up factor vs. the number of MPI processes, which is defined as $\frac{f(x)}{f(1) \cdot x}$ with $f(x)$ the average evaluation rate and x the number of MPI processes. **Dotted line:** The minimum acceptable relative speed-up factor, which is always $1/x$. **Dashed line:** A straight line with a slope of ~ 0.645 , connecting the lowest and highest evaluation rates. The tests were performed using the *OzSTAR computing facility* at the Swinburne University of Technology, Melbourne, Australia.

In Fig. 3.2, we show the results of the performed MPI scaling tests. On the left y-axis, the average evaluation rate vs. the number of MPI processes the test ran on is plotted, while the relative speed-up factor vs. the number of MPI processes is plotted on the right y-axis. The relative speed-up factor is defined as $f(x)/(f(1) \cdot x)$ with $f(x)$ the average

evaluation rate and x the number of MPI processes. The ideal MPI scaling would correspond to a relative speed-up factor of unity for all x .

In this figure, we can see the effect of the high-level MPI implementation. Because the emulator systems are spread out over the available MPI processes, the evaluation rate is mostly determined by the runtime of the MPI process with the highest number of systems assigned. Therefore, if the number of emulator systems (32 in this case) cannot be divided by the number of available MPI processes, the speed gain is reduced, leading to the plateaus like the one between $x = 16$ and $x = 31$. Due to the emulator systems not being the same, their individual evaluation rates are different such that a different evaluation rate has a bigger effect on the average evaluation rate of the emulator the more MPI processes there are. This is shown by the straight dashed line drawn between $f(1)$ and $f(32)$, which has a slope of ~ 0.645 .

The relative speed-up factor shows the efficiency of every individual MPI process in a specific run, compared to using a single MPI process. This also shows the effect of the high-level MPI implementation, giving peaks when the maximum number of emulator systems per MPI process has decreased. The dotted line shows the minimum acceptable relative speed-up factor, which is always defined as $1/x$. On this line, the average evaluation rate $f(x)$ for any given number of MPI processes is always equal to $f(1)$.

4.1 Writing a ModelLink subclass

In *Minimal example*, a description is given of how to initialize the *Pipeline* class using a default *ModelLink* subclass. Here, the basic steps for making a custom *ModelLink* subclass are shown.

Lst. 4.1: example_link.py

```
# -*- coding: utf-8 -*-

# Future imports
from __future__ import absolute_import, division, print_function

# Package imports
import numpy as np

# PRISM imports
from prism.modellink import ModelLink

# ExampleLink class definition
class ExampleLink(ModelLink):
    def __init__(self, *args, **kwargs):
        # Perform any custom operations here
        pass

        # Call ModelLink's __init__()
        super().__init__(*args, **kwargs)

    def get_default_model_parameters(self):
        # Define default parameters (not required)
        par_dict = {}
```

(continues on next page)

(continued from previous page)

```

    return(par_dict)

def get_default_model_data(self):
    # Define default data (not required)
    data_dict = {}
    return(data_dict)

# Override call_model method
def call_model(self, emul_i, par_set, data_idx):
    # Perform operations for obtaining the model output
    # Following is provided:
    # 'emul_i': Requested iteration
    # 'par_set': Requested sample(s)
    # 'data_idx': Requested data point(s)
    pass

# Override get_md_var method
def get_md_var(self, emul_i, par_set, data_idx):
    # Perform operations for obtaining the model discrepancy variance
    # Following is provided:
    # 'emul_i': Requested iteration
    # 'par_set': Requested sample
    # 'data_idx': Requested data point(s)
    pass

```

In the *example_link.py* file above, a minimal example of a *ModelLink* subclass is shown. It has two abstract methods that need to be overridden; *call_model()* (wrapper function for calling the model) and *get_md_var()* (calculates the model discrepancy variance). A *ModelLink* subclass cannot be initialized if either method has not been overridden. Given the importance of both methods, detailed descriptions are given in *Wrapping a model (call_model)* and *Model discrepancy variance (md_var)*, respectively.

Since every model is different, with some requiring preparations in order to work properly, the *__init__()* constructor method may be extended to include any custom code to be executed when the subclass is initialized. The superclass version of the *__init__()* method must always be called, as it sets several important flags and properties, but the time at which this is done does not matter. During the initialization of the *Emulator* class, it is checked whether or not the superclass constructor of a provided *ModelLink* instance was called (to avoid this from being forgotten).

Besides executing custom code, three properties/flags can be set in *__init__()*, which have the following default values if the extended constructor does not set them:

```

self.name = self.__class__.__name__ # Set instance name to the name of the class
self.call_type = 'single'            # Request single model calls
self.MPI_call = False                # Request only controller calls

```

The first property, *name*, defines the name of the *ModelLink* instance. This name is used by the *Emulator* class during initialization to check if a constructed emulator is linked to the proper *ModelLink* instance, in order to avoid causing mismatches. If one wants to use the same *ModelLink* subclass for different models (like, using different parameter spaces), it is recommended to add an identifier for this to this name. An example of this can be found in the definition of the *GaussianLink* class, which adds the number of Gaussians in the model to its *name* property.

The other two properties, *call_type* and *MPI_call*, are flags that tell *PRISM* how the *call_model()* method should be used. By default, every model evaluation sample is requested individually in serial, since this would be the most expected behavior. However, this is most likely not enough for sophisticated models, as they can require some preparation (e.g., having to read in data files) or more than a single core (in MPI) to function. Therefore,

`call_type` can be set to accept solely individual samples (`'single'`), solely entire sample sets (`'multi'`) or both (`'hybrid'`). In the same way, `MPI_call` can be set to `True` or `False` to identify that the model needs to be executed in serial or in MPI.

Note: If a model uses OpenMP parallelization, it is recommended to set `MPI_call` to `False` in the `ModelLink` subclass. This allows for all worker ranks to be used in OpenMP threads, while only the controller rank calls the model.

Finally, the `ModelLink` class has three methods that can be overridden for adding utility to the class (of which two are shown in [example_link.py](#)). The `get_default_model_parameters()` and `get_default_model_data()` methods return dictionaries containing the default model parameters and model data to use in this class instance, respectively. By overriding these methods, one can hard-code the use of specific parameters or comparison data, avoiding having to provide them when initializing the `ModelLink` subclass. Additionally, if a default parameter or data point is also provided during initialization, the provided information will override the defaults.

Example

The `GaussianLink` class has default parameters defined:

```
>>> from prism.modellink import GaussianLink
>>> model_data = {3: [3.0, 0.1]}
>>> modellink_obj = GaussianLink(model_data=model_data)
>>> modellink_obj
GaussianLink(model_parameters={'A1': [1.0, 10.0, 5.0], 'B1': [0.0, 10.0, 5.0],
                              'C1': [0.0, 5.0, 2.0]},
              model_data={3: [3.0, 0.1]})
```

Providing a custom set of parameters will override the coded defaults:

```
>>> model_parameters = {'A1': [-5, 7, 2]}
>>> modellink_obj = GaussianLink(model_parameters=model_parameters, model_data=model_
↳ data)
>>> modellink_obj
GaussianLink(model_parameters={'A1': [-5.0, 7.0, 2.0], 'B1': [0.0, 10.0, 5.0],
                              'C1': [0.0, 5.0, 2.0]},
              model_data={3: [3.0, 0.1]})
```

The third method, `get_str_repr()`, is a simple function that returns a list containing the representations of all non-default input arguments the `ModelLink` subclass takes. It can be overridden to add the missing input arguments to the full representation of the class, which is automatically called whenever the representation is requested. The `GaussianLink` class overrides this method to add its `n_gaussians` input argument.

Lst. 4.2: line_link.py

```
# -*- coding: utf-8 -*-

# Future imports
from __future__ import absolute_import, division, print_function

# Package imports
import numpy as np
```

(continues on next page)

(continued from previous page)

```

# PRISM imports
from prism.modellink import ModelLink

# LineLink class definition
class LineLink(ModelLink):
    def __init__(self, *args, **kwargs):
        # No custom operations required
        pass

        # Call ModelLink's __init__()
        super().__init__(*args, **kwargs)

    def get_default_model_parameters(self):
        # Define default parameters (not required)
        par_dict = {
            'A': [-10, 10, 3], # Intercept in [-10, 10] with estimate of 3
            'B': [0, 5, 1.5]} # Slope in [0, 5] with estimate of 1.5
        return(par_dict)

    def get_default_model_data(self):
        # Define default data (not required)
        data_dict = {
            1: [4.5, 0.1], # f(1) = 4.5 +- 0.1
            2.5: [6.8, 0.1], # f(2.5) = 6.8 +- 0.1
            -2: [0, 0.1]} # f(-2) = 0 +- 0.1
        return(data_dict)

    # Override call_model method
    def call_model(self, emul_i, par_set, data_idx):
        # Calculate the value on a straight line for requested data_idx
        vals = par_set['A'] + np.array(data_idx) * par_set['B']
        return(vals)

    # Override get_md_var method
    def get_md_var(self, emul_i, par_set, data_idx):
        # Calculate the model discrepancy variance
        # For a straight line, this value can be set to a constant
        return(1e-4 * np.ones_like(data_idx))

```

Using all the information above and the template given in [example_link.py](#), a *ModelLink* subclass can be written for a straight line model, shown in the [line_link.py](#) file above. Here, all methods discussed before (besides the `get_str_repr()` method, since no additional input arguments are used) have been overridden. Given that this model is very simple, no changes have been made to the instance constructor, `__init__()`. Therefore, only single evaluation samples in serial are requested.

PRISM provides the `test_subclass()` function that allows the user to check if a *ModelLink* subclass is properly written. It returns an instance of the subclass if the test passes, or raises a specific error if not. We can use this function to initialize our newly written subclass:

```

>>> from line_link import LineLink
>>> from prism.modellink import test_subclass
>>> modellink_obj = test_subclass(LineLink)
>>> modellink_obj

```

(continues on next page)

(continued from previous page)

```
LineLink(model_parameters={'A': [-10.0, 10.0, 3.0], 'B': [0.0, 5.0, 1.5]},
         model_data={2.5: [6.8, 0.1], -2: [0.0, 0.1], 1: [4.5, 0.1]})
```

Since no errors were raised, we can now use the initialized *ModelLink* subclass to initialize the *Pipeline* class:

```
>>> from prism import Pipeline
>>> pipe = Pipeline(modellink_obj)
```

4.2 Data identifiers (*data_idx*)

The comparison data points that are given to the *ModelLink* class each require a unique data point identifier, allowing *PRISM* to distinguish between them. This data identifier (called *data_idx*) can be used by the model wrapped in the *call_model()* method as a description of how to calculate/extract the data point. It can be provided as a non-mutable sequence (a Python tuple) of a combination of integers, floats and strings, each element describing a part of the operations required. The data identifier sequence can be of any length, and the length can differ between data points.

Note: If a data identifier is given as a single element, then the identifier is saved as that single element instead of a tuple. For example, *data_idx* = [(1), (2), (3, 4), ...] would be saved as *data_idx* = [1, 2, (3, 4), ...].

In its simplest form, the data identifier is a single value that is given to a function $f(x)$, which is a function that is defined for a given model parameter set and returns the function value belonging to the input x . This is the way the data identifier works for the two standard *ModelLink* subclasses, the *SineWaveLink* and *GaussianLink* classes. It is also used in the *LineLink* class described in the *line_link.py* file above.

For more sophisticated models, a single value/element is not enough to uniquely identify a data point. A simple example of this would be if the model generates a two-dimensional array of values, where one specific value needs to be returned. Then, the data identifier can be given as a tuple of two integers, like *data_idx* = [(1, 1), (4, 8), ...]. In the case that the model also generates several two-dimensional arrays which are named, an extra string could be used to identify this array first: *data_idx* = [('array1', 1, 1), ('array4', 4, 8), ...].

An even more complex example is when a data point needs to be retrieved from a specific named data set at a certain point in a model simulation, after which an operation needs to be carried out (like, making a histogram of the results) and the resulting data point is then found at a specific value in that histogram. The histogram here might only be necessary to make for specific data sets, while different operations are required for others. *PRISM* allows for such complex data identifiers to be given, as it treats every sequence of data identifier elements as separated. Two different data identifiers working as described above can for example be written as *data_idx* = [(14, 'array1', 'histogram', 7.5), (17, 'array7', 'average'), ...], where the first data point requires an extra (float) value for the histogram and the second does not. In order to do this, one would of course be required to make sure that the *call_model()* method can perform these operations when provided with the proper data identifier.

4.3 Wrapping a model (*call_model*)

The *call_model()* method is the most important method in the entire *PRISM* package. It provides the *Pipeline* instance with a way to call the model that is wrapped in the user-defined *ModelLink* subclass. For *PRISM*, this method is a black box: it takes a parameter/sample set, performs a series of unknown operations and returns the values corresponding to the requested data points and sample(s). Therefore, the *call_model()* method must be written with great care.

4.3.1 Input arguments

Depending on the values of the `multi_call` and `MPI_call` flags (where the first is set by the `call_type` flag), the `Pipeline` instance will use the `call_model()` method differently. As explained in [Writing a ModelLink subclass](#), every model evaluation sample is requested individually in serial by default, which corresponds to `multi_call` is `False` and `MPI_call` is `False`. When single-calling a model, *PRISM* expects an array-like container back with shape `(n_data)`, where the order of the elements is the same as the order of the requested `data_idx`. If we assume that we have an instance of the `LineLink` class (introduced in [line_link.py](#)) called `modellink_obj` and want to evaluate the model three times for all data points, then the model would be called as (solely by the controller rank):

```
# Get emul_i, sam_set and data_idx
emul_i = 1
sam_set = np.random.rand(3, modellink_obj.n_par)
data_idx = modellink_obj.data_idx

# Evaluate model
mod_set = np.zeros([sam_set.shape[0], len(data_idx)])
for i, par_set in enumerate(sam_set):
    par_dict = sdict(zip(modellink_obj.par_name, par_set))
    mod_set[i] = modellink_obj.call_model(emul_i=emul_i,
                                         par_set=par_dict,
                                         data_idx=data_idx)
```

Here, we looped through the entire sample set one-by-one, converted every individual sample to a (sorted) dict and called the model with it. The emulator iteration is given as a normal integer and the data identifiers `data_idx` is provided as a list of individual data identifiers (which are either single elements or tuples of elements, as described in [Data identifiers \(data_idx\)](#)). The requested data identifiers are not necessarily the same as those given in `data_idx`. An individual sample provided in this way will be of the form:

```
par_dict = {'par_1_name': par_1_val,
            'par_2_name': par_2_val,
            ...,
            'par_n_name': par_n_val}
```

An example of this would be `par_dict = {'A': 1.0, 'B': 2.0}` for the `LineLink` class. This works very well for models that do not require any preparation before they can start evaluating and requires a minimal amount of effort to implement. However, if the sample set is very large, then evaluating the model in this fashion can be inefficient due to many memory look-ups.

Therefore, the `GaussianLink` class accepts both single and multi-calls. When multi-calling a model, *PRISM* expects an array-like container back with shape `(n_sam, n_data)`, where the order of the columns is the same as the order of the requested `data_idx`. So, if we use the same example again, but this time have an instance of the `GaussianLink` class with `multi_call` is `True`, then the model would be called as (again solely by the controller rank):

```
# Get emul_i, sam_set and data_idx
emul_i = 1
sam_set = np.random.rand(3, modellink_obj.n_par)
data_idx = modellink_obj.data_idx

# Evaluate model
sam_dict = sdict(zip(modellink_obj.par_name, sam_set.T))
mod_set = modellink_obj.call_model(emul_i=emul_i,
                                   par_set=sam_dict,
                                   data_idx=data_idx)
```

This call is roughly the same as before, but this time the entire sample set is provided as a (sorted) dict instead of individual samples. The lay-out of this sample dict is of the form:

```
sam_dict = {'par_1_name': [par_1_val_1, par_1_val_2, ..., par_1_val_m],
            'par_2_name': [par_2_val_1, par_2_val_2, ..., par_2_val_m],
            ...,
            'par_n_name': [par_n_val_1, par_n_val_2, ..., par_n_val_m]}
```

Again, in the case of the *GaussianLink* class, this sample dict could look like `sam_dict = {'A1': [1.0, 5.5, 10.0], 'B1': [0.0, 5.0, 10.0], 'C1': [0.0, 2.5, 5.0]}`. This can be used when the model requires some kind of preparation before being able to perform evaluations, or when it is simply more efficient to provide all requested samples at once (like for the *GaussianLink* class).

Note: If a model uses OpenMP parallelization, it is recommended to set `MPI_call` to `False` in the *ModelLink* subclass. This allows for all worker ranks to be used in OpenMP threads, while only the controller rank calls the model.

Note: If one wishes to transform the received `sam_dict` back into a normal NumPy array of shape `(n_sam, n_par)`, this can be done quite easily by executing `sam_set = np.array(par_set.values()).T`, where `par_set` is the `sam_dict` provided to the `call_model()` method. Keep in mind that doing so means that the columns are sorted on the names of the model parameters.

New in version 1.1.2: It is also possible to make `call_model()` return a dict instead, where it has the identifiers in the requested `data_idx` as its keys and scalars (single-call) or 1D array-likes of shape `(n_sam)` (multi-call) as its values. *PRISM* will automatically convert the dict back to the array-like container format that is normally expected.

When the `MPI_call` flag is set to `True`, the calls to the `call_model()` method are almost the same as described above. The only difference is that all ranks call the method (each providing the same `emul_i`, `par_dict/sam_dict` and `data_idx`) instead of just the controller rank.

4.3.2 Multi-calling

When the `multi_call` flag is set to `False`, the `call_model()` method is most likely nothing more than a simple function. But, when `multi_call` is set to `True`, `call_model()` can be a lot more complex. An example of this would be if we tried to make an emulator of an emulator (which is possible, but completely pointless). In this case, it would be necessary for the “model” (as we are going to call the emulated emulator from now on) to be loaded into memory first before it can be evaluated. Although loading an emulator into memory usually does not take that long, we do not want to do this for every single “model” evaluation. Besides, evaluating an emulator is much quicker when all samples are evaluated at once (due to the way the `_evaluate_sam_set()` method is written).

So, therefore, it is necessary to use `multi_call` is `True` for this “model”. If we assume that we have already made an emulator of the *LineLink* class, then, the `call_model()` method could be written as:

```
def call_model(self, emul_i, par_set, data_idx):
    # Initialize Pipeline object as a model
    modellink_obj = LineLink()
    pipe_model = Pipeline(modellink_obj, working_dir='linelink_0')

    # Call pipe_model
    mod_set = pipe_model.evaluate(par_set, emul_i)['adj_exp_val']

    # Make sure only the requested data points are kept
```

(continues on next page)

(continued from previous page)

```

req_idx = [pipe_model.emulator._data_idx[emul_i].index(idx) for idx in data_idx]
mod_set = mod_set[:, req_idx]

# Return mod_set
return(mod_set)

```

Here, we only initialize the “model” once per model call, and then evaluate all samples in it by using the `evaluate()` method (which can take sample dicts as a valid input argument). This returns a dict of the evaluation results, where we are only interested in the adjusted expectation values. Note that making an emulator of an emulator is pointless, but used here as an example.

Note: Due to the way *PRISM* is written, it is technically speaking not necessary to reinitialize the *Pipeline* class every time that `call_model()` is called. It is possible to initialize it when the corresponding *ModelLink* subclass is initialized and keep it in memory. The code above would however be necessary if the “model” works in the same way as *PRISM*’s *worker_mode*, where all worker ranks are listening for calls until the “model” is finalized. This finalization would be required in order to give *PRISM* control back over all ranks.

4.3.3 Backing up progress

New in version 1.1.1.

Warning: This feature is still experimental and it may see significant changes or be (re)moved in the future.

In *PRISM*, an emulator system is constructed by calculating all required components individually. This means that the construction process of an emulator iteration can easily be interrupted and restored at a later time, only losing the progress that was made in the current step (e.g., interrupting construction during the calculation of the covariance matrix will lose progress made there, but not the already previously finished steps). This system was implemented to accommodate for *PRISM* running on clusters, where the construction is more prone to interruptions due to, for example, jobs timing out, and to allow for *PRISM* to be loaded unto any number of MPI processes.

However, the biggest step in the construction of all emulator systems, is the evaluation of the model. Since the evaluation of the model is carried out by the `call_model()` method, *PRISM* has no control over what is happening until this method gives control back to the *Pipeline* instance (by returning the requested data points). Therefore, automated backups of already calculated data points cannot be performed by *PRISM* itself, running the risk that many CPU hours are wasted if a job on a cluster takes longer than initially expected and times out. While this could be avoided if the user writes its own backup system, this would require more work from the user, which clashes with *PRISM*’s ease-of-use policy.

Therefore, the *ModelLink* class implements its own (experimental) backup system based on the *hickle* package, given by the `_make_backup()` and `_read_backup()` methods. This backup system is best used for models that are multi-called (`multi_call` set to `True`), as made backups will replace previous ones (of the same type). The `_make_backup()` method is meant to be used from within the `call_model()` method and will not work if called anywhere else. Attempting to call it incorrectly (e.g., not from within `call_model()` or with incorrect arguments), will raise a *RequestWarning* and simply return without doing anything, rather than raising a *RequestError*. This is to make sure that using it incorrectly does not disrupt the `call_model()` call, as that has the exact opposite effect of what the backup system tries to achieve.

The `_make_backup()` method takes two arguments, `*args` and `**kwargs`, of which at least one is required. Calling it from within the `call_model()` method will produce an HDF5-file containing the `emul_i`, `par_set` and `data_idx` argument values that were used to call `call_model()` with, and the supplied `*args` and

****kwargs**. The name of the HDF5-file contains the values of `emul_i` and `name`, and will be saved in the current working directory (NOT the emulator working directory, as the `ModelLink` instance has no access to its path). The backup can be read in by passing the value of `emul_i` to the `_read_backup()` method of the corresponding `ModelLink` instance, which will return a dict containing the values of the five arguments that were saved to the file.

Backups can be made at any point during the execution of `call_model()`, and basically all types of objects are compatible and can be viewed freely in the HDF5-file. It is possible that instances of certain custom classes may not be supported by the `hickle` package, in which case they will be pickled and saved as a string, causing them to not be able to be viewed freely (but they can still be backed up). Depending on the size of the data provided, it can sometimes take a little while before a backup is made. Therefore, it is probably best to trigger making backups at specified progress points in `call_model()`.

To illustrate how this backup system can be used, assume that we have written a `ModelLink` subclass, which requires some preparation before it can start evaluating the wrapped model. Here, we will assume that this preparation is provided by a function called `prepare_model()`, which returns an instance of some class that can be used to evaluate the model after the preparation is completed. Then, we could incorporate the backup system by writing a `call_model()` method like this:

```
def call_model(self, emul_i, par_set, data_idx):
    # Prepare the model for evaluation
    model = prepare_model()

    # Controller performs evaluations
    if model.is_controller:
        # Initialize empty array of results
        mod_set = np.zeros([len(par_set['par1']), len(data_idx)])

        # Convert par_set to a generator
        sam_set = map(lambda *args: args, *par_set.values())

        # Call model for every individual sample in sam_set
        for i, sam in enumerate(sam_set):
            mod_set[i] = model.evaluate(sam, data_idx)

            # Make a backup every 500 evaluations
            if not ((i+1) % 500):
                self._make_backup(mod_set=mod_set[:i])

    # Finalize the model
    model.finalize()

    # Return the results on the controller
    if model.is_controller:
        return(mod_set)
```

The code above shows an example of a model that needs to be initialized before it can be multi-called in MPI, and needs to be finalized afterward. Since such a model is probably quite complex, it may be a good idea to make a backup every once in a while. Therefore, whenever 500 evaluations have been done, a backup is made of all results gained up to that point. This means that whenever the model evaluation process is interrupted, a maximum of the last 500 evaluations is lost. The evaluations that are not lost can be loaded back in by using the `_read_backup()` method, and potentially (after a bit of formatting) be passed to the `ext_real_set` input argument of the `construct()` method when attempting to construct the emulator iteration again.

Note that if `model.evaluate()` was implemented such that it takes the entire sample set at once rather than one at a time, calling `_make_backup()` in `model.evaluate()` works perfectly fine, as long as `model.evaluate()` is always called by `call_model()` or any other function for which this is true. Put a little bit more simple: `_make_backup()` must be called either directly or indirectly by `call_model()`, as shown in the following example.

Example

```
def call_model(self, emul_i, par_set, data_idx):
    # Call a function A and return its output
    # This function does not require emul_i, so do not provide it
    return A(self, par_set, data_idx)

def A(modellink_obj, par_set, data_idx):
    # Prepare model
    model = prepare_model()

    # Prepare par_set for evaluation
    sam_set = map(lambda *args: args, *par_set.values())

    # Call a function B
    mod_set = B(modellink_obj, model, sam_set, data_idx)

    # Finalize the model
    model.finalize()

    # Return the results
    return(mod_set)

def B(modellink_obj, model_obj, sam_set, data_idx):
    # Prepare mod_set
    mod_set = np.zeros([np.shape(sam_set)[0], len(data_idx)])

    # Call model for every individual sample in sam_set
    for i, sam in enumerate(sam_set):
        mod_set[i] = model_obj.evaluate(sam, data_idx)

        # Make a backup every 500 evaluations
        if not((i+1) % 500):
            modellink_obj._make_backup(mod_set=mod_set[:i+1])

    # Return mod_set
    return(mod_set)
```

4.4 Model discrepancy variance (*md_var*)

Of the three different variances that are used for calculating the implausibility values of a parameter set, the *model discrepancy variance* is by far the most important. The model discrepancy variance describes all uncertainty about the correctness of the model output that is caused by the model itself. This includes the accuracy of the code implementation, completeness of the inclusion of the involved physics, made assumptions and the accuracy of the output itself, amongst others. It therefore acts as a measure of the quality of the model that is being emulated by *PRISM*, and as with `call_model()`, must be handled with great care.

4.4.1 Theory

When *PRISM* constructs an emulator, it attempts to make a perfect approximation of the model that covers the absolute plausible regions of parameter space. This perfect approximation would be reached if the adjusted emulator variance (*adj_var*) is zero for all samples. In this case, the emulator has the same variance associated with it as the model, which

is given by the model discrepancy variance. Therefore, if the model discrepancy variance is determined incorrectly, the emulator itself will be incorrect as well.

The reason for this is as follows. The implausibility value of a parameter set states how many standard deviations the emulator system expects the model realization corresponding to this parameter set, to be away from explaining the model comparison data. When the total variance increases, the implausibility value decreases (since less standard deviations fit in the total difference). For an emulator system that is still very inaccurate (e.g., first iteration), the adjusted emulator variance dominates over the other two variances. However, later on, the adjusted emulator variance becomes less and less dominant, causing the other two variances to start playing a role. In most cases, it is safe to assume that the model discrepancy variance is higher than the observational variance, since a model would be fitting noise if this was not the case. Therefore, there is going to be a moment when the model discrepancy variance starts being close to the adjusted emulator variance.

When this happens, the plausible region of parameter space starts being determined by the model discrepancy variance. If the model discrepancy variance is generally higher than it should be, then this will often result into the emulator system not converging as far as it could have, since parts of parameter space are still marked as plausible. The opposite however (the model discrepancy variance generally being lower than it should be) can mark parts of parameter space as implausible while they are not. This means that these parts are removed from the emulator.

From the above, it becomes clear that overestimating the model discrepancy variance is much less costly than underestimating its value. It is therefore important that this variance is properly described at all times. However, since the description of the model discrepancy variance can take a large amount of time, *PRISM* uses its own default description in case none was provided, which is defined as $\text{Var}(\epsilon_{\text{md},i}) = (z_i/6)^2$, where $\text{Var}(\epsilon_{\text{md},i})$ is the model discrepancy variance of a specified model comparison data point i and z_i is the corresponding data value. If one assumes that a model output within half of the data is considered to be acceptable, with acceptable being defined as the 3σ -interval, then the model discrepancy variance is obtained as:

$$\begin{aligned} [z_i - 3\sigma, z_i + 3\sigma] &= \left[\frac{1}{2}z_i, \frac{3}{2}z_i \right], \\ 6\sigma &= z_i, \\ \sigma &= \frac{z_i}{6}, \\ \text{Var}(\epsilon_{\text{md},i}) &= \sigma^2 = \left(\frac{z_i}{6} \right)^2. \end{aligned}$$

This description of the model discrepancy variance usually works well for simple models, and acts as a starting point within *PRISM*. When models become bigger and more complex, it is likely that such a description is not enough. Given that the model discrepancy variance is unique to every model and might even be different for every model output, *PRISM* cannot possibly cover all scenarios. It is therefore advised that the model discrepancy variance is provided externally by the user.

4.4.2 Implementation

The model discrepancy variance is given by the `get_md_var()` method. This method is, like `call_model()`, an abstract method and must be overridden by the `ModelLink` subclass before it can be initialized. The `get_md_var()` method is called every time the implausibility value of an emulator evaluation sample is determined. Unlike the `call_model()` method, the `get_md_var()` method is called by individual emulator systems, as they determine implausibility values individually.

For this reason, the `get_md_var()` method is provided with the emulator iteration `emul_i`, a single parameter set `par_set` and the data identifiers requested by the emulator system `data_idx`. The `call_type` and `MPI_call` flags have no influence on the way the `get_md_var()` method is used, as it is always called in serial for a single parameter set. When it is called, *PRISM* expects an array-like container back with shape `(n_data)` (if 1σ -interval is centered) or shape `(n_data, 2)` (if 1σ -interval is given by upper and lower errors), where the order of the elements is the same as the order of the requested `data_idx`. The default model discrepancy variance description given above is used if the `get_md_var()` method raises a `NotImplementedError`, but this is discouraged.

Warning: Because the `get_md_var()` method is always called for single parameter sets, it is important that it can be called without requiring any preparation of data or models.

New in version 1.1.2: It is also possible to make `get_md_var()` return a dict instead, where it has the identifiers in the requested `data_idx` as its keys and scalars (centered) or 1D array-likes of shape `(2)` (non-centered) as its values. *PRISM* will automatically convert the dict back to the array-like container format that is normally expected.

Here, various different aspects of how the *PRISM* package can be used are described.

5.1 Minimal example

A minimal example on how to initialize and use the *PRISM* pipeline is shown here. First, one has to import the *Pipeline* class and a *ModelLink* subclass:

```
>>> from prism import Pipeline
>>> from prism.modellink import GaussianLink
```

Normally, one would import a custom-made *ModelLink* subclass, but for this example one of the two *ModelLink* subclasses that come with the *PRISM* package is used (see [Writing a ModelLink subclass](#) for the basic structure of writing a custom *ModelLink* subclass).

Next, the *ModelLink* should be initialized, which is the *GaussianLink* class in this case. In addition to user-defined arguments, every *ModelLink* subclass takes two optional arguments, *model_parameters* and *model_data*. The use of either one will add the provided parameters/data to the default parameters/data defined in the class. Since the *GaussianLink* class does not have default data defined, it is required to supply it with some data during initialization (using an array, dict or external file):

```
>>> # f(3) = 3.0 +- 0.1, f(5) = 5.0 +- 0.1, f(7) = 3.0 +- 0.1
>>> model_data = {3: [3.0, 0.1], 5: [5.0, 0.1], 7: [3.0, 0.1]}
>>> modellink_obj = GaussianLink(model_data=model_data)
```

Here, the *GaussianLink* class was initialized by giving it three custom data points and using its default parameters. One can check this by looking at the representation of this *GaussianLink* object:

```
>>> modellink_obj
GaussianLink(model_parameters={'A1': [1.0, 10.0, 5.0], 'B1': [0.0, 10.0, 5.0],
                              'C1': [0.0, 5.0, 2.0]}),
              model_data={7: [3.0, 0.1], 5: [5.0, 0.1], 3: [3.0, 0.1]})
```

The *Pipeline* class takes several optional arguments, which are mostly paths and the type of *Emulator* class that must be used. It also takes one mandatory argument, which is an instance of the *ModelLink* subclass to use. Since it has already been initialized above, the *Pipeline* class can be initialized:

```
>>> pipe = pipeline(modellink_obj)
>>> pipe
Pipeline(GaussianLink(model_parameters={'A1': [1.0, 10.0, 5.0], 'B1': [0.0, 10.0, 5.
→0],
                                'C1': [0.0, 5.0, 2.0]}),
          model_data={7: [3.0, 0.1], 5: [5.0, 0.1], 3: [3.0, 0.1]}),
          working_dir='prism_0')
```

Since no working directory was provided to the *Pipeline* class and none already existed, it automatically created one (*prism_0*).

PRISM is now completely ready to start emulating the model. The *Pipeline* allows for all steps in a full cycle (see *PRISM pipeline*) to be executed automatically:

```
>>> pipe.run()
```

which is equivalent to:

```
>>> pipe.construct(analyze=False)
>>> pipe.analyze()
>>> pipe.project()
```

This will construct the next iteration (first in this case) of the emulator, analyze it to check if it contains plausible regions and make projections of all active parameters. The current state of the *Pipeline* object can be viewed by calling the *details()* method (called automatically after most user-methods), which gives an overview of many properties that the *Pipeline* object currently has.

This is all that is required to construct an emulator of the model of choice. All user-methods, with one exception (*evaluate()*), solely take optional arguments and perform the operations that make the most sense given the current state of the *Pipeline* object if no arguments are given. These arguments allow for one to modify the performed operations, like reconstructing/reanalyzing previous iterations, projecting specific parameters, evaluating the emulator and more.

5.2 Projections

After having made an emulator of a given model, *PRISM* can show the user the knowledge it has about the behavior of this model by making *projections* of the active parameters in a specific emulator iteration. These projections are created by the *project()* method, which has many different properties and options. For showing them below, the same emulator as the one in *Minimal example* is used.

5.2.1 Properties

Projections (and their figures) are made by analyzing a large set of evaluations samples. For 3D projections, this set is made up of a grid of *proj_res* x *proj_res* samples for the plotted (active) parameters, where the values for the remaining parameters in every individual grid point are given by an LHD of *proj_depth* samples. This gives the total number of analyzed samples as *proj_res* x *proj_res* x *proj_depth*.

Every sample in the sample set is then analyzed in the emulator, saving whether or not this sample is plausible and what the implausibility value at the first cut-off is (the first value in *impl_cut*). This yields *proj_depth* results per grid point, which can be used to determine the fraction of samples that is plausible and the minimum implausibility

value at the first cut-off in this point. Doing this for the entire grid and interpolating them, creates a map of results that is independent of the values of the non-plotted parameters. For 2D projections, it works the same way, except that only a single active parameter is plotted.

Note: When using a 2D model, the *projection depth* used to make a 2D projection will be *proj_depth*, which is to be expected. However, when using an nD model, the projection depth of a 2D projection is equal to *proj_res* x *proj_depth*. This is to make sure that for an nD model, the density of samples in a 2D projection is the same as in a 3D projection.

The *project()* method solely takes optional arguments. Calling it without any arguments will produce six projection figures: three 2D projections and three 3D projections. One of each type is shown below.

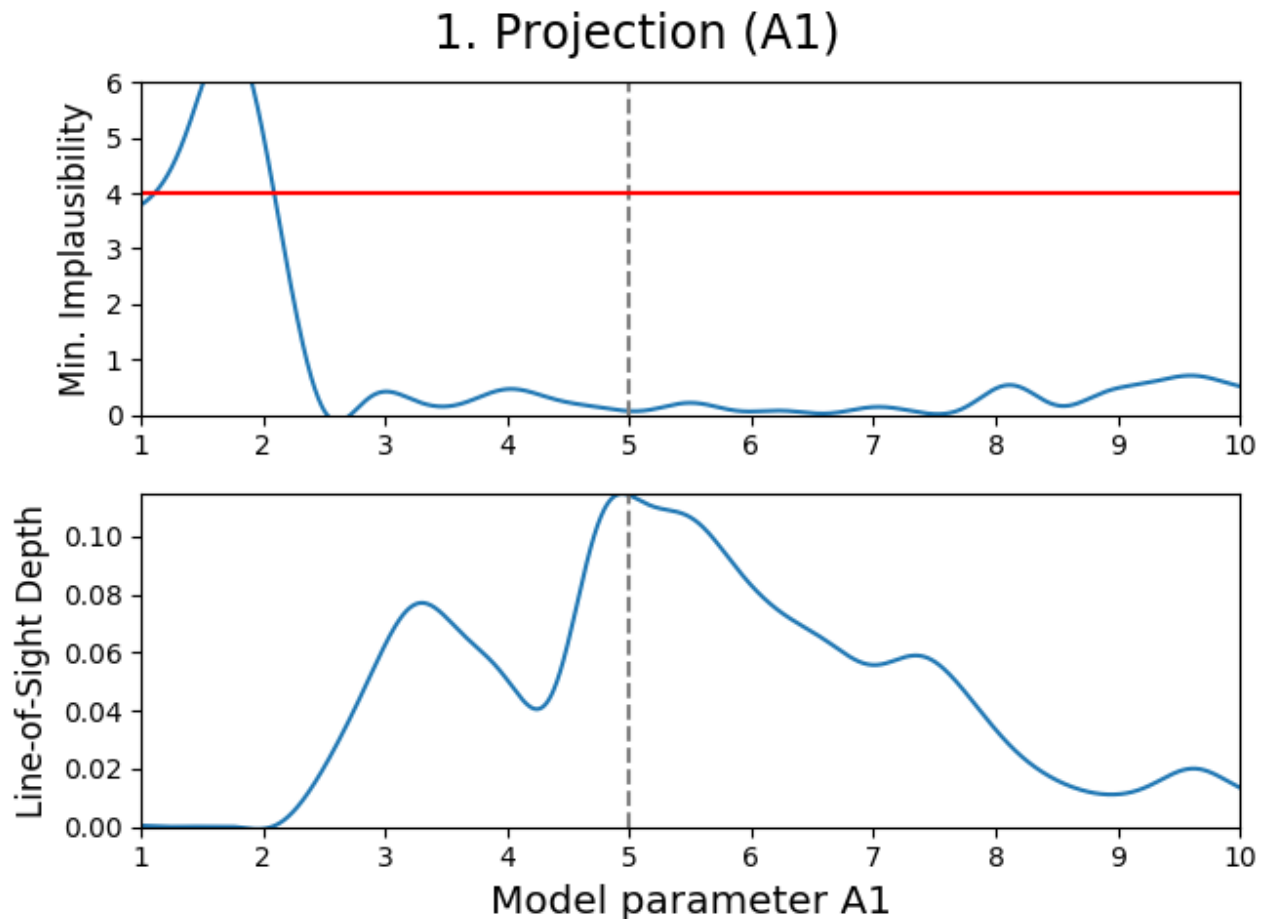


Fig. 5.1: 2D projection figure of model parameter A_1 . The vertical dashed line shows the parameter estimate of A_1 , whereas the horizontal red line shows the first implausibility cut-off value.

A projection figure is made up of two subplots. The upper subplot shows a map of minimum implausibility values that can be reached for any given value (combination) of the plotted parameter(s). The lower subplot gives a map of the fraction of samples that is plausible in a specified point on the grid (called “line-of-sight depth” due to the way it is calculated). Another way of describing this map is that it gives the probability that a parameter set with given plotted value(s) is plausible.

Both projection types have a different purpose. A 3D projection gives insight into what the dependencies (or correlations) are between the two plotted parameters, by showing where the best (top) and most (bottom) plausible samples

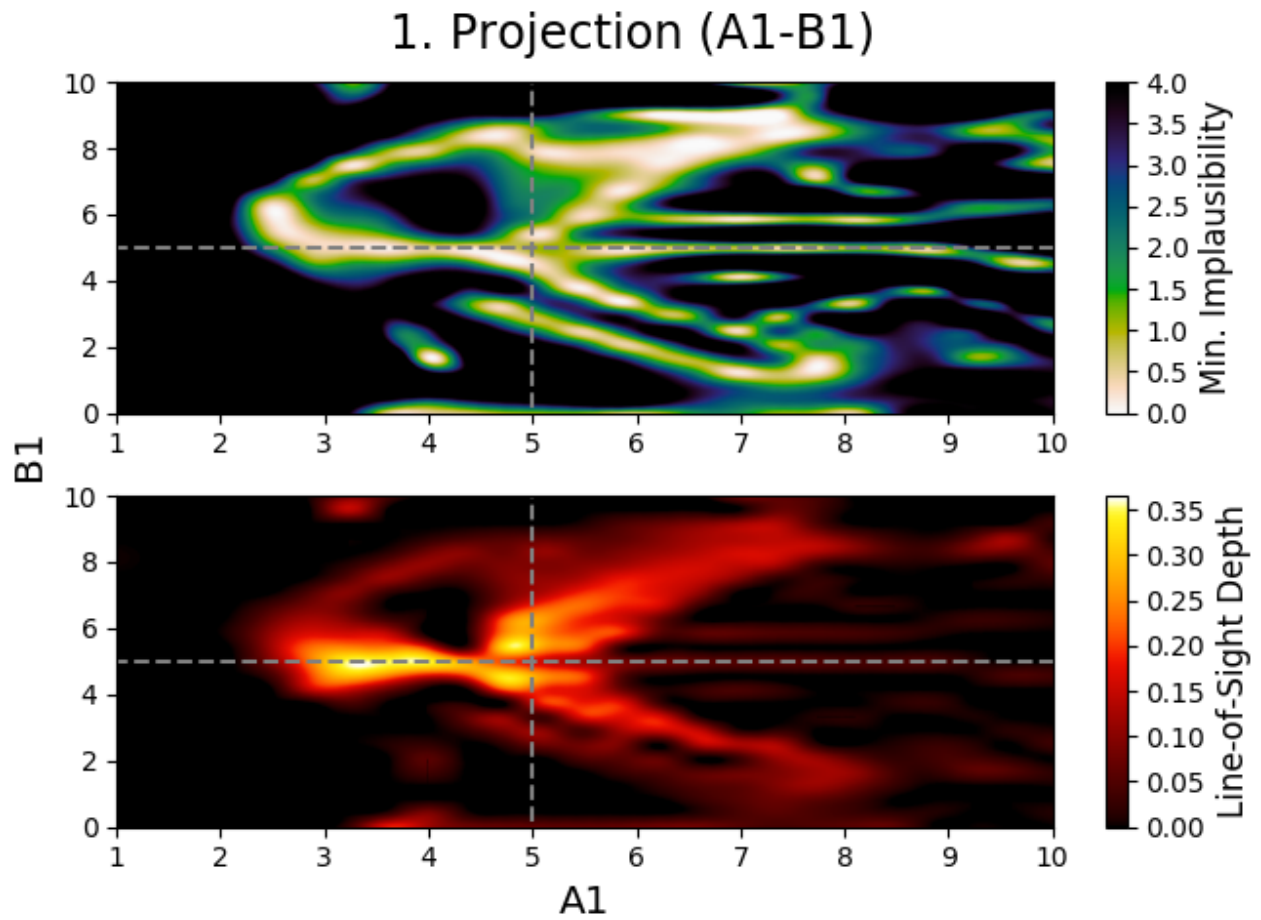


Fig. 5.2: 3D projection figure of model parameters A_1 and B_1 . The dashed lines show the estimates of both parameters.

can be found. On the other hand, a 2D projection is quite similar in meaning to a *maximum likelihood* optimization performed by MCMC methods, with the difference being that the projection is based on expectations rather than real model output. A combination of both subplots allows for many model properties to be derived, especially when they do not agree with each other.

5.2.2 Options

The `project()` method takes two (optional) arguments, `emul_i` and `proj_par`. The first controls which emulator iteration should be used, while the latter provides the model parameters of which projections need to be made. Since it only makes sense to make projections of active parameters, all passive parameters are filtered out of `proj_par`. The remaining parameters are then used to determine which projections are required (which also depends on the requested projection types). For example, if one wishes to only obtain projections of the A_1 and B_1 parameters (which are both active) in iteration 1, then this can be done with:

```
>>> pipe.project(1, ('A1', 'B1'))
```

This would generate the figures shown above, as well as the 2D projection figure of B_1 . By default, the last constructed emulator iteration and all model parameters are requested.

The remaining input arguments can only be given as keyword arguments, since they control many different aspects of the `project()` method. The `proj_type` argument controls which projection types to make. For 2D models, this is always `'2D'` and cannot be modified. However, for nD models, this can be set to `'2D'` (only 2D projections), `'3D'` (only 3D projections) or `'both'` (both 2D and 3D projections). By default, it is set to `'both'`.

The `figure` argument is a bool, that determines whether or not the projection figures should be created after calculating the projection data. If `True`, the projection figures will be created and saved, which is done by default. If `False`, the data that is contained within the projection figures will be calculated and returned in a dict. This allows the user to either let *PRISM* create the projection figures using the standard template or create the figures themselves.

The `align` argument controls the alignment of the subplots in every projection figure. By default, it aligns the subplots in a column (`'col'`), as shown in the figures above. Aligning the subplots in a row (`'row'`) would give Fig. 5.1 as the figure below.

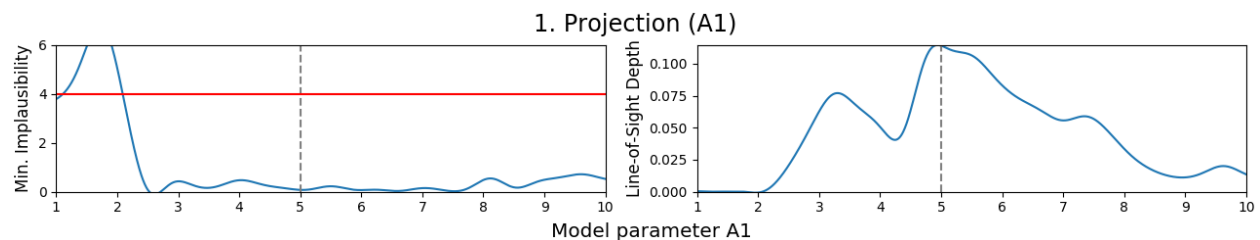


Fig. 5.3: 2D projection figure of model parameter A_1 with the `'row'` alignment.

New in version 1.1.2: The `show_cuts` argument is also a bool, that determines whether to show all implausibility cut-off values in 2D projections (`True`) or only the first cut-off value (`False`, default). In some cases, this may be useful when the first cut-off is not definitive in accepting or rejecting parameter values (as explained below for the *smooth* parameter).

The `smooth` argument is yet another bool, that determines what to do if a grid point in the projection figure contains no plausible samples, but does contain a minimum implausibility value below the first non-wildcard cut-off. If `False`, which is the default, these values are kept in the figure, which may show up as artifact-like features. If `True`, these values are set to the first cut-off, basically removing them from the projection figure. This may however also remove interesting features. Below are two identical projections, one that is smoothed and one that is not, to showcase this difference (these projections are from the second iteration, since this effect rarely occurs in the first iteration).

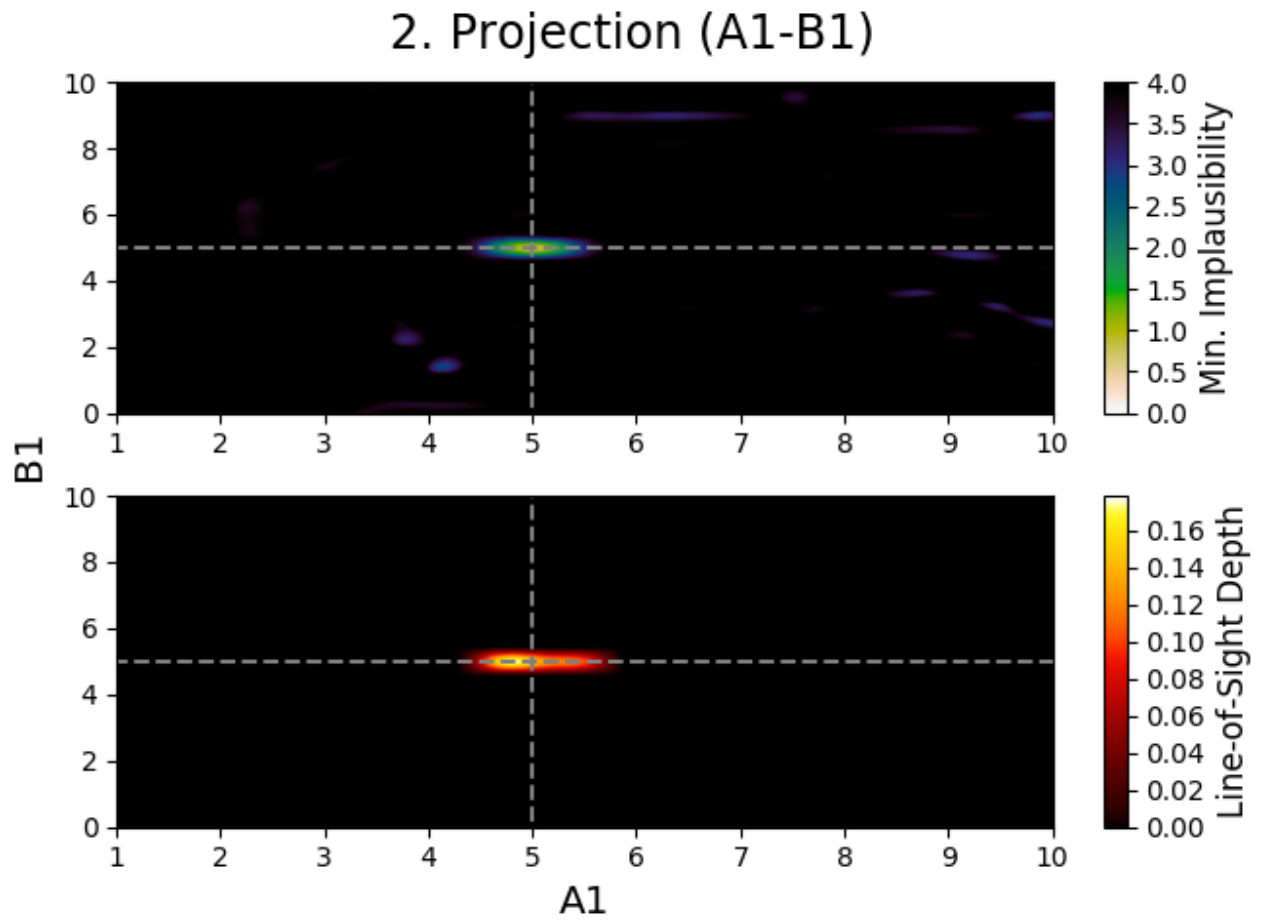


Fig. 5.4: Non-smoothed 3D projection figure of model parameters A_1 and B_1 .

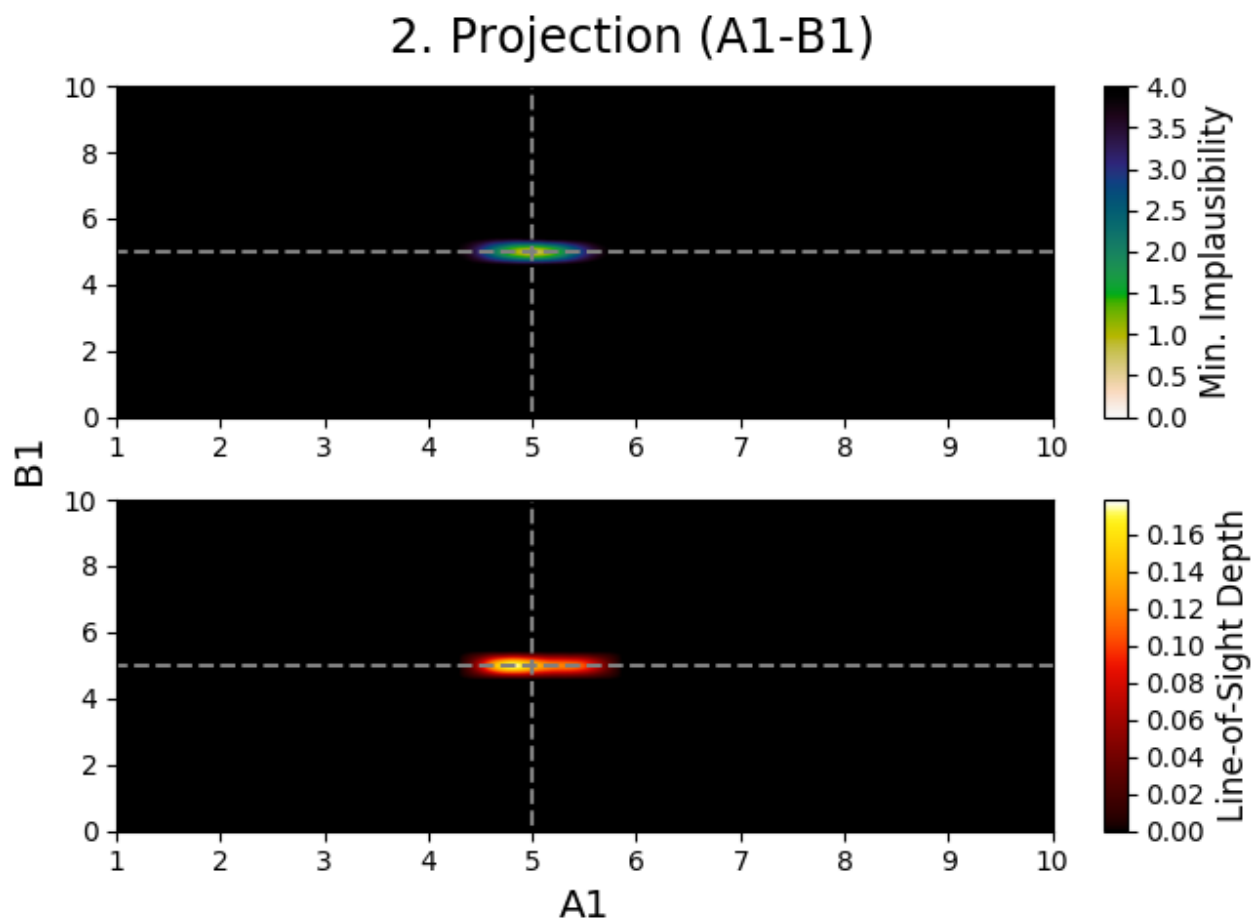


Fig. 5.5: Smoothed 3D projection figure of model parameters A_1 and B_1 .

In these figures, one can see that the non-smoothed projection shows many features in the upper subplot that look like artifacts. These features are however not artifacts, but caused by a sample (or samples) having its highest implausibility value being below the first implausibility cut-off, but still being implausible due to failing a later cut-off. For example, if the implausibility cut-offs are `[4.0, 3.7, 3.5]` and a sample has implausibility values `[3.9, 3.8, 3.2]`, it is found implausible due to failing to meet the second cut-off. However, since the first value is still the highest implausibility value, that value is used in the projection figure. Smoothing figures usually allows for 3D projections (2D projections rarely show this) to become less crowded, but they do throw away information. It should therefore only be used when necessary.

The *force* argument is a bool, which controls what to do if a projection is requested for which data already exists. If `False` (default), it will use the previously acquired projection data to create the projection figure if it does not exist, skip if it does or return the figure data if *figure* is `False`. If `True`, the projection data and all associated projection figures will be deleted, and the projection will be recalculated.

The remaining seven arguments are keyword argument dicts, that need to be passed to the various different plotting functions that are used for creating the projection figures. The *fig_kwargs* dict is passed to the `figure()` function when creating the projection figure instance. The *impl_kwargs_2D* and *los_kwargs_2D* dicts are passed to the `plot()` function when making the minimum implausibility and line-of-sight depth subplots, respectively, for the 2D projections. Similarly, the *impl_kwargs_3D* and *los_kwargs_3D* dicts are passed to the `hexbin()` function for 3D projections. And, finally, the *line_kwargs_est* and *line_kwargs_cut* dicts are passed to the `draw()` function for drawing the parameter estimate and implausibility cut-off lines.

5.3 Dual nature (normal/worker mode)

PRISM features a high-level MPI implementation, as described in [MPI implementation](#): all user-methods and most major methods are to be executed by all MPI ranks at the same time, and *PRISM* will automatically distribute the work among the available ranks within this function/method. This allows for *PRISM* to be used with both serial and parallel models, by setting the *MPI_call* flag accordingly, while also allowing for the same code to be used in serial and parallel. However, given that the emulator of *PRISM* can be very useful for usage in other routines, like [Hybrid sampling](#), an external code will call *PRISM*'s methods. In order to use *PRISM* in parallel with a parallelized model, this code would have to call *PRISM* with all MPI ranks simultaneously at all times, which may not always be possible (e.g., when using MCMC methods).

Therefore, *PRISM* has a *dual execution/call nature*, where it can be switched between two different modes. In the default mode, *PRISM* works as described before, where all MPI ranks call the same user-code. However, by using the *worker_mode* context manager, all code within will be executed in worker mode. When in worker mode, all worker ranks are continuously listening for calls from the controller rank, made with the `_make_call()` method. They will continue to do so until the controller makes the `pipe._make_call(None)` call, which is done automatically when the controller exits *worker_mode*. Manually making this call should solely be done in advanced use-cases.

In *worker_mode*, one uses the following structure (assuming that the *Pipeline* instance is called *pipe*):

```
# Code to be executed in default mode

with pipe.worker_mode:
    if pipe.is_controller:
        # Code to be executed in worker mode

# More code to be executed in default mode
```

Note: All code that is inside the *worker_mode* context manager should solely be executed by the controller rank. If not, all worker ranks will execute this code after the controller ranks exits the context manager. Currently, it is not possible to make a context manager handle this automatically (the rejected [PEP 377](#) describes this perfectly).

The `_make_call()` method accepts almost anything that can be called. It can also be used when not in `worker_mode`, in which case it works the exact same way for all MPI ranks. Its sole limitation is that all supplied arguments must be pickleable (e.g., compiled code objects are NOT pickleable due to safety reasons), both when used in `worker_mode` and outside of it. The `copyreg` module can be used to register specific objects to become pickleable (including compiled code objects).

The `worker_mode` can be used in a variety of ways, as described below. It can be used to access any attribute of the `Pipeline` instance:

```
with pipe.worker_mode:
    if pipe.is_controller:
        # Construct first emulator iteration
        pipe._make_call('construct', 1)

        # Print latest constructed emulator iteration
        print(pipe._make_call('emulator._get_emul_i', 1, 0))

        # Make a specific projection with the 'row' alignment
        pipe._make_call('project', 1, (0, 1), align='row')
```

which is equivalent to:

```
# Construct first emulator iteration
pipe.construct(1)

# Print latest constructed emulator iteration
print(pipe.emulator._get_emul_i(1, 0))

# Make a specific projection with the 'row' alignment
pipe.project(1, (0, 1), align='row')
```

The above two code snippets are equal to each other, and the `worker_mode` will most likely be used very rarely in this fashion. However, by supplying the `_make_call()` method with a callable function (that can be pickled), externally defined functions can be executed:

```
# Enable worker mode
with pipe.worker_mode:
    if pipe.is_controller:
        # Import print function that prepends MPI rank to message
        from prism._internal import rprint

        # Make call to use this function
        # Equivalent to 'rprint("Reporting in.")'
        pipe._make_call(rprint, "Reporting in.")
```

This is especially useful when one combines a serial code with *PRISM*, but wants *PRISM* to execute in MPI. An application example of this is *Hybrid sampling*.

5.4 Hybrid sampling

A common problem when using MCMC methods is that it can often take a very long time for MCMC to find its way on the posterior probability distribution function, which is often referred to as the *burn-in phase*. This is because, when considering a parameter set, there is usually no prior information that this parameter set is (un)likely to result into a desirable model realization. This means that such a parameter set must first be evaluated in the model before any probabilities can be calculated. However, by constructing an emulator of the model, one can use it as an additional prior for the posterior probability calculation. Therefore, although *PRISM* is primarily designed to make analyzing

models much more efficient and accessible than normal MCMC methods, it is also very capable of enhancing them. This process is called *hybrid sampling*, which can be performed easily with the `utils` module and will be explained below.

5.4.1 Algorithm

Hybrid sampling allows one to use *PRISM* to first analyze a model's behavior, and later use the gathered information to speed up parameter estimations (by using the emulator as an additional prior in a Bayesian analysis). Hybrid sampling works in the following way:

1. Whenever an MCMC walker proposes a new sample, it is first passed to the emulator of the model;
2. If the sample is not within the defined parameter space, it automatically receives a prior probability of zero (or $-\infty$ in case of logarithmic probabilities). Else, it will be evaluated in the emulator;
3. If the sample is labeled as *implausible* by the emulator, it also receives a prior probability of zero. If it is plausible, the sample is evaluated in the same way as for normal sampling;
4. Optionally, a scaled value of the first implausibility cut-off is used as an exploratory method by adding an additional (non-zero) prior probability. This can be enabled by using the *hybrid* input argument for the `get_inpost_fn()` function.

Since the emulator that *PRISM* makes of a model is not defined outside of the parameter space given by `par_rng`, the second step is necessary to make sure the results are valid. There are several advantages of using hybrid sampling over normal sampling:

- Acceptable samples are guaranteed to be within plausible space;
- This in turn makes sure that the model is only evaluated for plausible samples, which heavily reduces the number of required evaluations;
- No burn-in phase is required, as the starting positions of the MCMC walkers are chosen to be in plausible space;
- As a consequence, varying the number of walkers tends to have a much lower negative impact on the convergence probability and speed;
- Samples with low implausibility values can optionally be favored.

5.4.2 Usage

In order to help the user with combining *PRISM* with MCMC to use hybrid sampling, the `utils` module provides two functions: `get_walkers()` and `get_inpost_fn()`. The `get_walkers()` function analyzes a set of proposed `init_walkers` and returns the positions that are valid (and the number of positions that are valid). By default, it uses the available `impl_sam` of the last constructed iteration, but it can also be supplied with a custom set of proposed walkers or an integer stating how many proposed walkers the function should check:

```
>>> # Use impl_sam if it is available
>>> n, p0 = get_walkers(pipe)

>>> # Request 2000 proposed samples
>>> n_walkers = 2000
>>> n, p0 = get_walkers(pipe, init_walkers=n_walkers)

>>> # Use custom init_walkers
>>> from numpy.random import rand
>>> init_walkers = rand(n_walkers, pipe.modellink.n_par)
>>> n, p0 = get_walkers(pipe, init_walkers=init_walkers)
```

Given that most sampling methods operate in unit space (all parameter values are between zero and unity), the `get_walkers()` function automatically assumes that all samples are defined in unit space (and therefore using `rand()` works). It can however be disabled by setting the `unit_space` input argument to `False`. One has to keep in mind that, because of the way the emulator works, there is no guarantee for a specific number of valid walkers to be obtained. Having the desired emulator iteration already analyzed may give an indication how many samples in total need to be proposed to be left with a specific number.

When the initial positions of the MCMC walkers have been determined, one can use them in an MCMC parameter estimation algorithm, avoiding the burn-in phase. This in itself can already be very useful, but it does not allow for hybrid sampling yet. Most MCMC methods require the definition of an `lnpost` function, which takes a parameter set and returns the corresponding natural logarithm of the posterior probability. In order to do hybrid sampling, this `lnpost` function must have the algorithm described above implemented.

The `get_lnpost_fn()` function factory provides exactly that. It takes a user-defined `ext_lnpost` function and a `Pipeline` object, and returns a function definition `get_lnpost(par_set, *args, **kwargs)`. This `get_lnpost` function first analyzes a proposed `par_set` in the emulator, passes `par_set` (along with any additional arguments) to `ext_lnpost` if the sample is plausible, or returns $-\infty$ if it is not. The return-value of the `ext_lnpost` function is then returned by the `get_lnpost` function as well. To make sure that the `get_lnpost` function can be used in both execution modes (see *Dual nature (normal/worker mode)*), all parallel calls to the `Pipeline` object are done with the `_make_call()` method.

The use of a function factory here allows for all input arguments to be validated once and then saved as local variables for the `get_lnpost` function. Not only does this avoid that all arguments have to be provided and validated for every individual call, but it also ensures that the same arguments are used every time, as local variables of a function cannot be modified by anything. Since users most likely use `get_walkers()` and `get_lnpost_fn()` frequently together, the `get_walkers()` function allows for the `ext_lnpost` argument to be supplied to it. This will automatically call the `get_lnpost_fn()` function factory using the provided `ext_lnpost` and the same input arguments given to `get_walkers()`, and return the obtained `get_lnpost` function in addition to the initial MCMC walkers.

5.4.3 Application

Using the information above, using hybrid sampling on a model of choice can be done quite easily. For performing the MCMC analysis, we will be using the `emcee` package in this example.

Assume that we want to first analyze and then optimize the Gaussian model given by the `GaussianLink` class. So, we first have to make an emulator of the model:

```
>>> from prism import Pipeline
>>> from prism.modellink import GaussianLink
>>> model_data = {3: [3.0, 0.1], 5: [5.0, 0.1], 7: [3.0, 0.1]}
>>> modellink_obj = GaussianLink(model_data=model_data)
>>> pipe = Pipeline(modellink_obj)
>>> pipe.construct()
```

Using the constructed emulator, we can perform a model parameter optimization using hybrid sampling. For this, we need to define an `ext_lnpost` function, for which we will use a simple Gaussian probability function:

```
def ext_lnpost(par_set, pipe):
    # Create parameter dict for call_model
    par_dict = dict(zip(pipe.modellink.par_name, par_set))

    # Use wrapped model to obtain model output
    mod_out = pipe.modellink.call_model(pipe.emulator.emul_i,
                                       par_dict,
                                       pipe.modellink.data_idx)
```

(continues on next page)

(continued from previous page)

```

# Get the model and data variances
# Since the value space is linear, the data error is centered
md_var = pipe.modellink.get_md_var(pipe.emulator.emul_i,
                                   par_dict,
                                   pipe.modellink.data_idx)

data_var = [err[0]**2 for err in pipe.modellink.data_err]

# Calculate the posterior probability and return it
sigma_2 = md_var+data_var
diff = pipe.modellink.data_val-mod_out
return (-0.5*(np.sum(diff**2/sigma2)) )

```

Since the *Pipeline* object already has the model wrapped and linked, we used that to evaluate the model. The *GaussianLink* class has a centered data error, therefore we can take the upper bound for every error when calculating the variance. However, for more complex models, this is probably not true.

Next, we have to obtain the starting positions for the MCMC walkers. Since we want to do hybrid sampling, we can obtain the *get_lnpost* function at the same time as well:

```

>>> from prism.utils import get_walkers
>>> n, p0, get_lnpost = get_walkers(pipe, unit_space=False,
                                   ext_lnpost=ext_lnpost, hybrid=True)

```

By setting *hybrid* to True, we use the implausibility cut-off value as an additional prior. Now we only still need the *EnsembleSampler* class and NumPy (for the *ext_lnpost* function):

```

>>> import numpy as np
>>> from emcee import EnsembleSampler

```

Now we have everything that is required to perform a hybrid sampling analysis. In most cases, MCMC methods require to be executed on only a single MPI rank, so we will use the *worker_mode*:

```

# Activate worker mode
with pipe.worker_mode:
    if pipe.is_controller:
        # Create EnsembleSampler object
        sampler = EnsembleSampler(n, pipe.modellink.n_par,
                                  get_lnpost, args=[pipe])

        # Run mcmc for 1000 iterations
        sampler.run_mcmc(p0, 1000)

        # Execute any custom operations here
        # For example, saving the chain data or plotting the results

```

And that is basically all that is required for using *PRISM* together with MCMC. For a normal MCMC approach, the same code can be used, except that one has to use *ext_lnpost* instead of *get_lnpost* (and, obtain the starting positions of the walkers in a different way).

5.5 General usage rules

Below is a list of general usage rules that apply to *PRISM*.

- Unless specified otherwise in the documentation, any input argument in the *PRISM* package that accepts...
 - a bool (True/False) also accepts 0/1 as a valid input;

- `None` indicates a default value or operation for obtaining this input argument. In most of these cases, the default value depends on the current state of the *PRISM* pipeline, and therefore a small operation is required for obtaining this value;

Example

Providing `None` to `pot_active_par`, where it indicates that all model parameters should be potentially active.

- the names of model parameters also accepts the internal indices of these model parameters. The index is the order in which the parameter names appear in the `par_name` list or as they appear in the output of the `details()` method;
- a sequence of integers, floats and/or strings will accept (almost) any formatting including most special characters as separators as long as they do not have any meaning (like a dot for floats or valid escape sequences for strings). Keep in mind that providing `'1e3'` (or equivalent) will be converted to `1000.0`, as per Python standards;

Example

The following sequences are equal:

```
* A, 1, 20.0, B;
* [A,1,2e1,B];
* "A 1 20. B";
* "' [' (A / }| \n; <1{}} \,,'">20.000000 !! \t< )?%\B \."
```

- the path to a data file (*PRISM* parameters, model parameters, model data) will read in all the data from that file as a Python dict, with a colon `:` acting as the separator between the key and value.
- Depending on the used emulator type, state of loaded emulator and the *PRISM* parameter values, it is possible that providing values for certain *PRISM* parameters has no influence on the outcome of the pipeline. This can be either because they have non-changeable default values or are simply not used anywhere (given the current state of the pipeline);

Examples

- If `method != 'gaussian'`, it causes `sigma` to have no use in the pipeline;
 - Switching the bool value for `use_mock` while loading a constructed emulator has no effect, since the mock data is generated (or not) when constructing a new emulator and cannot be changed or swapped out afterward.
-

- All docstrings in *PRISM* are written in **RTF** (Rich Text Format) and are therefore best viewed in an editor that supports it (like **Spyder**).

5.6 External data files

When using *PRISM*, there are three different cases where the path to an external data file can be provided. As mentioned in *General usage rules*, all external files are read-in as a Python dict, with the colon being the separator between the key and value. Additionally, all lines are read as strings and converted back when assigned in memory, to allow

for many different mark-ups to be used. Depending on which of the three files is read-in, the keys and values have different meanings. Here, the three different files are described.

5.6.1 PRISM parameters file

This file contains the non-default values that must be used for the *PRISM parameters*. These parameters control various different functionalities of *PRISM*. It is provided as the *prism_par* argument when initializing the *Pipeline* class and stored in the *prism_dict* property (a dict or array-like can be provided instead as well). When certain parameters are set depends on their type:

- Emulator parameters: Whenever a new emulator is created;
- Pipeline parameters: When the *Pipeline* class is initialized;
- Implausibility parameters: When the *analyze()* method is called (saved to HDF5) or when an emulator iteration is loaded that has not been analyzed yet (not saved to HDF5);
- Projection parameters: When the *project()* method is called.

The default *PRISM* parameters file can be found in the *prism/data* folder and is shown below:

```
n_sam_init      : 500                # Number of initial model evaluation_
↳samples
proj_res       : 25                  # Number of projected grid points per_
↳model parameter
proj_depth     : 250                 # Number of emulator evaluation samples_
↳per projected grid point
base_eval_sam  : 800                 # Base number for growth in number of_
↳model evaluation samples
sigma          : 0.8                 # Gaussian sigma/standard deviation (only_
↳required if method == 'gaussian')
l_corr         : 0.3                 # Gaussian correlation length(s)
impl_cut       : [0.0, 4.0, 3.8, 3.5] # List of implausibility cut-off values
criterion      : None                # Criterion for constructing LHDs
method         : 'full'              # Method used for constructing the_
↳emulator
use_regr_cov   : False               # Use regression covariance
poly_order     : 3                   # Polynomial order for regression
n_cross_val    : 5                   # Number of cross-validations for_
↳regression
do_active_anal : True                # Perform active parameter analysis
freeze_active_par : True              # Active parameters always stay active
pot_active_par : None                # List of potentially active parameters
use_mock       : False               # Use mock data
```

In this file, the key is the name of the parameter that needs to be changed, and the value what it needs to be changed to. *PRISM* itself does not require this default file, as all of the default values are hard-coded, and is therefore never read-in. An externally provided *PRISM* parameters file is only required to have the non-default values. The contents of this file is equal to providing the following as *prism_par*:

```
# As a dict
prism_par = {'n_sam_init': 500,
             'proj_res': 25,
             'proj_depth': 250,
             'base_eval_sam': 800,
             'sigma': 0.8,
             'l_corr': 0.3,
             'impl_cut': [0.0, 4.0, 3.8, 3.5],
```

(continues on next page)

(continued from previous page)

```

        'criterion': None,
        'method': 'full',
        'use_regr_cov': False,
        'poly_order': 3,
        'n_cross_val': 5,
        'do_active_anal': True,
        'freeze_active_par': True,
        'pot_active_par': None,
        'use_mock': False}

# As an array_like
prism_par = [['n_sam_init', 500],
             ['proj_res', 25],
             ['proj_depth', 250],
             ['base_eval_sam', 800],
             ['sigma', 0.8],
             ['l_corr', 0.3],
             ['impl_cut', [0.0, 4.0, 3.8, 3.5]],
             ['criterion', None],
             ['method', 'full'],
             ['use_regr_cov', False],
             ['poly_order', 3],
             ['n_cross_val', 5],
             ['do_active_anal', True],
             ['freeze_active_par', True],
             ['pot_active_par', None],
             ['use_mock', False]]

```

Note that it is also possible to set any parameter besides *Emulator* parameters by using the corresponding class property.

5.6.2 Model parameters file

This file contains the non-default model parameters to use for a model. It is provided as the *model_parameters* input argument when initializing the *ModelLink* subclass (a dict or array-like can be provided instead as well). Keep in mind that the *ModelLink* subclass may not have default model parameters defined.

An example of the various different ways model parameter information can be provided is given below:

#	name	: lower_bndry	upper_bndry	estimate
	A	: 1	5	3
	Bravo	: 2	7	None
	C42	: 3	6.74	

In this file, the key is the name of the model parameter and the value is a sequence of integers or floats, specifying the *lower* and *upper* boundaries of the parameter and, optionally, its estimate. Similarly to the *PRISM* parameters, one can provide the following equivalent as *model_parameters* during initialization of a *ModelLink* subclass:

```

# As a dict
model_parameters = {'A': [1, 5, 3],
                   'Bravo': [2, 7, None],
                   'C42': [3, 6.74]}

# As an array_like
model_parameters = [['A', [1, 5, 3]],

```

(continues on next page)

(continued from previous page)

```

        ['Bravo', [2, 7, None]],
        ['C42', [3, 6.74]]]

# As two array_likes zipped
model_parameters = zip(['A', 'Bravo', 'C42'],
                       [[1, 5, 3], [2, 7, None], [3, 6.74]])

```

Providing `None` as the parameter estimate or not providing it at all, implies that no parameter estimate (for the corresponding parameter) should be used in the projection figures. If required, one can use the `convert_parameters()` function to validate their parameters formatting before using it to initialize a `ModelLink` subclass.

5.6.3 Model data file

This file contains the non-default model comparison data points to use for a model. It is provided as the `model_data` input argument when initializing the `ModelLink` subclass (a dict or array-like can be provided instead as well). Keep in mind that the `ModelLink` subclass may not have default model comparison data defined.

An example of the various different ways model comparison data information can be provided is given below:

#	data_idx	:	data_val		data_err		data_spc
	1, 2	:	1		0.05	0.05	'lin'
	3.0	:	2		0.05		'log'
	['A']	:	3		0.05	0.15	
	1, A, 1.0	:	4		0.05		

Here, the key is the full sequence of the data identifier of a data point, where any character that is not a letter, number, minus/plus or period acts as a separator between the elements of the data identifier. The corresponding value specifies the data value, data error(s) and data value space. Braces, parentheses, brackets and many other characters can be used as mark-up in the data identifier, to make it easier for the user to find a suitable file lay-out. A full list of all characters that can be used for this can be found in `prism.aux_char_set` and can be freely edited.

Similarly to the model parameters, the following is equal to the contents of this file:

```

# As a dict
model_data = {(1, 2): [1, 0.05, 0.05, 'lin'],
              3.0: [2, 0.05, 'log'],
              ('A'): [3, 0.05, 0.15],
              (1, 'A', 1.0): [4, 0.05]}

# As an array_like
model_data = [(1, 2), [1, 0.05, 0.05, 'lin']],
              [3.0, [2, 0.05, 'log']],
              [('A'), [3, 0.05, 0.15]],
              [(1, 'A', 1.0), [4, 0.05]]]

# As two array_likes zipped
model_data = zip([(1, 2), 3.0, ('A'), (1, 'A', 1.0)],
                 [[1, 0.05, 0.05, 'lin'], [2, 0.05, 'log'], [3, 0.05, 0.15], [4, 0.
→05]])

```

It is necessary for the data value to be provided at all times. The data error can be given as either a single value, where it assumed that the data point has a centered 1σ -confidence interval, or as two values, where they describe the *upper* and *lower* bounds of the 1σ -confidence interval. The data value space can be given as a string or omitted, in which case it is assumed that the value space is linear. Keep in mind that, as mentioned in [Data identifiers \(data_idx\)](#), providing

a single element data identifier causes it to be saved as a scalar instead of a tuple. Therefore, ['A'] or ('A') is the same as 'A' . If required, one can use the `convert_data()` function to validate their data formatting before using it to initialize a `ModelLink` subclass.

Note: The parameter value bounds are given as [*lower bound*, *upper bound*], whereas the data errors are given as [*upper error*, *lower error*]. The reason for this is that, individually, the order for either makes the most sense. Together however, it may cause some confusion, so extra care needs to be taken.

6.1 Terminology

Below is a list of the most commonly used terms/abbreviations in *PRISM* and their meaning.

Active emulator system An emulator system that has a data point assigned to it.

Active parameters The set of model parameters that are considered to have significant influence on the output of the model and contribute at least one polynomial term to one/the regression function.

Adjusted expectation The prior expectation of a parameter set, with the adjustment term taken into account. It is equal to the prior expectation if the emulator system has perfect accuracy.

Adjusted values The adjusted expectation and variance values of a parameter set.

Adjusted variance The prior variance of a parameter set, with the adjustment term taken into account. It is zero if the emulator system has perfect accuracy.

Adjustment term The extra term (as determined by the BLA) that is added to the prior expectation and variance values that describes all additional correlation knowledge between model realization samples.

Analysis

Analyze The process of evaluating a set of emulator evaluation samples in the last emulator iteration and determining which samples should be used to construct the next iteration.

BLA Abbreviation of *Bayes linear approach*.

Construct

Construction The process of calculating all necessary components to describe an iteration of the emulator.

Construction check A list of keywords determining which components of which emulator systems are still required to finish the construction of a specified emulator iteration.

Controller

Controller rank An MPI process that controls the flow of operations in *PRISM* and distributes work to all workers and itself. By default, a controller also behaves like a worker, although is not identified as such.

Covariance matrix

Inverted covariance matrix The (inverted) matrix of prior covariances between all model realization samples and itself.

Covariance vector The vector of prior covariances between all model realization samples and a given parameter set.

Data error The 1σ -confidence interval of a model comparison data point, often a measured/calculated observational error.

Data identifier

Data point identifier The unique identifier of a model comparison data point, often a sequence of integers, floats and strings that describe the operations required to extract it.

Data point A collection of all the details (value, error, space and identifier) about a specific model comparison data point that is used to constrain the model with.

Data space

Data value space The value space (linear, logarithmic or exponential) in which a model comparison data point is defined.

Data value The value of a model comparison data point, often an observed/measured value.

Emulation method The specific method (Gaussian, regression or both) that needs to be used to construct an emulator.

Emulator The collection of all emulator systems together, provided by an *Emulator* object.

Emulator evaluation samples The sample set (to be) used for evaluating the emulator.

Emulator iteration

Iteration A single, specified step in the construction of the emulator.

Emulator system The emulated version of a single model output/comparison data point in a single iteration.

Emulator type The type of emulator that needs to be constructed. This is used to make sure different emulator types are not mixed together by accident.

Evaluate

Evaluation The process of calculating the adjusted values of a parameter set in all emulator systems starting at the first iteration, determining the corresponding implausibility values and performing an implausibility check. This process is repeated in the next iteration if the check was successful and the requested iteration has not been reached.

External model realization set A set of externally calculated and provided model realization samples and their outputs.

Frozen parameters

Frozen active parameters The set of model parameters that, once considered active, will always stay active if possible.

FSLR Abbreviation of *forward stepwise linear regression*.

Gaussian correlation length The maximum distance between two values of a specific model parameter within which the Gaussian contribution to the correlation between the values is still significant.

Gaussian sigma The standard deviation of the Gaussian function. It is not required if regression is used.

HDF5 Abbreviation of *Hierarchical Data Format version 5*.

Hybrid sampling The process of doing a best parameter estimation of a model with MCMC sampling, while using its emulator as an additional Bayesian prior. This process is explained in [Hybrid sampling](#).

Implausibility check

Implausibility cut-off check The process of determining whether or not a given set of implausibility values satisfy the implausibility cut-offs of a specific emulator iteration.

Implausibility cut-offs The maximum implausibility values an evaluated parameter set is allowed to generate, to be considered plausible in a specific emulator iteration.

Implausibility value

Univariate implausibility value The number of sigmas an emulator system expects the (real) model output corresponding to a given parameter set, to be away from the data point it is compared against, given its adjusted values. It takes into account all variances associated with the parameter set, which are the observational variance (given by *data_err*), adjusted emulator variance (*adj_var*) and the model discrepancy variance (*md_var*).

Implausibility wildcard A maximum implausibility value, preceding the implausibility cut-offs, that is not taken into account during the implausibility cut-off check. It is denoted as 0 in provided implausibility cut-off lists.

LHD Abbreviation of *Latin-Hypercube design*.

Master file

Master HDF5 file (Path to) The HDF5-file in which all important data about the currently loaded emulator is stored. A master file is usually accompanied by several emulator system (HDF5) files, which store emulator system specific data and are externally linked to the master file.

MCMC Abbreviation of *Markov chain Monte Carlo*.

Mock data The set of comparison data points that has been generated by evaluating the model for a random parameter set and perturbing the output by the model discrepancy variance.

Model A *black box* that takes a parameter set, performs a sequence of operations and returns a unique collection of values corresponding to the provided parameter set.

Note: This is how *PRISM* ‘sees’ a model, not the used definition of one.

2D model A model that has/takes 2 model parameters.

2+D model

nD model A model that has/takes more than 2 model parameters.

ModelLink

ModelLink subclass The user-provided wrapper around the model that needs to be emulated, provided by a *ModelLink* object.

Model discrepancy variance A user-defined value that includes all contributions to the overall variance on a model output that is created/caused by the model itself. More information on this can be found in [Model discrepancy variance](#) (*md_var*).

Model evaluation samples The sample set (to be) used for evaluating the model.

Model output

Model outputs The model output(s) corresponding to a single (set of) model realization/evaluation sample(s).

Model realization samples Same as model evaluation samples.

Model realizations

Model realization set The combination of model realization/evaluation samples and their corresponding model outputs.

MPI Abbreviation of *Message Passing Interface*.

MPI rank An MPI process that is used by any *PRISM* operation, either being a controller or a worker.

MSE Abbreviation of *mean squared error*.

OLS Abbreviation of *ordinary least-squares*.

Parameter set

Sample A single combination/set of model parameter values, used to evaluate the emulator/model once.

Passive parameters The set of model parameters that are not considered active, and therefore are considered to not have a significant influence on the output of the model.

Pipeline

PRISM Pipeline The main *PRISM* framework that orchestrates all operations, provided by a *Pipeline* object.

Plausible region The region of model parameter space that still contains plausible samples.

Plausible samples A subset of a set of emulator evaluation samples that satisfied the implausibility checks.

Polynomial order Up to which order polynomial terms need to be taken into account for all regression processes.

Potentially active parameters A user-provided set of model parameters that are allowed to become active. Any model parameter that is not potentially active will never become active, even if it should.

PRISM The acronym for *Probabilistic Regression Instrument for Simulating Models*. It is also a one-word description of what *PRISM* does (splitting up a model into individually emulated model outputs).

Prior covariance The covariance value between two parameter sets as determined by an emulator system.

Prior expectation The expectation value of a parameter set as determined by an emulator system, without taking the adjustment term (from the BLA) into account. It is a measure of how much information is captured by an emulator system. It is zero if regression is not used, as no information is captured.

Prior variance The variance value of a parameter set as determined by an emulator system, without taking the adjustment term (from the BLA) into account.

Project

Projection The process of analyzing a specific set of active parameters in an iteration to determine the correlation between the parameters.

Projection figure The visual representation of a projection.

Regression The process of determining the important polynomial terms of the active parameters and their coefficients, by using an FSLR algorithm.

Regression covariances The covariances between all polynomial coefficients of the regression function. By default, they are not calculated and it is empty if regression is not used.

Residual variance The variance that has not been captured during the regression process. It is empty if regression is not used.

Root directory (Path to) The directory/folder on the current machine in which all *PRISM* working directories are located. It also acts as the base for all relative paths.

Sample set

Evaluation set A set of samples.

Worker

Worker rank An MPI process that receives its calls/orders from a controller and performs the heavy-duty operations in *PRISM*.

Working directory (Path to) The directory/folder on the current machine in which the *PRISM* master file, log-file and all projection figures of the currently loaded emulator are stored.

Worker mode A mode initialized by *worker_mode*, where all workers are continuously listening for calls made by the controller rank and execute the received messages. This allows for serial codes to be combined more easily with *PRISM*. See *Dual nature (normal/worker mode)* for more information.

6.2 PRISM parameters

Below are descriptions of all the parameters that can be provided to *PRISM* in a text-file when initializing the *Pipeline* class (using the *prism_par* input argument).

Changed in version 1.1.2: Input argument *prism_file* was renamed to *prism_par*. A dictionary with *PRISM* parameters instead of a file can additionally be provided to the *Pipeline* class. All *Pipeline* parameters can also be changed by setting the corresponding class property.

n_sam_init (Default: 500) Number of model evaluation samples that is used to construct the first iteration of the emulator. This value must be a positive integer.

proj_res (Default: 25) Number of emulator evaluation samples that is used to generate the grid for the projection figures (it defines the resolution of the projection). This value must be a positive integer.

proj_depth (Default: 250) Number of emulator evaluation samples that is used to generate the samples in every projected grid point (it defines the accuracy/depth of the projection). This value must be a positive integer.

base_eval_sam (Default: 800) Base number of emulator evaluation samples that is used to analyze an iteration of the emulator. It is multiplied by the iteration number and the number of model parameters to generate the true number of emulator evaluations, in order to ensure an increase in emulator accuracy. This value must be a positive integer.

sigma (Default: 0.8) The Gaussian sigma/standard deviation that is used when determining the Gaussian contribution to the overall emulator variance. This value is only required when *method* == 'gaussian', as the Gaussian sigma is obtained from the residual variance left after the regression optimization if regression is included. This value must be non-zero.

l_corr (Default: 0.3) The normalized amplitude(s) of the Gaussian correlation length. This number is multiplied by the difference between the upper and lower value boundaries of the model parameters to obtain the Gaussian correlation length for every model parameter. This value must be positive, normalized and either a scalar or a list of *n_par* scalars (where the values correspond to the sorted list of model parameters).

impl_cut (Default: [0.0, 4.0, 3.8, 3.5]) A list of implausibility cut-off values that specifies the maximum implausibility values a parameter set is allowed to have to be considered 'plausible'. A zero can be used as a filler value, either taking on the preceding value or acting as a wildcard if the preceding value is a wildcard or non-existent. Zeros are appended at the end of the list if the length is less than the number of comparison data points, while extra values are ignored if the length is more. This must be a sorted list of positive values (excluding zeros).

criterion (Default: None) The criterion to use for determining the quality of the LHDs that are used, represented by an integer, float, string or None. This parameter is the only non-*PRISM* parameter. Instead, it is used in the *lhd()*-function of the *e13Tools* package. By default, None is used.

method (Default: 'full') The method to use for constructing the emulator. 'gaussian' will only include Gaussian processes (no regression), which is much faster, but also less accurate. 'regression' will only include regression processes (no Gaussian), which is more accurate than Gaussian only, but underestimates the emulator

variance by multiple orders of magnitude. 'full' includes both Gaussian and regression processes, which is slower than Gaussian only, but by far the most accurate both in terms of expectation and variance values.

'gaussian' can be used for faster exploration especially for simple models. 'regression' should only be used when the polynomial representation of a model is important and enough model realizations are available. 'full' should be used by default.

Warning: When using *PRISM* on a model that can be described completely by the regression function (anything that has an analytical, polynomial form up to order `poly_order` like a straight line or a quadratic function), use the 'gaussian' method unless unavoidable (in which case `n_sam_init` and `base_eval_sam` must be set to very low values).

When using the regression method on such a model, *PRISM* will be able to capture the behavior of the model perfectly given enough samples, in which case the residual (unexplained) variance will be approximately zero and therefore `sigma` as well. This can occasionally cause floating point errors when calculating emulator variances, which in turn can give unexplainable artifacts in the adjustment terms, therefore causing answers to be incorrect.

Since *PRISM*'s purpose is to identify the characteristics of a model and therefore it does not know anything about its workings, it is not possible to automatically detect such problems.

`use_regr_cov` (Default: False) Whether or not the regression variance should be taken into account for the variance calculations. The regression variance is the variance on the regression process itself and is only significant if a low number of model realizations (`n_sam_init` and `base_eval_sam`) is used to construct the emulator systems. Including it usually only has a very small effect on the overall variance value, while it can slow down the emulator evaluation rate by as much as a factor of 3. This value is not required if `method == 'gaussian'` and is automatically set to True if `method == 'regression'`. This value must be a bool.

`poly_order` (Default: 3) Up to which order all polynomial terms of all model parameters should be included in the active parameters and regression processes. This value is not required if `method == 'gaussian'` and `do_active_anal` is False. This value must be a positive integer.

`n_cross_val` (Default: 5) Number of (k-fold) cross-validations that must be used for determining the quality of the active parameters analysis and regression process fits. If this parameter is zero, cross-validations are not used. This value is not required if `method == 'gaussian'` and `do_active_anal` is False. This value must be a non-negative integer and not equal to 1.

`do_active_anal` (Default: True) Whether or not an active parameters analysis must be carried out for every iteration of every emulator system. If False, all potentially active parameters listed in `pot_active_par` will be active. This value must be a bool.

`freeze_active_par` (Default: True) Whether or not active parameters should be frozen in their active state. If True, parameters that have been considered active in a previous iteration of an emulator system, will automatically be active again (and skip any active parameters analysis). This value must be a bool.

`pot_active_par` (Default: None) A list of parameter names that indicate which parameters are potentially active. Potentially active parameters are the only parameters that will enter the active parameters analysis (or will all be active if `do_active_anal` is False). Therefore, all parameters not listed will never be considered active. If all parameters should be potentially active, then a None can be given. This must either be a list of parameter names or None.

`use_mock` (Default: False) Whether or not mock data must be used as comparison data when constructing a new emulator. Mock data is calculated by evaluating the model for a specific set of parameter values, and adding the model discrepancy variances as noise to the returned data values. This set of parameter values is either the provided set, or a randomly chosen one if not. When using mock data for an emulator, it is not possible to change the comparison data in later emulator iterations. This value must be a bool or a list of `n_par` scalars (where the values correspond to the sorted list of model parameters).

6.3 HDF5

Whenever *PRISM* constructs an emulator, it automatically stores all the calculated data for it in an **HDF5-file** named 'prism.hdf5' in the designated working directory. This file contains all the data that is required in order to recreate all emulator systems that have been constructed for the emulator belonging to this run. If the *Pipeline* class is initialized by using an HDF5-file made by *PRISM*, it will load in this data and return a *Pipeline* object in the same state as described in the file.

Below is a short overview of all the data that can be found inside a *PRISM* master HDF5-file. HDF5-files can be viewed freely by the user using the **HDFView** application made available by **The HDFGroup**.

The general file contains:

- Attributes (9/10): Describe the general non-changeable properties of the emulator, which include:
 - Emulator type and method;
 - Gaussian parameters;
 - Name of used *ModelLink* subclass;
 - Used *PRISM* version;
 - Polynomial order;
 - Booleans for using mock data or regression covariance;
 - Mock data parameters if mock data was used.
 - Every emulator iteration has its own data group with the iteration number as its name. This data group stores all data/information specific to that iteration.
-

An iteration data group ('i') contains:

- Attributes (9): Describe the general properties and results of this iteration, including:
 - Active parameters for this emulator iteration;
 - Implausibility cut-off parameters;
 - Number of emulated data points, emulator systems, emulator evaluation samples, plausible samples and model realization samples;
 - Bool stating whether this emulator iteration used an external model realization set.
- 'emul_n': The data group that contains all data for a specific emulator system in this iteration. The value of 'n' indicates which emulator system it is, not the data point. See below for its contents;
- 'impl_sam': The set of emulator evaluation samples that survived the implausibility checks and will be used to construct the next iteration;
- 'proj_hcube': The data group that contains all data for the (created) projections for this iteration, if at least one has been made. See below for its contents;
- 'sam_set': The set of model realization samples that were used to construct this iteration. In every iteration after the first, this is the 'impl_sam' of the previous iteration;
- 'statistics': An empty data set that stores several different types of statistics as its attributes, including:
 - Size of the MPI communicator during various construction steps;

- Average evaluation rate/time of the emulator and model;
 - Total time cost of most construction steps (note that this value may be incorrect if a construction was interrupted);
 - Percentage of parameter space that is still plausible within the iteration.
-

An emulator system data group ('i/emul_n') contains:

- Attributes (5+): List the details about the model comparison data point used in this emulator system, including:
 - Active parameters for this emulator system;
 - Data errors, identifiers, value space and value;
 - Regression score and residual variance if regression was used.
 - 'cov_mat': The pre-calculated covariance matrix of all model evaluation samples in this emulator system. This data set is never used in *PRISM* and stored solely for user-convenience;
 - 'cov_mat_inv': The pre-calculated inverse of 'cov_mat';
 - 'exp_dot_term': The pre-calculated second expectation adjustment dot-term $(\text{Var}(D)^{-1} \cdot (D - E(D)))$ of all model evaluation samples in this emulator system.
 - 'mod_set': The model outputs for the data point in this emulator system corresponding to the 'sam_set' used in this iteration;
 - 'poly_coef' (if regression is used): The non-zero coefficients for the polynomial terms in the regression function in this emulator system;
 - 'poly_coef_cov' (if regression and regr_cov are used): The covariances for all polynomial coefficients 'poly_coef';
 - 'poly_idx' (if regression is used): The indices of the polynomial terms with non-zero coefficients if all active parameters are converted to polynomial terms;
 - 'poly_powers' (if regression is used): The powers of the polynomial terms corresponding to 'poly_idx'. Both 'poly_idx' and 'poly_powers' are required since different methods of calculating the polynomial terms are used depending on the number of required terms and samples;
 - 'prior_exp_sam_set': The pre-calculated prior expectation values of all model evaluation samples in this emulator system. This data set is also never used in *PRISM*.
-

A projections data group ('i/proj_hcube') contains individual projection data groups ('i/proj_hcube/<name>'), which

- Attributes (4): List the general properties with which this projection was made, including:
 - Implausibility cut-off parameters (they can differ from the iteration itself);
 - Projection depth and resolution.
- 'impl_los': The calculated line-of-sight depth for all grid points in this projection;
- 'impl_min': The calculated minimum implausibility values for all grid points in this projection.

7.1 How do I contribute?

Contributing to *PRISM* is done through pull requests in the [repository](#). If you have an idea on what to contribute, it is recommended to open a [GitHub issue](#) about it, such that the maintainers can give their advice or help. If you want to contribute but do not really know what, then you can take a look at the Todos that are scattered throughout the code. When making a contribution, please keep in mind that it must be compatible with all Python versions that *PRISM* supports (3.5+), and preferably with all operating systems as well.

7.2 How do I report a bug/problem?

By opening a [GitHub issue](#) and using the *Bug report* template.

7.3 What does the term ... mean?

A list of the most commonly used terms in *PRISM* can be found on the [Terminology](#) page.

7.4 What OS are supported?

PRISM should be compatible with all Windows, Mac OS and UNIX-based machines, as long as they support one of the required Python versions. Compatibility is currently tested for Linux 16.04 (Xenial)/Mac OS-X using [Travis CI](#) and Windows 32-bit/64-bit using [AppVeyor](#).

Community guidelines

PRISM is an open-source and free-to-use software package (and it always will be), provided under the BSD-3 license (see below for the full license).

Users are highly encouraged to make contributions to the package or request new features by opening a [GitHub issue](#). If you would like to contribute to the package, but do not know what, then there are quite a few Todos in the code that may give you some inspiration. As with contributions, if you find a problem or issue with *PRISM*, please do not hesitate to open a [GitHub issue](#) about it.

And, finally, if you use *PRISM* as part of your workflow in a scientific publication, please consider including an acknowledgement like “*Parts of the results in this work were derived using the PRISM Python package.*” and citing the *PRISM* pipeline paper:

```
@ARTICLE{2019ApJS..242...22V,
  author = {{van der Velden}, E. and {Duffy}, A.~R. and {Croton}, D. and
    {Mutch}, S.~J. and {Sinha}, M.},
  title = "{Model dispersion with PRISM; an alternative to MCMC for rapid analysis_
↪of models}",
  journal = {\apjs},
  keywords = {Astrophysics - Instrumentation and Methods for Astrophysics, Physics -
↪ Computational Physics},
  year = "2019",
  month = "Jun",
  volume = {242},
  number = {2},
  eid = {22},
  pages = {22},
  doi = {10.3847/1538-4365/ab1f7d},
  archivePrefix = {arXiv},
  eprint = {1901.08725},
  primaryClass = {astro-ph.IM},
  adsurl = {http://adsabs.harvard.edu/abs/2019ApJS..242...22V},
  adsnote = {Provided by the SAO/NASA Astrophysics Data System}
}
```

8.1 License

BSD 3-Clause License

Copyright (c) 2019, Ellert van der Velden
All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- * Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- * Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- * Neither the name of the copyright holder nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

8.2 Additions

Below are some bigger ideas/improvements that may be added to *PRISM* if there is demand:

- Add a developer's guide to the docs, describing the inner workings and structures of *PRISM*;
- Low-level MPI implementation (probably by using D2O);

With 6 emulator systems and 4 processes, the three different MPI levels would be:

- No level: 6-0-0-0;
 - High-level: 2-2-1-1;
 - Low-level: 1.5-1.5-1.5-1.5.
- Dynamic implausibility cut-offs;
 - Allow for a master projection figure to be made (kind of like a double corner plot);
 - Allow for user-provided methods in the *ModelLink* subclass to be executed at specific points in the emulator construction;
 - Implement multi-variate implausibilities;
 - Allow for no *ModelLink* object to be provided, which blocks construction but enables everything emulator-only related;

- Allow for old *PRISM* master files to be provided when making a new emulator, recycling work done previously;
- If `MPI_call` is `False` for the `ModelLink` subclass, use all MPI ranks to evaluate a part of `sam_set` in the model simultaneously. This will require a check or flag that the model can be called in multiple instances simultaneously (to accommodate for models that, for example, need to read files during evaluations). Added benefit of this is that it would become possible to add the option for the user to set a preferred number of MPI processes calling the model (in MPI), allowing *PRISM* to split up the available processes if more efficient;
- GPU acceleration;
- Adding the theory behind *PRISM* to the docs;
- Adding the possibility to evaluate the derivatives of the emulated model outputs, which could be used as approximations of the gradient field of a model for certain MCMC methods;
- Replace the *list of lists* data system with a *list of dicts* system. This would remove the need for converting global indices to/from local indices in several cases, and make it easier for users to understand. However, as indexing dicts is more complicated, this may require a lot of rewriting;
- Code objects can be made pickleable by importing the `codeutil` module. This package could be added to the requirements or an equivalent function could be written, which is then automatically imported/executed upon importing *PRISM*;


```
class prism.Pipeline(modellink_obj, *, root_dir=None, working_dir=None, prefix=None,
                    prism_par=None, emul_type=None, comm=None, **kwargs)
Defines the Pipeline class of the PRISM package.
```

The *Pipeline* class is the main user class of the *PRISM* package and provides a user-friendly environment that gives access to all operations within the package.

```
_Projection__analyze_proj_hcube(hcube)
Analyzes an emulator projection hypercube hcube.
```

Parameters *hcube* (1D array_like of int of length {1, 2}) – Array containing the internal integer identifiers of the main model parameters that require a projection hypercube.

Returns

- **impl_min_hcube** (1D ndarray object) – List containing the lowest implausibility value that can be reached in every single grid point on the given hypercube.
- **impl_los_hcube** (1D ndarray object) – List containing the fraction of the total amount of evaluated samples in every single grid point on the given hypercube, that still satisfied the implausibility cut-off criterion.

```
_Projection__draw_2D_proj_fig(hcube, impl_min, impl_los, proj_res)
Draws the 2D projection figure for the provided hcube, given the impl_min and impl_los values.
```

Parameters

- **hcube** (1D array_like of int of length 1) – Array containing the internal integer identifiers of the main model parameters that require a projection figure.
- **impl_min** (1D ndarray object) – List containing the lowest implausibility value that can be reached in every single grid point on the given hypercube.
- **impl_los** (1D ndarray object) – List containing the fraction of the total amount of evaluated samples in every single grid point on the given hypercube, that still satisfied the implausibility cut-off criterion.

- **proj_res** (*int*) – Number of emulator evaluations used to generate the grid for the given hypercube.

`_Projection__draw_3D_proj_fig` (*hcube*, *impl_min*, *impl_los*, *proj_res*)

Draws the 3D projection figure for the provided *hcube*, given the *impl_min* and *impl_los* values.

Parameters

- **hcube** (*1D array_like of int of length 2*) – Array containing the internal integer identifiers of the main model parameters that require a projection figure.
- **impl_min** (*1D ndarray object*) – List containing the lowest implausibility value that can be reached in every single grid point on the given hypercube.
- **impl_los** (*1D ndarray object*) – List containing the fraction of the total amount of evaluated samples in every single grid point on the given hypercube, that still satisfied the implausibility cut-off criterion.
- **proj_res** (*int*) – Number of emulator evaluations used to generate the grid for the given hypercube.

`_Projection__get_default_input_arguments` ()

Generates a dict containing default values for all input arguments.

Returns **kwargs_dict** (*dict*) – Dict containing all default input argument values.

`_Projection__get_default_parameters` ()

Generates a dict containing default values for all projection parameters.

Returns **par_dict** (*dict*) – Dict containing all default projection parameter values.

`_Projection__get_fig_path` (*hcube*, *emul_i=None*)

Determines the absolute path of a projection figure corresponding to a provided projection hypercube *hcube* in emulator iteration *emul_i* and returns it.

Parameters **hcube** (*1D array_like of int of length {1, 2} or str*) – Array containing the internal integer identifiers of the main model parameters that require a projection hypercube. If str, the name of *hcube* instead (`_Projection__get_hcube_name()`).

Other Parameters **emul_i** (*int or None. Default: None*) – If int, number indicating the requested emulator iteration. If *None*, the last iteration of the emulator will be used.

Returns

- **fig_path** (*str*) – The absolute path to the requested projection figure.
- **fig_path_s** (*str*) – The absolute path to the smoothed version.

`_Projection__get_hcube_name` (*hcube*)

Determines the name of a provided projection hypercube *hcube* and returns it.

Parameters **hcube** (*1D array_like of int of length {1, 2}*) – Array containing the internal integer identifiers of the main model parameters that require a projection hypercube.

Returns **hcube_name** (*str*) – The name of this projection hypercube.

`_Projection__get_proj_data` (*hcube*)

Returns the projection data belonging to the provided hypercube *hcube*.

Parameters **hcube** (*1D array_like of int of length {1, 2}*) – Array containing the internal integer identifiers of the main model parameters that require a projection hypercube.

Returns

- **impl_min_hcube** (1D `ndarray` object) – List containing the lowest implausibility value that can be reached in every single grid point on the given hypercube.
- **impl_los_hcube** (1D `ndarray` object) – List containing the fraction of the total amount of evaluated samples in every single grid point on the given hypercube, that still satisfied the implausibility cut-off criterion.
- **proj_res** (*int*) – Number of emulator evaluations used to generate the grid for the given hypercube.
- **proj_depth** (*int*) – Number of emulator evaluations used to generate the samples in every grid point for the given hypercube.

`_Projection__get_proj_hcube` (*hcube*)

Generates a projection hypercube *hcube* containing emulator evaluation samples. The output of this function depends on the requested projection type.

Parameters **hcube** (*1D array_like of int of length {1, 2}*) – Array containing the internal integer identifiers of the main model parameters that require a projection hypercube.

Returns **proj_hcube** (3D `ndarray` object) – 3D projection hypercube of emulator evaluation samples. For 3D projections, the grid uses matrix indexing (second parameter varies the fastest).

`_Projection__get_req_hcubes` (*proj_par*)

Determines which projection hypercubes have been requested by the user. Also checks if these projection hypercubes have been calculated before, and depending on the value of `force`, either skips them or recreates them.

Parameters **proj_par** (*1D array_like of {int, str} or None*) – For which model parameters to construct the projection figures. If 1D array_like, construct projection figures for all combinations of provided model parameters that are active, with a string referring to the name of the model parameter and an integer referring to the position in which the model parameter is shown in the `details()` method. If *None*, projection figures are made for all active model parameters.

Generates

hcubes [list of lists] List containing the parameter indices of the requested projection hypercubes.

create_hcubes [list of lists] List containing the parameter indices of the requested projection hypercubes that need to be created first.

`_Projection__prepare_projections` (*emul_i, proj_par, **kwargs*)

Prepares the pipeline for the creation of the requested projections.

Parameters

- **emul_i** (*int or None*) – Number indicating the requested emulator iteration.
- **proj_par** (*1D array_like of {int, str} or None*) – For which model parameters to construct the projection figures. If 1D array_like, construct projection figures for all combinations of provided model parameters that are active, with a string referring to the name of the model parameter and an integer referring to the position in which the model parameter is shown in the `details()` method. If *None*, projection figures are made for all active model parameters.
- **kwargs** (*dict*) – Dict of keyword arguments that were provided to `project()`.

`__Projection__process_input_arguments` (*emul_i*, ***kwargs*)

Processes the input arguments given to the `project()` method.

Parameters

- **`emul_i`** (*int or None*) – Number indicating the requested emulator iteration.
- **`kwargs`** (*dict*) – Dict of keyword arguments that were provided to `project()`.

Generates

All default and provided arguments are saved to their respective properties.

`__Projection__save_data` (*data_dict*)

Saves a given data dict {`keyword`: `data`} at the emulator iteration this class was initialized for, to the HDF5-file.

Parameters **`data_dict`** (*dict*) – Dict containing the data that needs to be saved to the HDF5-file.

Keyword Arguments

- **`keyword`** (*{'nD_proj_hcube'}*) – String specifying the type of data that needs to be saved.
- **`data`** (*{int, float, str, array_like} or dict*) – The actual data that needs to be saved at data keyword *keyword*. If dict, save every item individually.

Generates

The specified data is saved to the HDF5-file.

`__Projection__set_parameters` ()

Sets the Projection parameters from the `prism_dict` property and saves them in the current Projection instance.

`__init__` (*modellink_obj*, ***, *root_dir=None*, *working_dir=None*, *prefix=None*, *prism_par=None*, *emul_type=None*, *comm=None*, ***kwargs*)

Initialize an instance of the `Pipeline` class.

Parameters **`modellink_obj`** (*ModelLink* object) – Instance of the `ModelLink` subclass that links the emulated model to this `Pipeline` instance.

Other Parameters

- **`root_dir`** (*str or None. Default: None*) – String containing the absolute path of the root directory where all working directories are stored. If *None*, root directory will be set to the directory this class was initialized in.
- **`working_dir`** (*str, bool or None. Default: None*) – String containing the name of the working directory of the emulator in *root_dir*. If *True*, a new working directory will be created in *root_dir*. If *None* or *False*, working directory is set to the last one that was created in *root_dir* that starts with the given *prefix*. Note that providing an integer instead of a bool will not work here. If no directories are found, one will be created.
- **`prefix`** (*str or None. Default: None*) – String containing a prefix that is used for naming new working directories or scan for existing ones. If *None*, all directories in *root_dir* are considered working directories and 'prism_' will be used as a prefix for new ones.
- **`prism_par`** (*array_like, dict, str or None. Default: None*) – A dict containing the values for the *PRISM* parameters that need to be changed from their default values. If *array_like*, `dict(prism_par)` must generate a dict with the correct lay-out. If *str*, the string is the absolute or relative path to the file that contains the keys in the first column and the dict values

in the second column. If a relative path is given, its path must be relative to *root_dir* or the current directory. If *None*, no changes will be made to the default parameters.

- **emul_type** (*Emulator* subclass or *None*. Default: *None*) – The type of *Emulator* to use in this *Pipeline* instance. If *None*, use the default emulator instead.
- **comm** (*Intracomm* object or *None*. Default: *None*) – The MPI intra-communicator to use in this *Pipeline* instance or `MPI.COMM_WORLD` if *comm* is *None*.

_call_model (*emul_i*, *par_set*, *data_idx*)

Obtain the output corresponding to the provided *data_idx* that is generated by the model for a given model parameter sample *par_set*. The current emulator iteration *emul_i* is also provided in case it is required by the *ModelLink* subclass.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **par_set** (*1D array_like*) – Model parameter sample to calculate the model output for.
- **data_idx** (*list of tuples*) – The list of data identifiers for which the model is requested to return the corresponding data values.

Returns **mod_out** (*1D ndarray* object) – Model output corresponding to given *par_set*.

_compile_code_snippets ()

Compiles all pre-defined built-in code snippets to code objects and saves them to *code_objects*. These code objects are used for performing standard operations in the *_evaluate_sam_set* () method.

_do_impl_check (*emul_i*, *uni_impl_val*)

Performs an implausibility cut-off check on the provided implausibility values *uni_impl_val* at emulator iteration *emul_i*.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **uni_impl_val** (*1D array_like*) – Array containing all univariate implausibility values corresponding to a certain parameter set for all data points.

Returns

- **result** (*bool*) – 1 if check was successful, 0 if it was not.
- **impl_cut_val** (*float*) – Implausibility value at the first real implausibility cut-off.

_evaluate_model (*emul_i*, *sam_set*, *data_idx*)

Evaluates the model for provided evaluation sample set *sam_set* at given data points *data_idx*.

This method automatically distributes the samples according to the various flags set in the *ModelLink* subclass.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **sam_set** (*1D or 2D array_like*) – Parameter/sample set to evaluate in the model.
- **data_idx** (*list of tuples*) – The list of data identifiers for which the model is requested to return the corresponding data values.

Returns **mod_set** (*2D ndarray* object of shape (*n_sam*, *n_data*)) – Array containing the data values corresponding to the requested data points.

`_evaluate_sam_set` (*emul_i*, *sam_set*, *exec_code*)

Evaluates a provided set of emulator evaluation samples *sam_set* at a given emulator iteration *emul_i*. The provided tuple of code snippets *exec_code* are executed using Python's `exec()` function at specific points during the analysis.

Parameters

- **`emul_i`** (*int*) – Number indicating the requested emulator iteration.
- **`sam_set`** (2D `ndarray` object) – Array containing model parameter value sets to be evaluated in all emulator systems in emulator iteration *emul_i*.
- **`exec_code`** (`{'analyze', 'evaluate', 'hybrid', 'project'}` or *tuple*) – Tuple of five code snippets (*pre_code*, *eval_code*, *anal_code*, *post_code*, *exit_code*) to be executed at specific points during the analysis. If string, use one of the built-in tuples in `code_objects` instead.

Other Parameters

- **`pre_code`** (*str or code object*) – Code snippet to be executed before the evaluation of *sam_set* starts.
- **`eval_code`** (*str or code object*) – Code snippet to be executed after the evaluation of each sample in *sam_set*.
- **`anal_code`** (*str or code object*) – Code snippet to be executed after the analysis of each sample in *sam_set*. This code snippet is only executed by the controller.
- **`post_code`** (*str or code object*) – Code snippet to be executed after the evaluation of *sam_set* ends.
- **`exit_code`** (*str or code object*) – Code snippet to be executed before returning the results of the evaluation of *sam_set*. This code snippet is only executed by the controller.

Returns *results* (*object*) – The object that is assigned to `results`, which is defaulted to `None` if no code snippet changes it. Preferably, the execution of *post_code* and/or *exit_code* modifies *results*. All MPI ranks return it.

Notes

If any of the code snippets is provided as a string, it will be compiled into a code object before starting the evaluation.

`_get_default_parameters` ()

Generates a dict containing default values for all pipeline parameters.

Returns *par_dict* (*dict*) – Dict containing all default pipeline parameter values.

`_get_eval_sam_set` (*emul_i*)

Generates an emulator evaluation sample set to be used for analyzing an emulator iteration. Currently uses the `lhd()` function.

Parameters *emul_i* (*int*) – Number indicating the requested emulator iteration.

Returns *eval_sam_set* (2D `ndarray` object) – Array containing the evaluation samples.

`_get_ext_real_set` (*ext_real_set*)

Processes an externally provided model realization set *ext_real_set*, containing the used sample set and the corresponding data value set.

Parameters `ext_real_set` (*list, dict or None*) – List of arrays containing an externally calculated set of model evaluation samples and its data values, a dict with keys ['sam_set ', 'mod_set '] containing these arrays or *None* if no external set needs to be used.

Returns

- **ext_sam_set** (1D or 2D `ndarray` object) – Array containing the externally provided model evaluation samples.
- **ext_mod_set** (1D or 2D `ndarray` object) – Array containing the model outputs of all specified externally provided model evaluation samples.

_get_iteration_data (*emul_i, sam_set, ext_sam_set, ext_mod_set*)

Obtains the model realization data for given emulator iteration *emul_i* by evaluating the provided *sam_set* in the model and distributing model outputs to the correct emulator systems.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **sam_set** (2D `ndarray` object) – Array containing the model evaluation samples.
- **ext_sam_set** (1D or 2D `ndarray` object) – Array containing the externally provided model evaluation samples.
- **ext_mod_set** (1D or 2D `ndarray` object) – Array containing the model outputs of all specified externally provided model evaluation samples.

Generates

sam_set [2D `ndarray` object] Array containing the model evaluation samples for emulator iteration *emul_i*.

mod_set [2D `ndarray` object] Array containing the model outputs of all specified model evaluation samples for emulator iteration *emul_i*.

_get_md_var (*emul_i, par_set*)

Retrieves the model discrepancy variances, which includes all variances that are created by the model provided by the `ModelLink` instance. This method tries to call the `get_md_var()` method, and assumes a default model discrepancy variance of $1/6$ th the data value if it cannot be called. If the data value space is not linear, then this default value is calculated such to reflect that.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **par_set** (1D `ndarray` object) – Model parameter value set to calculate the model discrepancy variances for.

Returns **var_md** (2D `ndarray` object) – Variance of the model discrepancy.

_get_mock_data (*mock_par*)

Generates mock data and loads it into the `ModelLink` object that was provided during class initialization. This function overwrites the `ModelLink` properties holding the parameter estimates, data values and data errors.

Parameters **mock_par** (*1D array_like or None*) – If 1D *array_like*, use the provided parameter estimates to create the mock data. If *None*, a random parameter set will be generated as parameter estimates.

Generates

Overwrites the corresponding `ModelLink` class properties with the generated values.

`_get_n_eval_sam(emul_i)`

This function calculates the total number of emulator evaluation samples at a given emulator iteration `emul_i` from `base_eval_sam`.

Parameters `emul_i (int)` – Number indicating the requested emulator iteration.

Returns `n_eval_sam (int)` – Total number of emulator evaluation samples.

`_get_paths(root_dir, working_dir, prefix)`

Obtains the path for the root directory, working directory and parameters file for *PRISM*.

Parameters

- **root_dir (str or None)** – String containing the absolute path of the root directory where all working directories are stored. If *None*, root directory will be set to the directory this class was initialized in.
- **working_dir (str, bool or None)** – String containing the name of the working directory of the emulator in `root_dir`. If *True*, a new working directory will be created in `root_dir`. If *None* or *False*, working directory is set to the last one that was created in `root_dir` that starts with the given `prefix`. Note that providing an integer instead of a bool will not work here. If no directories are found, one will be created.
- **prefix (str or None)** – String containing a prefix that is used for naming new working directories or scan for existing ones. If *None*, all directories in `root_dir` are considered working directories and 'prism_' will be used as a prefix for new ones.

Generates

The absolute paths to the root directory, working directory, emulator master HDF5-file and *PRISM* parameters file.

`_get_uni_impl(emul_i, par_set, adj_exp_val, adj_var_val)`

Calculates the univariate implausibility values at a given emulator iteration `emul_i` for specified expectation and variance values `adj_exp_val` and `adj_var_val`, corresponding to given `par_set`.

Parameters

- **emul_i (int)** – Number indicating the requested emulator iteration.
- **par_set (1D ndarray object)** – Model parameter value set to calculate the univariate implausibility values for. Only used to pass to the `get_md_var()` method.
- **adj_exp_val, adj_var_val (1D array_like)** – The adjusted expectation and variance values to calculate the univariate implausibility for.

Returns `uni_impl_val (1D ndarray object)` – Univariate implausibility value for all requested emulator systems.

`_listen_for_calls()`

All worker ranks in the `comm` communicator start listening for calls from the corresponding controller rank and will attempt to execute the received message. Listening for calls continues until `_worker_mode` is set to *False*.

This method is automatically initialized and finalized when using the `worker_mode` context manager.

`_load_data()`

Loads in all the important pipeline data into memory for the controller rank. If it is detected that the last emulator iteration has not been analyzed yet, the implausibility analysis parameters are read in from the *PRISM* parameters file and temporarily stored in memory.

Generates

All relevant pipeline data up to the last emulator iteration is loaded into memory.

`_make_call(exec_fn, *args, **kwargs)`

Send the provided *exec_fn* to all worker ranks, if they are listening for calls, and tell them to execute it using the provided *args* and *kwargs*. All ranks that call this function will execute *exec_fn* as well.

If used within the *worker_mode* context manager, this function should only be called by the controller. If not, it should be called by all ranks that must execute *exec_fn*.

Parameters

- **`exec_fn`** (*str, callable or None*) – If string, a callable attribute of this *Pipeline* instance or a callable object that the workers should execute if not. If *None*, the workers stop listening for calls instead (disables worker mode).
- **`args`** (*tuple*) – Positional arguments that need to be provided to *exec_fn*.
- **`kwargs`** (*dict*) – Keyword arguments that need to be provided to *exec_fn*.

`_multi_call_model(emul_i, sam_set, data_idx)`

Obtain the output set corresponding to the provided *data_idx* that is generated by the model for a given model parameter sample set *sam_set*. The current emulator iteration *emul_i* is also provided in case it is required by the *ModelLink* subclass.

This is a multi-version of `_call_model()`.

Parameters

- **`emul_i`** (*int*) – Number indicating the requested emulator iteration.
- **`sam_set`** (*2D array_like*) – Model parameter sample set to calculate the model output for.
- **`data_idx`** (*list of tuples*) – The list of data identifiers for which the model is requested to return the corresponding data values.

Returns *mod_set* (2D *ndarray* object) – Model output set corresponding to given *sam_set*.

`_read_parameters(prism_par)`

Reads in all parameters in the provided *prism_par* and saves them as a dict in the current *Pipeline* instance.

`_save_data(data_dict)`

Saves a given data dict {*keyword*: *data*} at the last emulator iteration to the HDF5-file and as an data attribute to the current *Pipeline* instance.

Parameters *data_dict* (*dict*) – Dict containing the data that needs to be saved to the HDF5-file.

Keyword Arguments

- **`keyword`** (*{'impl_par', 'impl_sam', 'n_eval_sam'}*) – String specifying the type of data that needs to be saved.
- **`data`** (*{int, float, str, array_like} or dict*) – The actual data that needs to be saved at data keyword *keyword*. If dict, save every item individually.

Generates

The specified data is saved to the HDF5-file.

_save_statistics (*emul_i*, *stat_dict*)

Saves a given statistics dict {keyword: [value, unit]} at emulator iteration *emul_i* to the HDF5-file. The provided values are always saved as strings.

Parameters *emul_i* (*int*) – Number indicating the requested emulator iteration.

Keyword Arguments

- **keyword** (*str*) – String containing the name/keyword of the statistic that is being saved.
- **value** (*int*, *float* or *str*) – The value of the statistic.
- **unit** (*str*) – The unit of the statistic.

_set_impl_par (*impl_cut*)

Sets the *impl_cut* and *cut_idx* properties for implausibility evaluations using the read-in *PRISM* parameters file and the provided *impl_cut*.

Parameters *impl_cut* (*list of float* or *None*) – Incomplete, shortened *impl_cut*-offs list to be used during the analysis of this emulator iteration. If *None*, use the read-in *PRISM* parameters file instead.

Generates

impl_cut [1D *ndarray* object] Full list containing the *impl_cut*-offs for all data points provided to the emulator.

cut_idx [*int*] Index of the first *impl_cut*-off in the *impl_cut* list that is not 0.

_set_parameters ()

Sets the Pipeline parameters from the *prism_dict* property and saves them in the current *Pipeline* instance.

analyze (*, *impl_cut=None*)

Analyzes the emulator at the last emulator iteration for a large number of emulator evaluation samples. All samples that survive the implausibility checks set by the provided *impl_cut*, are used in the construction of the next emulator iteration.

Other Parameters *impl_cut* (*list of float* or *None*. *Default: None*) – Incomplete, shortened *impl_cut*-offs list to be used during the analysis of this emulator iteration. If *None*, the currently set implausibility cut-off values (*impl_cut*) will be used if this is the first analysis or the ‘*impl_cut*’ value in *prism_dict* is used if this is a reanalysis.

Generates

impl_sam [2D *ndarray* object] Array containing all emulator evaluation samples that survived the implausibility checks.

construct (*emul_i=None*, *, *analyze=True*, *ext_real_set=None*, *force=False*)

Constructs the emulator at the specified emulator iteration *emul_i*, and performs an implausibility analysis on the emulator iteration right afterward if requested (*analyze()*).

Other Parameters

- **emul_i** (*int or None. Default: None*) – If *int*, number indicating the requested emulator iteration. If *None*, the next iteration of the emulator will be constructed.
- **analyze** (*bool. Default: True*) – Bool indicating whether or not to perform an analysis after the specified emulator iteration has been successfully constructed, which is required for constructing the next iteration.
- **ext_real_set** (*list, dict or None. Default: None*) – List of arrays containing an externally calculated set of model evaluation samples and its data values, a dict with keys `['sam_set', 'mod_set']` containing these arrays or *None* if no external set needs to be used.
- **force** (*bool. Default: False*) – Controls what to do if the specified emulator iteration *emul_i* already (partly) exists. If *False*, finish construction of the specified iteration or skip it if already finished. If *True*, reconstruct the specified iteration entirely.

Generates

A new HDF5-group with the emulator iteration as its name, in the loaded emulator master file, containing emulator data required for this emulator iteration.

Notes

Using an emulator iteration that has been (partly) constructed before, will finish construction or skip it if already finished when *force* is *False*; or it will delete that and all following iterations, and reconstruct the specified iteration when *force* is *True*. Using *emul_i* = 1 and *force* is *True* is equivalent to reconstructing the entire emulator.

If no implausibility analysis is requested, then the implausibility parameters are read in from the *PRISM* parameters file and temporarily stored in memory in order to enable the usage of the `evaluate()` and `project()` methods.

details (*emul_i=None*)

Prints the details/properties of the currently loaded *Pipeline* instance at given emulator iteration *emul_i*. See `Props` for detailed descriptions of all printed properties.

Other Parameters **emul_i** (*int or None. Default: None*) – If *int*, number indicating the requested emulator iteration. If *None*, the last iteration of the emulator will be used.

Props

Working directory The relative path to the working directory of the emulator starting at the current working directory.

Emulator type The type of this emulator, corresponding to the *emul_type* of the provided *emul_type* during *Pipeline* initialization. If no emulator type was provided during initialization, this is ‘default’.

ModelLink subclass Name of the *ModelLink* subclass used to construct this emulator.

Emulation method Indicates the combination of regression and Gaussian emulation methods that have been used for this emulator.

Mock data used? Whether or not mock data has been used to construct this emulator. If so, the printed estimates for all model parameters are the parameter values used to create the mock data.

Emulator iteration The iteration of the emulator this details overview is about. By default, this is the last (partly) constructed iteration.

Construction completed? Whether or not the construction of this emulator iteration is completed. If not, the missing components for each emulator system are listed and the remaining information of this iteration is not printed.

Plausible regions? Whether or not plausible regions have been found during the analysis of this emulator iteration. If no analysis has been done yet, “N/A” will be printed.

Projections available? Whether or not projections have been created for this emulator iteration. If projections are available and analysis has been done, but with different implausibility cut-offs, a “desync” note is added. Also prints number of available projections versus maximum number of projections in parentheses.

of model evaluation samples The total number of model evaluation samples used to construct all emulator iterations up to this iteration, with the number for every individual iteration in parentheses.

of plausible/analyzed samples The number of emulator evaluation samples that passed the implausibility check out of the total number of analyzed samples in this emulator iteration. This is the number of model evaluation samples that was/will be used for the construction of the next emulator iteration. If no analysis has been done, the numbers show up as “-”.

% of parameter space remaining The percentage of the total number of analyzed samples that passed the implausibility check in this emulator iteration. If no analysis has been done, the number shows up as “-”.

of active/total parameters The number of model parameters that was considered active during the construction of this emulator iteration, compared to the total number of model parameters defined in the used *ModelLink* subclass.

of emulated data points The number of data points that have been emulated in this emulator iteration.

of emulator systems The total number of emulator systems that are required in this emulator. The number of active emulator systems is equal to the number of data points.

Parameter space Lists the name, lower and upper value boundaries and estimate (if provided) of all model parameters defined in the used *ModelLink* subclass. An asterisk is printed in front of the parameter name if this model parameter was considered active during the construction of this emulator iteration. A question mark is used instead if the construction of this emulator iteration is not finished.

evaluate (*sam_set*, *emul_i=None*)

Evaluates the given model parameter sample set *sam_set* up to given emulator iteration *emul_i*. The output of this function depends on the number of dimensions in *sam_set*. The output is always provided on the controller rank.

Parameters *sam_set* (*1D or 2D array_like or dict*) – Array containing model parameter value sets to be evaluated in the emulator up to emulator iteration *emul_i*.

Other Parameters *emul_i* (*int or None. Default: None*) – If int, number indicating the requested emulator iteration. If *None*, the last iteration of the emulator will be used.

Returns

- **impl_check** (*list of bool*) – List of bool indicating whether or not the given samples passed the implausibility check at the given emulator iteration *emul_i*.
- **emul_i_stop** (*list of int*) – List containing the last emulator iterations at which the given samples are still within the plausible region of the emulator.

- **adj_exp_val** (2D `ndarray` object) – Array containing the adjusted expectation values for all given samples.
- **adj_var_val** (2D `ndarray` object) – Array containing the adjusted variance values for all given samples.
- **uni_impl_val** (2D `ndarray` object) – Array containing the univariate implausibility values for all given samples.

Prints (if 1D *sam_set*)

emul_i_stop [int] Last emulator iteration at which the given sample is still within the plausible region of the emulator.

adj_exp_val [1D `ndarray` object] The adjusted expectation values for the given sample.

adj_var_val [1D `ndarray` object] The adjusted variance values for the given sample.

sigma_val [1D `ndarray` object] The corresponding sigma value for the given sample.

uni_impl_val [1D `ndarray` object] The univariate implausibility values for the given sample.

Notes

If given emulator iteration *emul_i* has been analyzed before, the implausibility parameters of the last analysis are used. If not, then the values are used that were read in when the emulator was loaded or that have been set by the user.

project (*emul_i=None, proj_par=None, **kwargs*)

Analyzes the emulator iteration *emul_i* and constructs a series of projection figures detailing the behavior of the model parameters corresponding to the given *proj_par*. The input and output depend on the number of model parameters *n_par*.

All optional keyword arguments (except *force*) control various aspects of drawing the projection figures and do not affect the projection data that is saved to HDF5. This is instead influenced by the *proj_res* and *proj_depth* properties.

Parameters

- **emul_i** (*int or None. Default: None*) – If int, number indicating the requested emulator iteration. If *None*, the last iteration of the emulator will be used.
- **proj_par** (*1D array_like of {int, str} or None. Default: None*) – For which model parameters to construct the projection figures. If 1D array_like, construct projection figures for all combinations of provided model parameters that are active, with a string referring to the name of the model parameter and an integer referring to the position in which the model parameter is shown in the *details()* method. If *None*, projection figures are made for all active model parameters.

Keyword Arguments

- **proj_type** (*{'2D', '3D', 'both'}. Default: '2D' (2D), 'both' (nD)*) – String indicating which projection type to create for all supplied active parameters. If *n_par == 2*, this is always '2D' (and cannot be modified).
- **figure** (*bool. Default: True*) – Whether or not to create the projection figures. If *True*, the figures are calculated, drawn and saved. If *False*, the figures are calculated and their data is returned in a dict.

- **align** (*{'row'/'horizontal', 'col'/'column'/'vertical'}*. *Default: 'col'*) – If *figure* is *True*, string indicating how to position the two subplots. If *'row'/'horizontal'*, the subplots are positioned on a single row. If *'col'/'column'/'vertical'*, the subplots are positioned on a single column.
- **show_cuts** (*bool. Default: False*) – If *figure* is *True* and *proj_type* is not *'3D'*, whether to show all implausibility cut-offs in the 2D projections (*True*) or only the first cut-off (*False*).
- **smooth** (*bool. Default: False*) – Controls what to do if a grid point contains no plausible samples, but does contain a minimum implausibility value below the first non-wildcard cut-off. If *False*, these values are kept, which can show up as artifact-like features in the projection figure. If *True*, these values are set to the first cut-off, removing them from the projection figure. Doing this may also remove interesting features. This does not affect the projection data saved to HDF5. Smoothed figures have an *'_s'* string appended to their filenames.
- **force** (*bool. Default: False*) – Controls what to do if a projection hypercube has been calculated at the emulator iteration *emul_i* before. If *False*, it will use the previously acquired projection data to create the projection figure. If *True*, it will recalculate all the data required to create the projection figure. Note that this will also delete all associated projection figures.
- **fig_kwargs** (*dict. Default: {'figsize': (6.4, 4.8), 'dpi': 100}*) – Dict of keyword arguments to be used when creating the subplots figure. It takes all arguments that can be provided to the `figure()` function.
- **impl_kwargs_2D** (*dict. Default: {}*) – Dict of keyword arguments to be used for making the minimum implausibility (top/left) plot in the 2D projection figures. It takes all arguments that can be provided to the `plot()` function.
- **impl_kwargs_3D** (*dict. Default: {'cmap': 'rainforest_r'}*) – Dict of keyword arguments to be used for making the minimum implausibility (top/left) plot in the 3D projection figures. It takes all arguments that can be provided to the `hexbin()` function.
- **los_kwargs_2D** (*dict. Default: {}*) – Dict of keyword arguments to be used for making the line-of-sight (bottom/right) plot in the 2D projection figures. It takes all arguments that can be provided to the `plot()` function.
- **los_kwargs_3D** (*dict. Default: {'cmap': 'blaze'}*) – Dict of keyword arguments to be used for making the line-of-sight (bottom/right) plot in the 3D projection figures. It takes all arguments that can be provided to the `hexbin()` function.
- **line_kwargs_est** (*dict. Default: {'linestyle': '-', 'color': 'grey'}*) – Dict of keyword arguments to be used for drawing the parameter estimate lines in both plots. It takes all arguments that can be provided to the `draw()` function.
- **line_kwargs_cut** (*dict. Default: {'color': 'r'}*) – Dict of keyword arguments to be used for drawing the implausibility cut-off line(s) in the top/left plot in the 2D projection figures. It takes all arguments that can be provided to the `draw()` function.

Returns **fig_data** (*dict of dicts*) – Dict containing the data for every requested projection figure, split up into the *'impl_min'* and *'impl_los'* dicts. For 2D projections, every dict contains a list with the x and y values. For 3D projections, it contains the x, y and z values. Note that due to the figures being interpolations, the y/z values can be below zero or the line-of-sight values being above unity.

Generates (if *figure* is *True*)

A series of projection figures detailing the behavior of the model. The lay-out and output of the projection figures depend on the type of figure:

2D projection figure: The output will feature a figure with two subplots for every active model parameter (`n_par`). Every figure gives details about the behavior of the corresponding model parameter, by showing the minimum implausibility value (top/left) and the line-of-sight depth (bottom/right) obtained at the specified parameter value, independent of the values of the other parameters.

3D projection figure (only if `n_par > 2`): The output will feature a figure with two subplots for every combination of two active model parameters that can be made ($n_par * (n_par - 1) / 2$). Every figure gives details about the behavior of the corresponding model parameters, as well as their dependency on each other. This is done by showing the minimum implausibility (top/left) and the line-of-sight depth (bottom/right) obtained at the specified parameter values, independent of the values of the remaining model parameters.

Notes

If given emulator iteration `emul_i` has been analyzed before, the implausibility parameters of the last analysis are used. If not, then the values are used that were read in when the emulator was loaded or that have been set by the user.

run (`emul_i=None, *, force=False`)

Calls the `construct()` method to start the construction of the given iteration of the emulator and creates the projection figures right afterward if this construction was successful.

Other Parameters

- **emul_i** (*int or None. Default: None*) – If int, number indicating the requested emulator iteration. If *None*, the next iteration of the emulator will be constructed.
- **force** (*bool. Default: False*) – Controls what to do if the specified emulator iteration `emul_i` already (partly) exists. If *False*, finish construction of the specified iteration or skip it if already finished. If *True*, reconstruct the specified iteration entirely.

File

Custom `File` class that has added logging and automatically uses `hdf5_file` as the HDF5-file to open.

Type `File`

base_eval_sam

Base number of emulator evaluations used to analyze the emulator systems. This number is scaled up by the number of model parameters and the current emulator iteration to generate the true number of emulator evaluations (`n_eval_sam`).

Type `int`

code_objects

Collection of pre-compiled built-in code snippets that are used in the `_evaluate_sam_set()` method.

Type dict of code objects

comm

The MPI intra-communicator that is used in this `Pipeline` instance. By default, this is `MPI.COMM_WORLD`.

Type `Intracomm`

criterion

Value indicating which criterion to use in the `lhd()` function.

Type `str`, `float` or `None`

cut_idx

The index of the first non-wildcard in a list of implausibility values. This is equivalent to the number of wildcards leading the cut-off values in `impl_cut`.

Type `int`

do_active_anal

Whether or not to do an active parameters analysis during the construction of the emulator systems.

Type `bool`

do_logging

Whether or not to save all logging messages. If *False*, all logging messages of level `INFO` and below are ignored. It also enables/disables the progress bar for making projections.

Type `bool`

emulator

The *Emulator* instance created during *Pipeline* initialization.

Type *Emulator*

freeze_active_par

Whether or not previously active parameters always stay active if possible.

Type `bool`

hdf5_file

Absolute path to the loaded master HDF5-file.

Type `str`

impl_cut

The non-wildcard univariate implausibility cut-off values for an emulator iteration. Setting it with the reduced implausibility cut-off list will change the values of `cut_idx` and this property at the last emulator iteration.

Type list of float

impl_sam

The model evaluation samples that will be added to the next emulator iteration.

Type `ndarray`

is_controller

Whether or not this MPI process is a controller rank. If no MPI is used, this is always *True*.

Type `bool`

is_worker

Whether or not this MPI process is a worker rank. If no MPI is used, this is always *False*.

Type `bool`

modellink

The *ModelLink* instance provided during *Pipeline* initialization.

Type *ModelLink*

n_eval_sam

The number of evaluation samples used to analyze an emulator iteration of the emulator systems. The number of plausible evaluation samples is stored in `n_impl_sam`. It is zero if the specified iteration has not been analyzed yet.

Type `int`

n_impl_sam

Number of model evaluation samples that passed the implausibility checks during the analysis of an emulator iteration. It is zero if the specified iteration has not been analyzed yet or has no plausible samples.

Type `int`

n_sam_init

Number of evaluation samples used to construct the first iteration of the emulator systems.

Type `int`

pot_active_par

The potentially active parameters. Only parameters from this list can become active during the active parameters analysis. If `do_active_anal` is `False`, all parameters in this list will be active.

Type `list of str`

prism_dict

Dictionary containing all *PRISM* parameters that were provided during *Pipeline* initialization.

Type `dict`

proj_depth

Number of emulator evaluations that will be used to generate the samples in every grid point for the projection figures. Note that when making 2D projections of nD models, the used depth will be this number multiplied by `proj_res`.

Type `int`

proj_res

Number of emulator evaluations that will be used to generate the grid for the projection figures.

Type `int`

rank

The rank of this MPI process in `comm`. If no MPI is used, this is always 0.

Type `int`

root_dir

Absolute path to the root directory.

Type `str`

size

The number of MPI processes in `comm`. If no MPI is used, this is always 1.

Type `int`

worker_mode

Special context manager within which all code is executed in worker mode. In worker mode, all worker ranks are continuously listening for calls from the controller rank made with `_make_call()`.

Note that all code within the context manager is executed by all ranks, with the worker ranks executing it after the controller rank exited. It is therefore advised to use an if-statement inside to make sure only the controller rank executes the code.

Using this context manager allows for easier use of *PRISM* in combination with serial/OpenMP codes (like MCMC methods).

Type `_GeneratorContextManager`

working_dir

Absolute path to the working directory.

Type `str`

10.1 Classes

10.1.1 Emulator

class prism.emulator.**Emulator** (*pipeline_obj*, *modellink_obj*)
Defines the *Emulator* base class of the *PRISM* package.

Description

The *Emulator* class is the backbone of the *PRISM* package, holding all tools necessary to construct, load, save and evaluate the emulator of a model. It performs many checks to see if the provided *ModelLink* object is compatible with the current emulator, advises the user on alternatives when certain operations are requested, automatically takes care of distributing emulator systems over MPI ranks and more.

Even though the purpose of the *Emulator* class is to hold only information about the emulator and therefore does not require any details about the provided *ModelLink* object, it will keep track of changes made to it. This is to allow the user to modify the properties of the *ModelLink* subclass without causing any desynchronization problems by accident.

The *Emulator* class requires to be linked to an instance of the *Pipeline* class and will automatically attempt to do so when initialized. By default, this class should only be initialized from within a *Pipeline* object.

__init__ (*pipeline_obj*, *modellink_obj*)
Initialize an instance of the *Emulator* class.

Parameters

- **pipeline_obj** (*Pipeline* object) – The *Pipeline* instance this *Emulator* instance should be linked to.
- **modellink_obj** (*ModelLink* object) – The *ModelLink* instance that should be linked to *pipeline_obj*.

`_assign_data_idx` (*emul_i*)

Determines the emulator system each data point in the provided emulator iteration *emul_i* should be assigned to, in order to make sure that recurring data points have the same emulator system index as in the previous emulator iteration. If multiple options are possible, data points are assigned such to spread them as much as possible.

Parameters *emul_i* (*int*) – Number indicating the requested emulator iteration.

Returns

- **`data_to_emul_s`** (*list of int*) – The index of the emulator system that each data point should be assigned to.
- **`n_emul_s`** (*int*) – The total number of active and passive emulator systems there will be in the provided emulator iteration.

Examples

If the number of data points is less than the previous iteration:

```
>>> emul_i = 2
>>> self._data_idx[emul_i-1]
['A', 'B', 'C', 'D', 'E']
>>> self._modellink._data_idx
['B', 'F', 'G', 'E']
>>> self._assign_data_idx(emul_i)
([1, 3, 2, 4], 5)
```

If the number of data points is more than the previous iteration:

```
>>> emul_i = 2
>>> self._data_idx[emul_i-1]
['A', 'B', 'C', 'D', 'E']
>>> self._modellink._data_idx
['B', 'F', 'G', 'E', 'A', 'C']
>>> self._assign_data_idx(emul_i)
([1, 5, 3, 4, 0, 2], 6)
```

If there is no previous iteration:

```
>>> emul_i = 1
>>> self._data_idx[emul_i-1]
[]
>>> self._modellink._data_idx
['B', 'F', 'G', 'E', 'A', 'C']
>>> self._assign_data_idx(emul_i)
([5, 4, 3, 2, 1, 0], 6)
```

`_assign_emul_s` (*emul_i*)

Determines which emulator systems (files) should be assigned to which MPI rank in order to balance the number of active emulator systems on every rank for every iteration up to the provided emulator iteration *emul_i*. If multiple choices can achieve this, the emulator systems are automatically spread out such that the total number of active emulator systems on a single rank is also balanced as much as possible.

Parameters *emul_i* (*int*) – Number indicating the requested emulator iteration.

Returns *emul_s_to_core* (*list of lists*) – A list containing the emulator systems that have been assigned to the corresponding MPI rank by the controller.

Notes

Currently, this function only uses high-level MPI. Additional speed can be obtained by also implementing low-level MPI, which will potentially be done in the future.

`_cleanup_emul_files (emul_i)`

Opens all emulator HDF5-files and removes the provided emulator iteration *emul_i* and subsequent iterations from them. Also removes any related projection figures that have default names. If *emul_i* == 1, all emulator HDF5-files are removed instead.

Parameters *emul_i* (*int*) – Number indicating the requested emulator iteration.

`_construct_iteration (emul_i)`

Constructs the emulator iteration corresponding to the provided *emul_i*, by performing the given emulation method and pre-calculating the prior expectation and variance values of the used model evaluation samples.

Parameters *emul_i* (*int*) – Number indicating the requested emulator iteration.

Generates

All data sets that are required to evaluate the emulator at the constructed iteration.

`_create_new_emulator ()`

Creates a new master HDF5-file that holds all the information of a new emulator and writes all important emulator details to it. Afterwards, resets all loaded emulator data and prepares the HDF5-file and emulator for the construction of the first emulator iteration.

Generates

A new master HDF5-file ‘prism.hdf5’ contained in the working directory specified in the *Pipeline* instance, holding all information required to construct the first iteration of the emulator.

`_do_regression (emul_i, emul_s_seq)`

Performs a forward stepwise linear regression for all requested emulator systems *emul_s_seq* in the provided emulator iteration *emul_i*. Calculates what the expectation values of all polynomial coefficients are. The polynomial order that is used in the regression depends on *poly_order*.

Parameters

- *emul_i* (*int*) – Number indicating the requested emulator iteration.
- *emul_s_seq* (*list of int*) – List of numbers indicating the requested emulator systems.

Generates (for every emulator system)

rsdl_var [float] Residual variance of the regression function.

regr_score [float] Fit-score of the regression function.

poly_coef [1D *ndarray* object] Array containing the expectation values of the non-zero polynomial coefficients.

poly_powers [2D *ndarray* object] Array containing the powers of the non-zero polynomial terms in the regression function.

poly_idx [1D *ndarray* object] Array containing the indices of the non-zero polynomial terms in the regression function.

poly_coef_cov [1D `ndarray` object (if `use_regr_cov` is `True`)] Array containing the covariance values of the non-zero polynomial coefficients.

_evaluate (*emul_i*, *par_set*)

Evaluates the emulator systems *emul_s_seq* at iteration *emul_i* for given *par_set*.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **emul_s_seq** (*list of int*) – List of numbers indicating the requested emulator systems.
- **par_set** (1D `ndarray` object) – Model parameter value set to evaluate the emulator at.

Returns

- **adj_exp_val** (1D `ndarray` object) – Adjusted emulator expectation value for all requested emulator systems on this MPI rank.
- **adj_var_val** (1D `ndarray` object) – Adjusted emulator variance value for all requested emulator systems on this MPI rank.

_get_active_par (*emul_i*, *emul_s_seq*)

Determines the active parameters to be used for every emulator system listed in *emul_s_seq* in the provided emulator iteration *emul_i*. Uses backwards stepwise elimination to determine the set of active parameters. The polynomial order that is used in the stepwise elimination depends on *poly_order*.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **emul_s_seq** (*list of int*) – List of numbers indicating the requested emulator systems.

Generates (for every emulator system)

active_par_data [1D `ndarray` object] Array containing the indices of all the parameters that are active in the emulator iteration *emul_i*.

_get_adj_exp (*emul_i*, *emul_s_seq*, *par_set*, *cov_vec*)

Calculates the adjusted emulator expectation values for requested emulator systems *emul_s_seq* at a given emulator iteration *emul_i* for specified parameter set *par_set* and corresponding covariance vector *cov_vec*.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **emul_s_seq** (*list of int*) – List of numbers indicating the requested emulator systems.
- **par_set** (1D `ndarray` object) – Model parameter value set to calculate the adjusted emulator expectation for.
- **cov_vec** (2D `ndarray` object) – Covariance vector corresponding to *par_set*.

Returns **adj_exp_val** (1D `ndarray` object) – Adjusted emulator expectation value for all requested emulator systems on this MPI rank.

_get_adj_var (*emul_i*, *emul_s_seq*, *par_set*, *cov_vec*)

Calculates the adjusted emulator variance values for requested emulator systems *emul_s_seq* at a given emulator iteration *emul_i* for specified parameter set *par_set* and corresponding covariance vector *cov_vec*.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.

- **emul_s_seq** (*list of int*) – List of numbers indicating the requested emulator systems.
- **par_set** (1D `ndarray` object) – Model parameter value set to calculate the adjusted emulator variance for.
- **cov_vec** (2D `ndarray` object) – Covariance vector corresponding to *par_set*.

Returns **adj_var_val** (1D `ndarray` object) – Adjusted emulator variance value for all requested emulator systems on this MPI rank.

_get_cov (*emul_i, emul_s_seq, par_set1, par_set2*)

Calculates the full emulator covariances for requested emulator systems *emul_s_seq* at emulator iteration *emul_i* for given parameter sets *par_set1* and *par_set2*. The contributions to these covariances depend on *method*.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **emul_s_seq** (*list of int*) – List of numbers indicating the requested emulator systems.
- **par_set1, par_set2** (1D `ndarray` object or *None*) – If *par_set1* and *par_set2* are both not *None*, calculate covariances for *par_set1* with *par_set2*. If *par_set1* is not *None* and *par_set2* is *None*, calculate covariances for *par_set1* with *sam_set* (covariance vector). If *par_set1* and *par_set2* are both *None*, calculate covariances for *sam_set* (covariance matrix). When not *None*, *par_set* is the model parameter value set to calculate the covariances for.

Returns **cov** (1D, 2D or 3D `ndarray` object) – Depending on the arguments provided, a covariance value, vector or matrix for requested emulator systems.

_get_cov_matrix (*emul_i, emul_s_seq*)

Calculates the (inverse) matrix of covariances between known model evaluation samples for requested emulator systems *emul_s_seq* at emulator iteration *emul_i*.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **emul_s_seq** (*list of int*) – List of numbers indicating the requested emulator systems.

Generates

cov_mat [3D `ndarray` object] Matrix containing the covariances between all known model evaluation samples for requested emulator systems.

cov_mat_inv [3D `ndarray` object] Inverse of covariance matrix for requested emulator systems.

_get_default_parameters ()

Generates a dict containing default values for all emulator parameters.

Returns **par_dict** (*dict*) – Dict containing all default emulator parameter values.

_get_emul_i (*emul_i, cur_iter*)

Checks if the provided emulator iteration *emul_i* can be requested or replaces it if *None* was provided. This method requires all MPI ranks to call it simultaneously.

Parameters

- **emul_i** (*int or None*) – Number indicating the requested emulator iteration.
- **cur_iter** (*bool*) – Bool determining whether the current (*True*) or the next (*False*) emulator iteration is requested.

Returns `emul_i` (*int*) – The requested emulator iteration that passed the check.

`_get_exp_dot_term` (*emul_i*, *emul_s_seq*)

Pre-calculates the second expectation adjustment dot-term for requested emulator systems *emul_s_seq* at a given emulator iteration *emul_i* for all model evaluation samples and saves it for later use.

Parameters

- **`emul_i`** (*int*) – Number indicating the requested emulator iteration.
- **`emul_s_seq`** (*list of int*) – List of numbers indicating the requested emulator systems.

Generates

`exp_dot_term` [2D `ndarray` object] 2D array containing the pre-calculated values for the second adjustment dot-term of the adjusted expectation for requested emulator systems.

`_get_inv_matrix` (*matrix*)

Calculates the inverse of a given *matrix*. Right now only uses the `inv()` function.

Parameters **`matrix`** (2D array_like) – Matrix to be inverted.

Returns **`matrix_inv`** (2D `ndarray` object) – Inverse of the given *matrix*.

`_get_prior_exp` (*emul_i*, *emul_s_seq*, *par_set*)

Calculates the prior expectation value for requested emulator systems *emul_s_seq* at a given emulator iteration *emul_i* for specified parameter set *par_set*. This expectation depends on *method*.

Parameters

- **`emul_i`** (*int*) – Number indicating the requested emulator iteration.
- **`emul_s_seq`** (*list of int*) – List of numbers indicating the requested emulator systems.
- **`par_set`** (1D `ndarray` object or *None*) – If *None*, calculate the prior expectation values of *sam_set*. If not *None*, calculate the prior expectation value for the given model parameter value set.

Returns **`prior_exp`** (1D or 2D `ndarray` object) – Prior expectation values for either *sam_set* or *par_set* for requested emulator systems.

`_get_regr_cov` (*emul_i*, *emul_s_seq*, *par_set1*, *par_set2*)

Calculates the covariances of the regression function for requested emulator systems *emul_s_seq* at emulator iteration *emul_i* for given parameter sets *par_set1* and *par_set2*.

Parameters

- **`emul_i`** (*int*) – Number indicating the requested emulator iteration.
- **`emul_s_seq`** (*list of int*) – List of numbers indicating the requested emulator systems.
- **`par_set1`, `par_set2`** (1D `ndarray` object or *None*) – If *par_set1* and *par_set2* are both not *None*, calculate regression covariances for *par_set1* with *par_set2*. If *par_set1* is not *None* and *par_set2* is *None*, calculate regression covariances for *par_set1* with *sam_set* (covariance vector). If *par_set1* and *par_set2* are both *None*, calculate regression covariances for *sam_set* (covariance matrix). When not *None*, *par_set* is the model parameter value set to calculate the regression covariances for.

Returns **`regr_cov`** (1D, 2D or 3D `ndarray` object) – Depending on the arguments provided, a regression covariance value, vector or matrix for requested emulator systems.

`_load_data` (*emul_i*)

Loads in all the important emulator data up to emulator iteration *emul_i* into memory.

Parameters `emul_i` (*int*) – Number indicating the requested emulator iteration.

Generates

All relevant emulator data up to emulator iteration `emul_i` is loaded into memory.

`_load_emulator` (*modellink_obj*)

Checks if the provided working directory contains a constructed emulator and loads in the emulator data accordingly.

Parameters `modellink_obj` (*ModelLink* object) – Instance of the *ModelLink* class that links the emulated model to this *Pipeline* object.

`_prepare_new_iteration` (*emul_i*)

Prepares the emulator for the construction of a new iteration `emul_i`. Checks if this iteration can be prepared or if it has been prepared before, and acts accordingly.

Parameters `emul_i` (*int*) – Number indicating the requested emulator iteration.

Returns `reload` (*bool*) – Bool indicating whether or not the controller rank of the *Pipeline* instance needs to reload its data.

Generates

A new group in the master HDF5-file with the emulator iteration as its name, containing subgroups corresponding to all emulator systems that will be used in this iteration.

Notes

Preparing an iteration that has been prepared before, causes that and all subsequent iterations of the emulator to be deleted. A check is carried out to see if it was necessary to reprepare the requested iteration and a warning is given if this check fails.

`_read_data_idx` (*emul_s_group*)

Reads in and combines the parts of the data point identifier that is assigned to the provided *emul_s_group*.

Parameters `emul_s_group` (*Group* object) – The HDF5-group from which the data point identifier needs to be read in.

Returns `data_idx` (*tuple of {int, float, str}*) – The combined data point identifier.

`_retrieve_parameters` ()

Reads in the emulator parameters from the provided working directory and saves them in the current *Emulator* instance.

`_save_data` (*emul_i*, *lemul_s*, *data_dict*)

Saves a given data dict {`keyword:` `data`} at the given emulator iteration `emul_i` and local emulator system `lemul_s` to the HDF5-file and as an data attribute to the current *Emulator* instance.

Parameters

- `emul_i` (*int*) – Number indicating the requested emulator iteration.
- `lemul_s` (*int or None*) – Number indicating the requested local emulator system. If *None*, use the master emulator file instead.
- `data_dict` (*dict*) – Dict containing the data that needs to be saved to the HDF5-file.

Keyword Arguments

- **keyword** (*{'active_par', 'active_par_data', 'cov_mat', 'exp_dot_term', 'mod_real_set', 'regression'}*) – String specifying the type of data that needs to be saved.
- **data** (*{int, float, str, array_like} or dict*) – The actual data that needs to be saved at data keyword *keyword*. If dict, save every item individually.

Generates

The specified data is saved to the HDF5-file.

`_set_mock_data()`

Loads previously used mock data into the *ModelLink* object, overwriting the parameter estimates, data values, data errors, data spaces and data identifiers with their mock equivalents.

Generates

Overwrites the corresponding *ModelLink* class properties with the previously used values (taken from the first emulator iteration).

`_set_modellink(modellink_obj, modellink_loaded)`

Sets the *ModelLink* object that will be used for constructing this emulator. If a constructed emulator is present, checks if provided *modellink_obj* argument matches the *ModelLink* subclass used to construct it.

Parameters

- **modellink_obj** (*ModelLink* object) – Instance of the *ModelLink* class that links the emulated model to this *Pipeline* object. The provided *ModelLink* object must match the one used to construct the loaded emulator.
- **modellink_loaded** (*str or None*) – If str, the name of the *ModelLink* subclass that was used to construct the loaded emulator. If *None*, no emulator is loaded.

`_set_parameters()`

Sets the Emulator parameters from the *prism_dict* property and saves them in the current *Emulator* instance.

`_write_data_idx(emul_s_group, data_idx)`

Splits a given *data_idx* up into individual parts and saves it as an attribute to the provided *emul_s_group*.

Parameters

- **emul_s_group** (*Group* object) – The HDF5-group to which the data point identifier needs to be saved.
- **data_idx** (*tuple of {int, float, str}*) – The data point identifier to be saved.

`active_emul_s`

The indices of the emulator systems on this MPI rank that are active.

Type list of int

`active_par`

The model parameter names that are considered active. Only available on the controller rank.

Type list of str

`active_par_data`

The model parameter names that are considered active for every emulator system on this MPI rank.

Type list of str

ccheck

The emulator system components that are still required to complete the construction of an emulator iteration on this MPI rank. The controller rank additionally lists the required components that are emulator iteration specific ('mod_real_set' and 'active_par').

Type list of str

cov_mat_inv

The inverses of the covariance matrices for every emulator system on this MPI rank.

Type list of ndarray

data_idx_to_core

List of the data identifiers that were assigned to the emulator systems listed in *emul_s_to_core*. Only available on the controller rank.

Type list of lists

emul_i

The last emulator iteration that is fully constructed for all emulator systems on this MPI rank.

Type int

emul_load

Whether or not a previously constructed emulator is currently loaded.

Type bool

emul_s

The indices of the emulator systems that are assigned to this MPI rank.

Type list of int

emul_s_to_core

List of the indices of the emulator systems that are assigned to every MPI rank. Only available on the controller rank.

Type list of lists

emul_type

The type of emulator that is currently loaded. This determines the way in which the *Pipeline* instance will treat this *Emulator* instance.

Type str

exp_dot_term

The second expectation adjustment dot-term values of all model evaluation samples for every emulator system on this MPI rank.

Type list of ndarray

l_corr

The Gaussian correlation lengths for all model parameters, which is defined as the maximum distance between two values of a specific model parameter within which the Gaussian contribution to the correlation between the values is still significant.

Type ndarray

method

The emulation method to use for constructing the emulator. Possible are 'gaussian', 'regression' and 'full'.

Type str

mod_set

The model outputs corresponding to the samples in *sam_set* for every emulator system on this MPI rank.

Type list of `ndarray`

n_cross_val

Number of (k-fold) cross-validations that are used for determining the quality of the regression process. It is set to zero if cross-validations are not used. If `method == 'gaussian'` and `do_active_anal` is `False`, this number is not required.

Type `int`

n_emul_s

Number of emulator systems assigned to this MPI rank.

Type `int`

n_emul_s_tot

Total number of emulator systems assigned to all MPI ranks combined. Only available on the controller rank.

Type `int`

n_sam

Number of model evaluation samples that have been/will be used to construct an emulator iteration.

Type `int`

poly_coef

The non-zero coefficients for the polynomial terms in the regression function for every emulator system on this MPI rank. Empty if `method == 'gaussian'`.

Type list of `ndarray`

poly_coef_cov

The covariances for all coefficients in `poly_coef` for every emulator system on this MPI rank. Empty if `method == 'gaussian'` or `use_regr_cov` is `False`.

Type list of `ndarray`

poly_idx

The indices for all polynomial terms with non-zero coefficients in the regression function for every emulator system on this MPI rank. Empty if `method == 'gaussian'`.

Type list of `ndarray`

poly_order

Polynomial order that is considered for the regression process. If `method == 'gaussian'` and `do_active_anal` is `False`, this number is not required.

Type `int`

poly_powers

The powers for all polynomial terms with non-zero coefficients in the regression function for every emulator system on this MPI rank. Empty if `method == 'gaussian'`.

Type list of `ndarray`

rsdl_var

The residual variances for every emulator system on this MPI rank. Obtained from regression process and replaces the Gaussian sigma. Empty if `method == 'gaussian'`.

Type list of `float`

sam_set

The model evaluation samples that have been/will be used to construct the specified emulator iteration.

Type `ndarray`

sigma

Value of the Gaussian sigma. If *method* != 'gaussian', this value is not required, since it is obtained from the regression process instead.

Type float

use_mock

Whether or not mock data has been used for the construction of this emulator instead of actual data. If *True*, changes made to the data in the provided *ModelLink* object are ignored.

Type bool

use_regr_cov

Whether or not to take into account the regression covariance when calculating the covariance of the emulator, in addition to the Gaussian covariance. If *method* == 'gaussian', this bool is not required. If *method* == 'regression', this bool is always set to *True*.

Type bool

11.1 Classes

11.1.1 GaussianLink

class prism.modellink.**GaussianLink** (*n_gaussians=1, *args, **kwargs*)
ModelLink class wrapper for a simple Gaussian model, used for testing the functionality of the *PRISM* pipeline.

Formatting data_idx

x [int] The value that needs to be used for x in the function $\sum_i A_i \exp\left(-\frac{(x-B_i)^2}{2C_i^2}\right)$ to obtain the data value.

__init__ (*n_gaussians=1, *args, **kwargs*)
 Initialize an instance of the *GaussianLink* class.

Other Parameters **n_gaussians** (*int. Default: 1*) – The number of Gaussians to use for the Gaussian model in this instance. The resulting number of model parameters *n_par* will be $3 * n_{gaussians}$.

n_gaussians
 Number of Gaussians used in this *GaussianLink* instance.

Type *int*

11.1.2 ModelLink

class prism.modellink.**ModelLink** (*, *model_parameters=None, model_data=None*)
 Provides an abstract base class definition that allows the *Pipeline* class to be linked to any model/test object of choice. Every model wrapper used in the *Pipeline* class must be an instance of the *ModelLink* class.

Description

The `ModelLink` class is an abstract base class, which forms the base for wrapping a model and allowing *PRISM* to use it effectively. Because it is mandatory for every model to be wrapped in a user-made `ModelLink` subclass, several tools are provided to the user to make this as versatile as possible.

The `ModelLink` class uses three properties that define the way the subclass will be used by *PRISM*: `name`, `call_type` and `MPI_call`. The first defines what the name of the subclass is, which is used by *PRISM* to identify the subclass with and check if one did not use a different subclass by accident. The other two are flags that determine how the `call_model()` method should be used. These three properties can be set anywhere during the initialization of the `ModelLink` subclass, or are set to a default value if they are not modified.

The model parameters and comparison data can be set in two different ways. They can be hard-coded into the `ModelLink` subclass by altering the `get_default_model_parameters()` and `get_default_model_data()` methods or set by providing them during class initialization. A combination of both is also possible. More details on this can be found in `__init__()`.

The `ModelLink` class has two abstract methods that must be overridden before the subclass can be initialized. The `call_model()` method is the most important method, as it provides *PRISM* with a way of calling the model wrapped in the `ModelLink` subclass. The `get_md_var()` method allows for *PRISM* to calculate the model discrepancy variance.

Notes

The `__init__()` method may be extended by the `ModelLink` subclass, but the superclass version must always be called.

If required, one can use the `test_subclass()` function to test a `ModelLink` subclass on correct functionality.

`__ModelLink__set_model_data(add_model_data)`

Generates the model data properties from the default model data and the additional input argument `add_model_data`.

Parameters `add_model_data` (*array_like, dict, str or None*) – Anything that can be converted to a dict that provides non-default model data information or *None* if only default data is used from `get_default_model_data()`.

Generates

n_data [int] Number of provided data points.

data_val [list] List with values of provided data points.

data_err [list of lists] List with upper and lower 1σ -confidence levels of provided data points.

data_spc [list] List with types of value space ({'lin', 'log', 'ln'}) of provided data points.

data_idx [list of tuples] List with user-defined data point identifiers.

`__ModelLink__set_model_parameters(add_model_parameters)`

Generates the model parameter properties from the default model parameters and the additional input argument `add_model_parameters`.

Parameters `add_model_parameters` (*array_like, dict, str or None*) – Anything that can be converted to a dict that provides non-default model parameters information or *None* if only default information is used from `get_default_model_parameters()`.

Generates

n_par [int] Number of model parameters.

par_name [list] List with model parameter names.

par_rng [ndarray object] Array containing the lower and upper values of the model parameters.

par_est [list] List containing user-defined estimated values of the model parameters. Contains *None* in places where estimates were not provided.

__init__ (*, *model_parameters=None, model_data=None*)

Initialize an instance of the *ModelLink* subclass.

Other Parameters **model_parameters**, **model_data** (*array_like, dict, str or None. Default: None*) – Anything that can be converted to a dict that provides non-default model parameters/data information or *None* if only default information is used from *get_default_model_parameters()* or *get_default_model_data()*. For more information on the lay-out of these dicts, see Notes.

If *array_like*, *dict(model_parameters/model_data)* must generate a dict with the correct lay-out. If *dict*, the dict itself must have the correct lay-out. If *str*, the string must be the path to a file containing the dict keys in the first column and the dict values in the second column, which combined generate a dict with the correct lay-out.

Notes (model_parameters)

The model parameters dict requires to have the name of the parameters as the keyword, and a 1D list containing the lower bound, the upper bound and, if applicable, the estimate of this parameter. It is not required to provide an estimate for every parameter. The estimates are used to draw illustrative lines when making projection figures. An example of a model parameters file can be found in the ‘data’ folder of the *PRISM* package. If required, one can use the *convert_parameters()* function to validate their formatting.

Formatting : {par_name: [lower_bnd, upper_bnd, par_est]}

Notes (model_data)

The model data dict requires to have the data identifiers (*data_idx*) as the keyword, and a 1D list containing the data value (*data_val*); the data errors (*data_err*) and the data space (*data_spc*).

If the data errors are given with one value, then the data points are assumed to have a centered 1σ -confidence interval. If the data errors are given with two values, then the data points are assumed to have a 1σ -confidence interval defined by the provided upper and lower errors.

The data spaces are one of five strings ({‘lin’, ‘log’ or ‘log_10’, ‘ln’ or ‘log_e’}) indicating in which of the three value spaces (linear, log, ln) the data values are. It defaults to ‘lin’ if it is not provided.

The data identifier is a sequence of ints, floats and strings that is unique for every data point. *PRISM* uses it to identify a data point with, which is required in some cases (like MPI), while the model itself can use it as a description of the operations required to extract the data point from the model output. It can be provided as any sequence of any length for any data point. If any sequence contains a single element, it is replaced by just that element instead of a tuple.

A simple example of a data identifier is *f(*

check_md_var (*md_var, name*)

Checks validity of provided set of model discrepancy variances *md_var* in this *ModelLink* instance.

Parameters

- **md_var** (*1D or 2D array_like or dict*) – Model discrepancy variance set to validate in this *ModelLink* instance.
- **name** (*str*) – The name of the model discrepancy set, which is used in the error message if the validation fails.

Returns **md_var** (2D *ndarray* object) – The (converted) provided *md_var* if the validation was successful. If *md_var* was a dict, it will be converted to a *ndarray* object.

_check_mod_set (*mod_set, name*)

Checks validity of provided set of model outputs *mod_set* in this *ModelLink* instance.

Parameters

- **mod_set** (*1D or 2D array_like or dict*) – Model output (set) to validate in this *ModelLink* instance.
- **name** (*str*) – The name of the model output (set), which is used in the error message if the validation fails.

Returns **mod_set** (1D or 2D *ndarray* object) – The provided *mod_set* if the validation was successful. If *mod_set* was a dict, it will be converted to a *ndarray* object.

_check_sam_set (*sam_set, name*)

Checks validity of provided set of model parameter samples *sam_set* in this *ModelLink* instance.

Parameters

- **sam_set** (*1D or 2D array_like or dict*) – Parameter/sample set to validate in this *ModelLink* instance.
- **name** (*str*) – The name of the parameter/sample set, which is used in the error message if the validation fails.

Returns **sam_set** (1D or 2D *ndarray* object) – The provided *sam_set* if the validation was successful. If *sam_set* was a dict, it will be converted to a *ndarray* object.

_get_backup_path (*emul_i, suffix*)

Returns the absolute path to a backup file made by this *ModelLink* instance, using the provided *emul_i* and *suffix*.

This method is used by the *_make_backup()* and *_read_backup()* methods, and should not be called directly.

Parameters

- **emul_i** (*int*) – The emulator iteration for which a backup filepath is needed.
- **suffix** (*str or None*) – If *str*, determine path to associated backup file using provided *suffix*. If *suffix* is empty, obtain last created backup file. If *None*, create a new path to a backup file.

Returns **filepath** (*str*) – Absolute path to requested backup file.

_get_model_par_seq (*par_seq, name*)

Converts a provided sequence *par_seq* of model parameter names and indices to a list of indices, removes duplicates and checks if every provided name/index is valid.

Parameters

- **par_seq** (*1D array_like of {int, str}*) – A sequence of integers and strings determining which model parameters need to be used for a certain operation.

- **name** (*str*) – A string stating the name of the variable the result of this method will be stored in. Used for error messages.

Returns `par_seq_conv` (*list of int*) – The provided sequence `par_seq` converted to a sorted list of model parameter indices.

`_make_backup` (**args, **kwargs*)

WARNING: This is an advanced utility method and probably will not work unless used properly. Use with caution!

Creates an HDF5-file backup of the provided *args* and *kwargs* when called by the `call_model()` method or any of its inner functions. Additionally, the backup will contain the *emul_i*, *par_set* and *data_idx* values that were passed to the `call_model()` method. It also contains the version of *PRISM* that made the backup. The backup can be restored using the `_read_backup()` method.

If it is detected that this method is used incorrectly, a `RequestWarning` is raised (and the method returns) rather than a `RequestError`, in order to not disrupt the call to `call_model()`.

Parameters

- **args** (*tuple*) – All positional arguments that must be stored in the backup file.
- **kwargs** (*dict*) – All keyword arguments that must be stored in the backup file.

Notes

The name of the created backup file contains the value of *emul_i*, *name* and a random string to avoid replacing an already existing backup file.

The saved *emul_i*, *par_set* and *data_idx* are the values these variables have locally in the `call_model()` method at the point this method is called. Because of this, making any changes to them may cause problems and is therefore heavily discouraged. If changes are necessary, it is advised to copy them to a different variable first.

`_read_backup` (*emul_i, *, suffix=None*)

Reads in a backup HDF5-file created by the `_make_backup()` method, using the provided *emul_i* and the value of *name*.

Parameters *emul_i* (*int*) – The emulator iteration that was provided to the `call_model()` method when the backup was made.

Other Parameters *suffix* (*str or None. Default: None*) – The suffix of the backup file (everything between parentheses) that needs to be read. If *None*, the last created backup will be read.

Returns

- **filename** (*str*) – The absolute path to the backup file that has been read.
- **data** (dict with keys ('*emul_i*', '*prism_version*', '*par_set*', '*data_idx*', '*args*', '*kwargs*'))
– A dict containing the data that was provided to the `_make_backup()` method.

`_to_par_space` (*sam_set*)

Converts provided *sam_set* from unit space ([0, 1]) to parameter space ([*lower_bnd*, *upper_bnd*]).

`_to_unit_space` (*sam_set*)

Converts provided *sam_set* from parameter space ([*lower_bnd*, *upper_bnd*]) to unit space ([0, 1]).

`call_model` (*emul_i, par_set, data_idx*)

Calls the model wrapped in this `ModelLink` subclass at emulator iteration *emul_i* for model parameter values *par_set* and returns the data points corresponding to *data_idx*.

This method is called with solely keyword arguments.

This is an abstract method and must be overridden by the *ModelLink* subclass.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **par_set** (dict of `float64`) – Dict containing the values for all model parameters corresponding to the requested model realization(s). If model is single-called, dict item is formatted as `{par_name: par_val}`. If multi-called, it is formatted as `{par_name: [par_val_1, par_val_2, ..., par_val_n]}`.
- **data_idx** (*list of tuples*) – List containing the user-defined data point identifiers corresponding to the requested data points.

Returns **data_val** (*1D or 2D array_like or dict*) – Array containing the data values corresponding to the requested data points generated by the requested model realization(s). If model is multi-called, *data_val* is of shape `(n_sam, n_data)`. If dict, it has the identifiers in *data_idx* as its keys with either scalars or 1D array_likes as its values.

get_default_model_data()

Returns the default model data to use for every instance of this *ModelLink* subclass. By default, returns *_default_model_data*.

get_default_model_parameters()

Returns the default model parameters to use for every instance of this *ModelLink* subclass. By default, returns *_default_model_parameters*.

get_md_var (*emul_i, par_set, data_idx*)

Calculates the linear model discrepancy variance at a given emulator iteration *emul_i* for model parameter values *par_set* and given data points *data_idx* for the model wrapped in this *ModelLink* subclass.

This method is always single-called by one MPI rank with solely keyword arguments.

This is an abstract method and must be overridden by the *ModelLink* subclass.

Parameters

- **emul_i** (*int*) – Number indicating the requested emulator iteration.
- **par_set** (dict of `float64`) – Dict containing the values for all model parameters corresponding to the requested model realization.
- **data_idx** (*list of tuples*) – List containing the user-defined data point identifiers corresponding to the requested data points.

Returns **md_var** (*1D or 2D array_like*) – Array containing the linear model discrepancy variance values corresponding to the requested data points. If 1D array_like, data is assumed to have a centered one sigma confidence interval. If 2D array_like, the values determine the upper and lower variances and the array is of shape `(n_data, 2)`. If dict, it has the identifiers in *data_idx* as its keys with either scalars or 1D array_likes of length 2 as its values.

Notes

The returned model discrepancy variance values must be of linear form, even for those data values that are returned in logarithmic form by the *call_model()* method. If not, the possibility exists that the emulation process will not converge properly.

get_str_repr()

Returns a list of string representations of all additional input arguments with which this *ModelLink* subclass was initialized.

MPI_call

Whether *call_model()* can/should be called by all MPI ranks simultaneously instead of by the controller. By default, only the controller rank calls the model (False).

Type bool

_default_model_data

The default model data to use for every instance of this *ModelLink* subclass.

Type dict

_default_model_parameters

The default model parameters to use for every instance of this *ModelLink* subclass.

Type dict

call_type

String indicating whether *call_model()* should be supplied with a single evaluation sample ('single') or a set of samples ('multi'), or can be supplied with both ('hybrid'). By default, single model calls are requested ('single').

Type str

data_err

The upper and lower 1σ -confidence levels of provided data points.

Type list of float

data_idx

The user-defined data point identifiers.

Type list of tuples

data_spc

The types of value space ({'lin', 'log', 'ln'}) of provided data points.

Type list of str

data_val

The values of provided data points.

Type list of float

multi_call

Whether *call_model()* can/should be supplied with a set of evaluation samples. At least one of *single_call* and *multi_call* must be *True*. By default, single model calls are requested (False).

Type bool

n_data

Number of provided data points.

Type int

n_par

Number of model parameters.

Type int

name

Name associated with an instance of this *ModelLink* subclass. By default, it is set to the name of this *ModelLink* subclass. Can be manually manipulated to allow for more user control.

Type `str`

par_est

The user-defined estimated values of the model parameters. Contains *None* in places where estimates were not provided.

Type `dict of {float, None}`

par_name

List with model parameter names.

Type `list of str`

par_rng

The lower and upper values of the model parameters.

Type `ndarray`

single_call

Whether *call_model()* can/should be supplied with a single evaluation sample. At least one of *single_call* and *multi_call* must be *True*. By default, single model calls are requested (*True*).

Type `bool`

11.1.3 SineWaveLink

class `prism.modellink.SineWaveLink(*args, **kwargs)`

ModelLink class wrapper for a simple sine wave model, used for testing the functionality of the *PRISM* pipeline.

Formatting data_idx

x [int] The value that needs to be used for *x* in the function $A + 0.1 * B * \sin(C * x + D)$ to obtain the data value.

11.2 Utilities

`prism.modellink.convert_data(model_data)`

Converts the provided *model_data* into a full data dict, taking into account all formatting options, and returns it.

This function can be used externally to check how the provided *model_data* would be interpreted when provided to the *ModelLink* subclass. Its output can be used for the ‘model_data’ input argument.

Parameters *model_data* (*array_like, dict or str*) – Anything that can be converted to a dict that provides model data information.

Returns *data_dict* (*dict*) – Dict with the provided *model_data* converted to its full format.

`prism.modellink.convert_parameters(model_parameters)`

Converts the provided *model_parameters* into a full parameters dict, taking into account all formatting options, and returns it.

This function can be used externally to check how the provided *model_parameters* would be interpreted when provided to the *ModelLink* subclass. Its output can be used for the ‘model_parameters’ input argument.

Parameters `model_parameters` (*array_like, dict or str*) – Anything that can be converted to a dict that provides model parameters information.

Returns `par_dict` (*dict*) – Dict with the provided *model_parameters* converted to its full format.

`prism.modellink.test_subclass` (*subclass, *args, **kwargs*)

Tests a provided *Modellink* subclass by initializing it with the given *args* and *kwargs* and checking if all required methods can be properly called.

This function needs to be called by all MPI ranks.

Parameters

- **subclass** (*Modellink* subclass) – The *Modellink* subclass that requires testing.
- **args** (*tuple*) – Positional arguments that need to be provided to the constructor of the *subclass*.
- **kwargs** (*dict*) – Keyword arguments that need to be provided to the constructor of the *subclass*.

Returns `modellink_obj` (*Modellink* object) – Instance of the provided *subclass* if all tests pass successfully. Specific exceptions are raised if a test fails.

Note: Depending on the complexity of the model wrapped in the given *subclass*, this function may take a while to execute.

Provides a collection of functions useful for using/mixing *PRISM* with other applications.

`prism.utils.get_lnpost_fn(ext_lnpost, pipeline_obj, *, emul_i=None, unit_space=True, hybrid=True)`

Returns a function definition `get_lnpost(par_set, *args, **kwargs)`.

This `get_lnpost` function can be used to calculate the natural logarithm of the posterior probability, which analyzes a given `par_set` first in the provided `pipeline_obj` at iteration `emul_i` and passes it to the `ext_lnpost` function if it is plausible.

This function needs to be called by all MPI ranks.

Parameters

- **ext_lnpost** (*function*) – Function definition that needs to be called if the provided `par_set` is plausible in iteration `emul_i` of `pipeline_obj`. The used call signature is `ext_lnpost(par_set, *args, **kwargs)`. All MPI ranks will call this function unless called within the `worker_mode` context manager.
- **pipeline_obj** (*Pipeline* object) – The instance of the `Pipeline` class that needs to be used for determining the validity of the proposed sampling step.

Other Parameters

- **emul_i** (*int or None. Default: None*) – If int, number indicating the requested emulator iteration. If `None`, the last iteration of the emulator will be used.
- **unit_space** (*bool. Default: True*) – Bool determining whether or not the provided sample will be given in unit space.
- **hybrid** (*bool. Default: True*) – Bool determining whether or not the `get_lnpost` function should use the implausibility values of a given `par_set` as an additional prior.

Returns `get_lnpost` (*function*) – Definition of the function `get_lnpost(par_set, *args, **kwargs)`.

See also:

`get_walkers()` Analyzes proposed `init_walkers` and returns valid `p0_walkers`.

`worker_mode` Special context manager within which all code is executed in worker mode.

Warning: Calling this function factory will disable all regular logging in `pipeline_obj` (`do_logging` set to `False`), in order to avoid having the same message being logged every time `get_lnpst` is called.

`prism.utils.get_walkers(pipeline_obj, *, emul_i=None, init_walkers=None, unit_space=True, ext_lnpst=None, **kwargs)`

Analyzes proposed `init_walkers` and returns valid `p0_walkers`.

Analyzes sample set `init_walkers` in the provided `pipeline_obj` at iteration `emul_i` and returns all samples that are plausible to be used as MCMC walkers. The provided samples and returned walkers should be/are given in unit space if `unit_space` is `True`.

If `init_walkers` is `None`, returns `impl_sam` instead if it is available.

This function needs to be called by all MPI ranks.

Parameters `pipeline_obj` (`Pipeline` object) – The instance of the `Pipeline` class that needs to be used for determining the validity of the proposed walkers.

Other Parameters

- **`emul_i`** (*int or None. Default: None*) – If int, number indicating the requested emulator iteration. If `None`, the last iteration of the emulator will be used.
- **`init_walkers`** (*2D array_like, int or None. Default: None*) – Sample set of proposed initial MCMC walker positions. All plausible samples in `init_walkers` will be returned. If int, generate an LHD of provided size and return all plausible samples. If `None`, return `impl_sam` corresponding to iteration `emul_i` instead.
- **`unit_space`** (*bool. Default: True*) – Bool determining whether or not the provided samples and returned walkers are given in unit space.
- **`ext_lnpst`** (*function or None. Default: None*) – If function, call `get_lnpst_fn()` function factory using `ext_lnpst` and the same values for `pipeline_obj`, `emul_i` and `unit_space`, and return the resulting function definition `get_lnpst`. Any additionally provided `kwargs` are also passed to it.

Returns

- **`n_walkers`** (*int*) – Number of returned MCMC walkers.
- **`p0_walkers`** (*2D ndarray object*) – Array containing starting positions of valid MCMC walkers.
- **`get_lnpst`** (*function (if `ext_lnpst` is a function)*) – The function returned by `get_lnpst_fn()` function factory using `ext_lnpst`, `pipeline_obj`, `emul_i`, `unit_space` and `kwargs` as the input values.

See also:

`get_lnpst_fn()` Returns a function definition `get_lnpst(par_set, *args, **kwargs)`.

`worker_mode` Special context manager within which all code is executed in worker mode.

Notes

If *init_walkers* is *None* and emulator iteration *emul_i* has not been analyzed yet, a *RequestError* will be raised.

Contains a collection of support classes/functions for the *PRISM* package.

exception `prism._internal.FeatureWarning`

Generic warning raised for experimental features in *PRISM*.

General purpose warning class, raised whenever a feature is used that should be considered experimental. Its behavior and API are subject to change, or the entire feature may be removed without a deprecation period.

exception `prism._internal.RequestError`

Generic exception raised for invalid action requests in the *PRISM* pipeline.

General purpose exception class, raised whenever a requested action cannot be executed due to it not being allowed or possible in the current state of the *Pipeline* instance.

exception `prism._internal.RequestWarning`

Generic warning raised for (future) action requests in the *PRISM* pipeline that may not be useful.

General purpose warning class, raised whenever a requested action may not produce appropriate or expected results due to the current state of the *Pipeline* instance. It is also raised if an obtained result can lead to such an action in the future.

class `prism._internal.CFilter (MPI_rank)`

Custom *Filter* class that only allows the controller rank to log messages to the logfile. Calls from worker ranks are ignored.

class `prism._internal.PRISM_Logger (*args, **kwargs)`

Special *Logger* class that allows for special filters to be set more easily.

class `prism._internal.RFilter (MPI_rank)`

Custom *Filter* class that prepends the rank of the MPI process that calls it to the logging message. If the size of the used MPI intra-communicator is 1, this filter does nothing.

`prism._internal.check_compatibility (emul_version)`

Checks if the provided *emul_version* is compatible with the current version of *PRISM*. Raises a *RequestError* if *False* and indicates which version of *PRISM* still supports the provided *emul_version*.

`prism._internal.check_vals (values, name, *args)`

Checks if all values in provided input argument *values* with *name* meet all criteria given in *args*. If no criteria

are given, it is checked if all values are finite. Returns *values* (0 or 1 in case of bool) if *True* and raises a *ValueError* or *TypeError* if *False*.

Parameters

- **values** (*array_like* of {*int*, *float*, *complex*, *str*, *bool*}) – The values to be checked against all given criteria in *args*. It must be possible to convert *values* to a *ndarray* object.
- **name** (*str*) – The name of the input argument, which is used in the error message if a criterion is not met.
- **args** (*tuple* of {'*bool*', '*complex*', '*float*', '*int*', '*neg*', '*nneg*', '*normal*', '*npos*', '*nzero*', '*pos*', '*str*'}) – Sequence of strings determining the criteria that *values* must meet. If *args* is empty, it is checked if *values* are finite.

Returns *return_values* (*array_like* of {*int*, *float*, *complex*, *str*}) – If *args* contained 'bool', returns 0s or 1s. Else, returns *values*.

Notes

If *values* contains integers, but *args* contains 'float', *return_values* will be cast as float.

`prism._internal.get_PRISM_File(prism_hdf5_file)`

Returns a class definition *PRISM_File*(*mode*, *emul_s*=None, **kwargs).

This class definition is a specialized version of the *File* class with the filename automatically set to *prism_hdf5_file* and added logging to the constructor and destructor methods.

Parameters *prism_hdf5_file* (*str*) – Absolute path to the master HDF5-file that is used in a *Pipeline* instance.

Returns *PRISM_File* (*class*) – Definition of the class *PRISM_File*(*mode*, *emul_s*=None, **kwargs).

`prism._internal.get_formatter()`

Returns a *Formatter* object containing the default logging formatting.

`prism._internal.get_handler(filename)`

Returns a *Handler* object containing the default logging handling settings.

`prism._internal.get_info()`

Prints a string that gives an overview of all information relevant to the PRISM package distribution.

`prism._internal.getCLogger(name=None)`

Creates a *PRISM_Logger* instance with *name*, adds the *CFilter* to it and returns it.

`prism._internal.getLogger(name=None, filters=None)`

Creates a *PRISM_Logger* instance with *name* and adds the provided *filters* to it. The returned *PRISM_Logger* instance is a child of the base *PRISM_Logger* created with *set_base_logger()*, but has its name changed (such that the parent name does not show up in the log-file).

Other Parameters

- **name** (*str* or *None*. Default: *None*) – The name of the *PRISM_Logger* instance to create. If *None*, initialize the base *PRISM_Logger* instead.
- **filters** (*list* of *str* or *None*. Default: *None*) – List of strings naming the filters that must be applied to the created *PRISM_Logger* instance. If *None* or the *PRISM_Logger* instance already existed, no filters will be applied.

Returns *logger* (*PRISM_Logger* object) – The created *PRISM_Logger* instance.

`prism._internal.getRLogger(name=None)`

Creates a *PRISM_Logger* instance with *name*, adds the *RFilter* to it and returns it.

`prism._internal.move_logger(working_dir)`

Moves the base *PRISM_Logger* from the current working directory to the given *working_dir*, and then restarts it again.

Parameters *working_dir* (*str*) – String containing the directory the log-file needs to be moved to.

`prism._internal.numpy_array(obj, *args, **kwargs)`

Returns `np.array(obj, *args, copy=False, **kwargs)`.

`prism._internal.set_base_logger(filename=None)`

Initializes the base *PRISM_Logger*, from which all other *PRISM_Logger* instances are derived.

Other Parameters *filename* (*str or None. Default: None*) – String containing the name of the log-file that is opened. If *None*, a new log-file will be created.

CHAPTER 14

Acknowledgements

Special thanks to Alan Duffy, Darren Croton, Simon Mutch and Manodeep Sinha for providing many valuable suggestions and constructive feedback points. Huge thanks to James Josephides for making the *PRISM* logo.

—
`prism.__internal`, [97](#)

u
`prism.utils`, [93](#)

Symbols

`_ModelLink__set_model_data()`
 (*prism.modellink.ModelLink* method), 86
`_ModelLink__set_model_parameters()`
 (*prism.modellink.ModelLink* method), 86
`_Projection__analyze_proj_hcube()`
 (*prism.Pipeline* method), 55
`_Projection__draw_2D_proj_fig()`
 (*prism.Pipeline* method), 55
`_Projection__draw_3D_proj_fig()`
 (*prism.Pipeline* method), 56
`_Projection__get_default_input_arguments()`
 (*prism.Pipeline* method), 56
`_Projection__get_default_parameters()`
 (*prism.Pipeline* method), 56
`_Projection__get_fig_path()` (*prism.Pipeline*
 method), 56
`_Projection__get_hcube_name()`
 (*prism.Pipeline* method), 56
`_Projection__get_proj_data()`
 (*prism.Pipeline* method), 56
`_Projection__get_proj_hcube()`
 (*prism.Pipeline* method), 57
`_Projection__get_req_hcubes()`
 (*prism.Pipeline* method), 57
`_Projection__prepare_projections()`
 (*prism.Pipeline* method), 57
`_Projection__process_input_arguments()`
 (*prism.Pipeline* method), 57
`_Projection__save_data()` (*prism.Pipeline*
 method), 58
`_Projection__set_parameters()`
 (*prism.Pipeline* method), 58
`__init__()` (*prism.Pipeline* method), 58
`__init__()` (*prism.emulator.Emulator* method), 73
`__init__()` (*prism.modellink.GaussianLink* method),
 85
`__init__()` (*prism.modellink.ModelLink* method), 87
`_assign_data_idx()` (*prism.emulator.Emulator*
 method), 73
`_assign_emul_s()` (*prism.emulator.Emulator*
 method), 74
`_call_model()` (*prism.Pipeline* method), 59
`_check_md_var()` (*prism.modellink.ModelLink*
 method), 87
`_check_mod_set()` (*prism.modellink.ModelLink*
 method), 88
`_check_sam_set()` (*prism.modellink.ModelLink*
 method), 88
`_cleanup_emul_files()`
 (*prism.emulator.Emulator* method), 75
`_compile_code_snippets()` (*prism.Pipeline*
 method), 59
`_construct_iteration()`
 (*prism.emulator.Emulator* method), 75
`_create_new_emulator()`
 (*prism.emulator.Emulator* method), 75
`_default_model_data` (*prism.modellink.ModelLink*
 attribute), 91
`_default_model_parameters`
 (*prism.modellink.ModelLink* attribute), 91
`_do_impl_check()` (*prism.Pipeline* method), 59
`_do_regression()` (*prism.emulator.Emulator*
 method), 75
`_evaluate()` (*prism.emulator.Emulator* method), 76
`_evaluate_model()` (*prism.Pipeline* method), 59
`_evaluate_sam_set()` (*prism.Pipeline* method), 59
`_get_active_par()` (*prism.emulator.Emulator*
 method), 76
`_get_adj_exp()` (*prism.emulator.Emulator* method),
 76
`_get_adj_var()` (*prism.emulator.Emulator* method),
 76
`_get_backup_path()` (*prism.modellink.ModelLink*
 method), 88
`_get_cov()` (*prism.emulator.Emulator* method), 77
`_get_cov_matrix()` (*prism.emulator.Emulator*
 method), 77
`_get_default_parameters()` (*prism.Pipeline*

`method`), 60
`_get_default_parameters()`
 (*prism.emulator.Emulator method*), 77
`_get_emul_i()` (*prism.emulator.Emulator method*),
 77
`_get_eval_sam_set()` (*prism.Pipeline method*), 60
`_get_exp_dot_term()` (*prism.emulator.Emulator*
 method), 78
`_get_ext_real_set()` (*prism.Pipeline method*), 60
`_get_inv_matrix()` (*prism.emulator.Emulator*
 method), 78
`_get_iteration_data()` (*prism.Pipeline method*),
 61
`_get_md_var()` (*prism.Pipeline method*), 61
`_get_mock_data()` (*prism.Pipeline method*), 61
`_get_model_par_seq()`
 (*prism.modellink.ModelLink method*), 88
`_get_n_eval_sam()` (*prism.Pipeline method*), 62
`_get_paths()` (*prism.Pipeline method*), 62
`_get_prior_exp()` (*prism.emulator.Emulator*
 method), 78
`_get_regr_cov()` (*prism.emulator.Emulator*
 method), 78
`_get_uni_impl()` (*prism.Pipeline method*), 62
`_listen_for_calls()` (*prism.Pipeline method*), 62
`_load_data()` (*prism.Pipeline method*), 62
`_load_data()` (*prism.emulator.Emulator method*), 78
`_load_emulator()` (*prism.emulator.Emulator*
 method), 79
`_make_backup()` (*prism.modellink.ModelLink*
 method), 89
`_make_call()` (*prism.Pipeline method*), 63
`_multi_call_model()` (*prism.Pipeline method*), 63
`_prepare_new_iteration()`
 (*prism.emulator.Emulator method*), 79
`_read_backup()` (*prism.modellink.ModelLink*
 method), 89
`_read_data_idx()` (*prism.emulator.Emulator*
 method), 79
`_read_parameters()` (*prism.Pipeline method*), 63
`_retrieve_parameters()`
 (*prism.emulator.Emulator method*), 79
`_save_data()` (*prism.Pipeline method*), 63
`_save_data()` (*prism.emulator.Emulator method*), 79
`_save_statistics()` (*prism.Pipeline method*), 64
`_set_impl_par()` (*prism.Pipeline method*), 64
`_set_mock_data()` (*prism.emulator.Emulator*
 method), 80
`_set_modellink()` (*prism.emulator.Emulator*
 method), 80
`_set_parameters()` (*prism.Pipeline method*), 64
`_set_parameters()` (*prism.emulator.Emulator*
 method), 80
`_to_par_space()` (*prism.modellink.ModelLink*
 method), 89
`_to_unit_space()` (*prism.modellink.ModelLink*
 method), 89
`_write_data_idx()` (*prism.emulator.Emulator*
 method), 80
 2+D model, 43
 2D model, 43

A

Active emulator system, 41
 Active parameters, 41
`active_emul_s` (*prism.emulator.Emulator attribute*),
 80
`active_par` (*prism.emulator.Emulator attribute*), 80
`active_par_data` (*prism.emulator.Emulator at-*
 tribute), 80
 Adjusted expectation, 41
 Adjusted values, 41
 Adjusted variance, 41
 Adjustment term, 41
 Analysis, 41
 Analyze, 41
`analyze()` (*prism.Pipeline method*), 64

B

`base_eval_sam` (*prism.Pipeline attribute*), 69
 BLA, 41

C

`call_model()` (*prism.modellink.ModelLink method*),
 89
`call_type` (*prism.modellink.ModelLink attribute*), 91
`ccheck` (*prism.emulator.Emulator attribute*), 80
`CFilter` (*class in prism._internal*), 99
`check_compatibility()` (*in module*
 prism._internal), 99
`check_vals()` (*in module prism._internal*), 99
`code_objects` (*prism.Pipeline attribute*), 69
`comm` (*prism.Pipeline attribute*), 69
 Construct, 41
`construct()` (*prism.Pipeline method*), 64
 Construction, 41
 Construction check, 41
 Controller, 41
 Controller rank, 42
`convert_data()` (*in module prism.modellink*), 92
`convert_parameters()` (*in module*
 prism.modellink), 92
`cov_mat_inv` (*prism.emulator.Emulator attribute*), 81
 Covariance matrix, 42
 Covariance vector, 42
`criterion` (*prism.Pipeline attribute*), 69
`cut_idx` (*prism.Pipeline attribute*), 70

D

Data error, [42](#)
 Data identifier, [42](#)
 Data point, [42](#)
 Data point identifier, [42](#)
 Data space, [42](#)
 Data value, [42](#)
 Data value space, [42](#)
 data_err (*prism.modellink.ModelLink* attribute), [91](#)
 data_idx (*prism.modellink.ModelLink* attribute), [91](#)
 data_idx_to_core (*prism.emulator.Emulator* attribute), [81](#)
 data_spc (*prism.modellink.ModelLink* attribute), [91](#)
 data_val (*prism.modellink.ModelLink* attribute), [91](#)
 details() (*prism.Pipeline* method), [65](#)
 do_active_anal (*prism.Pipeline* attribute), [70](#)
 do_logging (*prism.Pipeline* attribute), [70](#)

E

emul_i (*prism.emulator.Emulator* attribute), [81](#)
 emul_load (*prism.emulator.Emulator* attribute), [81](#)
 emul_s (*prism.emulator.Emulator* attribute), [81](#)
 emul_s_to_core (*prism.emulator.Emulator* attribute), [81](#)
 emul_type (*prism.emulator.Emulator* attribute), [81](#)
 Emulation method, [42](#)
 Emulator, [42](#)
 Emulator (*class in prism.emulator*), [73](#)
 emulator (*prism.Pipeline* attribute), [70](#)
 Emulator evaluation samples, [42](#)
 Emulator iteration, [42](#)
 Emulator system, [42](#)
 Emulator type, [42](#)
 Evaluate, [42](#)
 evaluate() (*prism.Pipeline* method), [66](#)
 Evaluation, [42](#)
 Evaluation set, [44](#)
 exp_dot_term (*prism.emulator.Emulator* attribute), [81](#)
 External model realization set, [42](#)

F

FeatureWarning, [99](#)
 File (*prism.Pipeline* attribute), [69](#)
 freeze_active_par (*prism.Pipeline* attribute), [70](#)
 Frozen active parameters, [42](#)
 Frozen parameters, [42](#)
 FSLR, [42](#)

G

Gaussian correlation length, [42](#)
 Gaussian sigma, [42](#)
 GaussianLink (*class in prism.modellink*), [85](#)

get_default_model_data() (*prism.modellink.ModelLink* method), [90](#)
 get_default_model_parameters() (*prism.modellink.ModelLink* method), [90](#)
 get_formatter() (*in module prism._internal*), [100](#)
 get_handler() (*in module prism._internal*), [100](#)
 get_info() (*in module prism._internal*), [100](#)
 get_lnpost_fn() (*in module prism.utils*), [95](#)
 get_md_var() (*prism.modellink.ModelLink* method), [90](#)
 get_PRISM_File() (*in module prism._internal*), [100](#)
 get_str_repr() (*prism.modellink.ModelLink* method), [90](#)
 get_walkers() (*in module prism.utils*), [96](#)
 getCLogger() (*in module prism._internal*), [100](#)
 getLogger() (*in module prism._internal*), [100](#)
 getRLogger() (*in module prism._internal*), [100](#)

H

HDF5, [42](#)
 hdf5_file (*prism.Pipeline* attribute), [70](#)
 Hybrid sampling, [43](#)

I

impl_cut (*prism.Pipeline* attribute), [70](#)
 impl_sam (*prism.Pipeline* attribute), [70](#)
 Implausibility check, [43](#)
 Implausibility cut-off check, [43](#)
 Implausibility cut-offs, [43](#)
 Implausibility value, [43](#)
 Implausibility wildcard, [43](#)
 Inverted covariance matrix, [42](#)
 is_controller (*prism.Pipeline* attribute), [70](#)
 is_worker (*prism.Pipeline* attribute), [70](#)
 Iteration, [42](#)

L

l_corr (*prism.emulator.Emulator* attribute), [81](#)
 LHD, [43](#)

M

Master file, [43](#)
 Master HDF5 file, [43](#)
 MCMC, [43](#)
 method (*prism.emulator.Emulator* attribute), [81](#)
 Mock data, [43](#)
 mod_set (*prism.emulator.Emulator* attribute), [81](#)
 Model, [43](#)
 Model discrepancy variance, [43](#)
 Model evaluation samples, [43](#)
 Model output, [43](#)
 Model outputs, [43](#)
 Model realization samples, [43](#)
 Model realization set, [44](#)

Model realizations, [43](#)
ModelLink, [43](#)
ModelLink (class in prism.modellink), [85](#)
modellink (prism.Pipeline attribute), [70](#)
ModelLink subclass, [43](#)
move_logger() (in module prism._internal), [101](#)
MPI, [44](#)
MPI rank, [44](#)
MPI_call (prism.modellink.ModelLink attribute), [91](#)
MSE, [44](#)
multi_call (prism.modellink.ModelLink attribute), [91](#)

N

n_cross_val (prism.emulator.Emulator attribute), [82](#)
n_data (prism.modellink.ModelLink attribute), [91](#)
n_emul_s (prism.emulator.Emulator attribute), [82](#)
n_emul_s_tot (prism.emulator.Emulator attribute), [82](#)
n_eval_sam (prism.Pipeline attribute), [70](#)
n_gaussians (prism.modellink.GaussianLink attribute), [85](#)
n_impl_sam (prism.Pipeline attribute), [71](#)
n_par (prism.modellink.ModelLink attribute), [91](#)
n_sam (prism.emulator.Emulator attribute), [82](#)
n_sam_init (prism.Pipeline attribute), [71](#)
name (prism.modellink.ModelLink attribute), [91](#)
nD model, [43](#)
np_array() (in module prism._internal), [101](#)

O

OLS, [44](#)

P

par_est (prism.modellink.ModelLink attribute), [92](#)
par_name (prism.modellink.ModelLink attribute), [92](#)
par_rng (prism.modellink.ModelLink attribute), [92](#)
Parameter set, [44](#)
Passive parameters, [44](#)
Pipeline, [44](#)
Pipeline (class in prism), [55](#)
Plausible region, [44](#)
Plausible samples, [44](#)
poly_coef (prism.emulator.Emulator attribute), [82](#)
poly_coef_cov (prism.emulator.Emulator attribute), [82](#)
poly_idx (prism.emulator.Emulator attribute), [82](#)
poly_order (prism.emulator.Emulator attribute), [82](#)
poly_powers (prism.emulator.Emulator attribute), [82](#)
Polynomial order, [44](#)
pot_active_par (prism.Pipeline attribute), [71](#)
Potentially active parameters, [44](#)
Prior covariance, [44](#)
Prior expectation, [44](#)
Prior variance, [44](#)

PRISM, [44](#)
PRISM Pipeline, [44](#)
prism._internal (module), [97](#)
prism.utils (module), [93](#)
prism_dict (prism.Pipeline attribute), [71](#)
PRISM_Logger (class in prism._internal), [99](#)
proj_depth (prism.Pipeline attribute), [71](#)
proj_res (prism.Pipeline attribute), [71](#)
Project, [44](#)
project() (prism.Pipeline method), [67](#)
Projection, [44](#)
Projection figure, [44](#)
Python Enhancement Proposals
PEP 377, [30](#)

R

rank (prism.Pipeline attribute), [71](#)
Regression, [44](#)
Regression covariances, [44](#)
RequestError, [99](#)
RequestWarning, [99](#)
Residual variance, [44](#)
RFilter (class in prism._internal), [99](#)
Root directory, [44](#)
root_dir (prism.Pipeline attribute), [71](#)
rsdl_var (prism.emulator.Emulator attribute), [82](#)
run() (prism.Pipeline method), [69](#)

S

sam_set (prism.emulator.Emulator attribute), [82](#)
Sample, [44](#)
Sample set, [44](#)
set_base_logger() (in module prism._internal), [101](#)
sigma (prism.emulator.Emulator attribute), [82](#)
SineWaveLink (class in prism.modellink), [92](#)
single_call (prism.modellink.ModelLink attribute), [92](#)
size (prism.Pipeline attribute), [71](#)

T

test_subclass() (in module prism.modellink), [93](#)

U

Univariate implausibility value, [43](#)
use_mock (prism.emulator.Emulator attribute), [83](#)
use_regr_cov (prism.emulator.Emulator attribute), [83](#)

W

Worker, [44](#)
Worker mode, [45](#)
Worker rank, [45](#)

`worker_mode` (*prism.Pipeline attribute*), [71](#)
`Working directory`, [45](#)
`working_dir` (*prism.Pipeline attribute*), [72](#)