
esem

Release v0.2.0

D Watson-Parris

Sep 06, 2021

CONTENTS:

1	Installing ESEm	1
1.1	Using PyPi	1
1.2	Using conda	1
1.3	Dependencies	1
2	What's new in ESEm	3
2.1	What's new in ESEm 1.1	3
3	Emulating with ESEm	5
3.1	Gaussian processes emulation	5
3.2	Neural network emulation	6
3.3	Random forest emulation	7
3.4	Data processing	7
3.5	Feature selection	7
4	Calibrating with ESEm	9
4.1	Using approximate Bayesian computation (ABC)	9
4.2	Using Markov-chain Monte-Carlo	10
5	Examples	13
5.1	Emulating using GPs	13
5.2	Emulating using CNNs	15
5.3	Random Forest Example: Cloud-resolving model sensitivity	22
5.4	Calibrating GPs using ABC	32
5.5	Calibrating GPs using MCMC	50
5.6	CMIP6 Emulation	57
5.7	Create paper emulation figure	68
6	API reference	77
6.1	Top-level functions	77
6.2	Emulator	79
6.3	Sampler	82
6.4	Wrappers	87
6.5	ModelAdaptor	90
6.6	DataProcessor	92
6.7	Utilities	95
7	ESEm design	99
7.1	Emulation	99
7.2	Calibration	99

8	Glossary	101
9	Indices and tables	103
	Index	105

INSTALLING ESEM

1.1 Using PyPi

It is straightforward to install esem using pip, this will automatically include tensorflow (with GPU support):

```
$ pip install esem
```

Optionally also install GPFlow, keras or scikit-learn

```
$ pip install esem[gpflow]
```

Or

```
$ pip install esem[gpflow,keras,scikit-learn]
```

1.2 Using conda

In order to make the most of the support for [Iris](#) and [CIS](#) creating a specific conda environment is recommended. If you don't already have conda, you must first download and install it. Anaconda is a free conda package that includes Python and many common scientific and data analysis libraries, and is available [here](#). Further documentation on using Anaconda and the features it provides can be found at <http://docs.continuum.io/anaconda/index.html>.

Having installed (mini-) conda - and ideally within a fresh environment - you can easily install CIS (and Iris) with the following command:

```
$ conda install -c conda-forge cis
```

It is then straightforward to install esem in to this environment using pip as above.

1.3 Dependencies

If you choose to install the dependencies yourself, use the following command to check the required dependencies are present:

```
$ python setup.py checkdep
```


WHAT'S NEW IN ESEM

2.1 What's new in ESEm 1.1

This page documents the new features added, and bugs fixed in ESEm since version 1.0. For more detail see all changes here: <https://github.com/duncanwp/ESEm/compare/1.0.0...1.1.0>

2.1.1 ESEm 1.1 features

- We have added this What's New page for tracking the latest developments in ESEm!
- We have dropped the mandatory requirement of Iris to make installation of ESEm easier. We have also added support for xarray DataArrays so that users can use their preferred library for data processing.
- The `esem.emulator.Emulator.predict()` and `esem.emulator.Emulator.batch_stats()` methods can now accept `pd.DataFrames` to match the training interface. The associated doc-strings and signatures have been extended to reflect this.

2.1.2 Bugs fixed

- Use `tqdm.auto` to automatically choose the appropriate progress bar for the context
- Fix `plot_validation` handling of masked data

EMULATING WITH ESEM

ESEm provides a simple and streamlined interface to emulate earth system datasets, denoted Y in the following documentation. These datasets can be provided as iris Cubes, xarray DataArrays or numpy arrays and the resulting emulation results will preserve any associated metadata. The corresponding predictors (X) can be provided as a numpy array or pandas *DataFrame*. This emulation is essentially just a regression estimating the functional form:

$$Y \approx f(X)$$

and can be performed using a variety of techniques using the same API.

3.1 Gaussian processes emulation

Gaussian processes (GPs) are a popular choice for model emulation due to their simple formulation and robust uncertainty estimates, particularly in cases where there is limited training data. Many excellent texts are available to describe their implementation and use (Rasmussen and Williams, 2005) and we only provide a short description here. Briefly, a GP is a stochastic process (a distribution of continuous functions) and can be thought of as an infinite dimensional normal distribution (hence the name).

The ESEm GP emulation module provides a thin wrapper around the [GPFlow](#) implementation. Please see their documentation for a detailed description of the implementation.

An important consideration when using GP regression is the form of the covariance matrix, or kernel. Typical kernels include: constant; linear; radial basis function (RBF; or squared exponential); and Matérn 3/2 and 5/2 which are only once and twice differentiable respectively. Kernels can also be designed to represent any aspect of the functions of interest such as non-stationarity or periodicity. This choice can often be informed by the physical setting and provides greater control and interpretability of the resulting model compared to e.g., Neural Networks.

Note: By default, ESEm uses a combination of linear, RBF and polynomial kernels which are suitable for the smooth and continuous parameter response expected for the examples used in this paper and related problems. However, given the importance of the kernel for determining the form of the functions generated by the GP we have also included the ability for users to specify combinations of other common kernels. See e.g., [Duvenaud, 2011](#) for a clear description of some common kernels and their combinations, as well as work towards automated methods for choosing them.

The framework provided by GPFlow also allows for multi-output GP regression and ESEm takes advantage of this to automatically provide regression over each of the output features provided in the training data. E.g. Y can be of arbitrary dimensionality. It will be automatically flattened and reshaped before being passed to GPFlow.

The most convenient way to setup a GPFlow emulator is using the `esem.gp_model()` function which can be imported directly:

```
from esim import gp_model
```

This creates a regression model with a default kernel as described above but provides a convenient interface for defining arbitrary kernels through addition and multiplication. For example, to initialize a model with a `Linear+Cosine()` kernel:

```
from esim import gp_model

# X_train and Y_train are our predictors and outputs, respectively.
model = gp_model(X_train, Y_train, kernel=['Linear', 'Cosine'], kernel_op='add')
```

Further details are described in the function description [`esim.gp_model\(\)`](#).

Examples of emulation using Gaussian processes can be found in [Emulating_using_GPs.ipynb](#) and [CMIP6_emulator.ipynb](#).

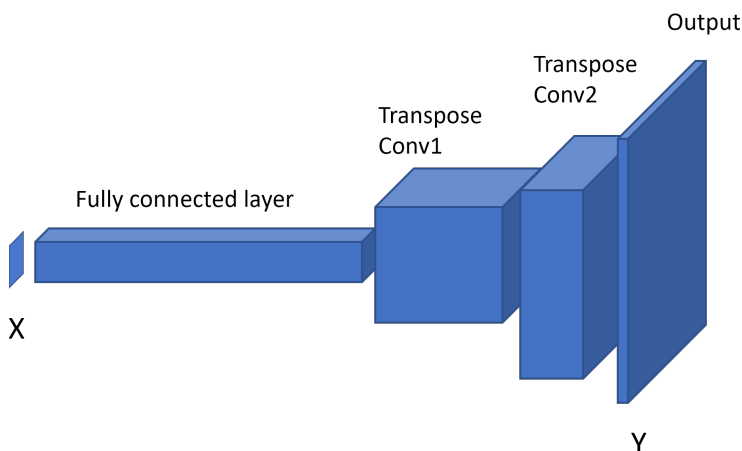
3.2 Neural network emulation

While fully connected neural networks have been used for many years, even in climate science, the recent surge in popularity has been powered by the increases in expressibility provided by deep, convolutional neural networks (CNNs) and the regularisation techniques which prevent these huge models from over-fitting the large amounts of training data required to train them.

Many excellent introductions can be found elsewhere but, briefly, a neural network consists of a network of nodes connecting (through a variety of architectures) the inputs to the target outputs via a series of weighted activation functions. The network architecture and activation functions are typically chosen a-priori and then the model weights are determined through a combination of back-propagation and (batch) gradient descent until the outputs match (defined by a given loss function) the provided training data. As previously discussed, the random dropping of nodes (by setting the weights to zero), termed dropout, can provide estimates of the prediction uncertainty of such networks.

The computational efficiency of such networks and the rich variety of architectures available have made them the tool of choice in many machine learning settings, and they are starting to be used in climate sciences for emulation, although the large amounts of training data required (especially compared to GPs) have so far limited their use somewhat.

ESEm provides a baseline CNN architecture based on the [Keras](#) library which essentially acts as a decoder - transforming input parameters into predicted 2(or 3) dimensional output fields:



This model can be easily constructed using the [`esim.cnn_model\(\)`](#) function. It is possible to use any Keras model in this way though and there are many potential ways of improving / developing this simple model.

An example of emulation using this convolution neural network can be found in [Emulating_using_ConvNets.ipynb](#).

3.3 Random forest emulation

ESEm also provides the option for emulation with Random Forests using the open-source implementation provided by scikit-learn. Random Forest estimators are comprised of an ensemble of decision trees; each decision tree is a recursive binary partition over the training data and the predictions are an average over the predictions of the decision trees.

As a result of this architecture, Random Forests (along with other algorithms built on decision trees) have two main attractions. Firstly, they require very little pre-processing of the inputs as the binary partitions are invariant to monotonic rescaling of the training data. Secondly, and of particular importance for climate problems, they are unable to extrapolate outside of their training data because the predictions are averages over subsets of the training dataset.

This model can be constructed using the `esem.rf_model()` function. All of the relevant scikit-learn arguments and keyword-arguments can be provided through this interface.

An example of emulation using the random forest can be found in [CRM_Emulation_with_RandomForest.ipynb](#).

3.4 Data processing

Many of the above approaches make assumptions about, or simply perform better when, the training data is structured or distributed in a certain way. These transformations are purely to help the emulator fit the training data, and can complicate comparison with e.g. observations during calibration. ESEm provides a simple and transparent way of transforming the datasets for training, and this automatically un-transforms the model predictions to aid in observational comparison.

Where these transformations are strictly necessary for a given model then it will be included in the wrapper function. Other choices are left to the user to apply as required.

For example, to ‘whiten’ the data (that is, remove the mean and normalise by the standard deviation):

```
import esim
from esim import gp_model

# X_train and Y_train are our predictors and outputs, respectively.
model = gp_model(X_train, Y_train, data_processors=[esim.data_processors.Whiten()])
```

A full list of the data processors can be found in the [API documentation](#).

3.5 Feature selection

ESEm includes a simple utility function that wraps the scikit-learn LassoLarsIC regression tool in order to enable an initial feature (parameter) selection. This can be useful to reduce the dimensionality of the input space. Either the Akaike information criterion (AIC) or the Bayes Information criterion (BIC) can be used, although BIC is the default.

For example,

```
from esim import gp_model
from esim.utils import get_param_mask

# X and Y are our model parameters and outputs respectively.
```

(continues on next page)

(continued from previous page)

```
active_params = get_param_mask(X, y)

# The model parameters can then be subsampled either directly
X_sub = X[:, active_params]

# Or by specifying the GP active_dims
active_dims, = np.where(active_params)
model = gp_model(X, y, active_dims=active_dims)
```

Note, this estimate only applies to one-dimensional outputs. Feature selection for higher dimension outputs is a much harder task beyond the scope of this package.

CALIBRATING WITH ESEM

Having trained a fast, robust emulator this can be used to calibrate our model against available observations. The following description closely follows that in our model description paper but with explicit links to the ESEM algorithms to aid clarity. Generally, this problem involves estimating the model parameters which could give rise to, or best match, the available observations.

In other words, we would like to know the posterior probability distribution of the input parameters: $p(\theta|Y^0)$.

Using Bayes' theorem, we can write this as:

$$p(\theta|Y^0) = p(Y^0|\theta) p(\theta) / p(Y^0) \quad (4.1)$$

Where the probability of an output given the input parameters, $p(Y^0|\theta)$, is referred to as the likelihood. While the model is capable of sampling this distribution, generally the full distribution is unknown and intractable, and we must approximate this likelihood. Depending on the purpose of the calibration and assumptions about the form of $p(Y^0|Y)$, different techniques can be used. In order to determine a (conservative) estimate of the parametric uncertainty in the model for example, we can use approximate Bayesian computation (ABC) to determine those parameters which are plausible given a set of observations. Alternatively, we may wish to know the optimal parameters to best match a set of observations and Markov-Chain Monte-Carlo based techniques might be more appropriate. Both of these sampling strategies are available in ESEM and we describe each of them here.

4.1 Using approximate Bayesian computation (ABC)

The simplest ABC approach seeks to approximate the likelihood using only samples from the simulator and a discrepancy function ρ :

$$p(\theta|Y^0) \propto p(Y^0|Y) p(Y|\theta) p(\theta) \approx \int \mathbb{I}(\rho(Y^0, Y) \leq \epsilon) p(Y|\theta) p(\theta) dY \quad (4.2)$$

where the indicator function $\mathbb{I}(x) = 1$ if x is true and $\mathbb{I}(x) = 0$ if x is false, and ϵ is a small discrepancy. This can then be integrated numerically using e.g., Monte-Carlo sampling of $p(\theta)$. Any of those parameters for which $\rho(Y^0, Y) \leq \epsilon$ are accepted and those which do not are rejected. As $\epsilon \rightarrow 0$ therefore, all parameters are accepted and we recover $p(\theta)$. For $\epsilon = 0$, it can be shown that we generate samples from the posterior $p(\theta|Y^0)$ exactly.

In practice however the simulator proposals will never exactly match the observations and we must make a pragmatic choice for both ρ and ϵ . ESEM includes an implementation of the 'implausibility metric' which defines the discrepancy in terms of the standardized Cartesian distance:

$$\rho(Y^0, Y(\theta)) = \frac{|Y^0 - Y(\theta)|}{\sqrt{\sigma_E^2 + \sigma_Y^2 + \sigma_R^2 + \sigma_S^2}} = \rho(Y^0, \theta) \quad (4.3)$$

where the total standard deviation is taken to be the squared sum of the emulator variance (σ_E^2) and the uncertainty in the observations (σ_Y^2) and due to representation (σ_R^2) and structural model uncertainties (σ_S^2) as described in the paper.

Framed in this way, ϵ , can be thought of as representing the number of standard deviations the (emulated) model value is from the observations. While this can be treated as a free parameter and may be specified in ESEm, it is common to choose $\epsilon = 3$ since it can be shown that for unimodal distributions values of 3σ correspond to a greater than 95% confidence bound. This approach is implemented in the `esem.abc_sampler.ABCSampler` class where ϵ is referred to as a threshold since it defines the cut-off for acceptance.

In the general case, multiple (\mathcal{N}) observations can be used and ρ can be written as a vector of implausibilities, $\rho(Y_i^O, \theta)$ or simply $\rho_i(\theta)$, and a modified method of rejection or acceptance must be used. A simple choice is to require $\rho_i < \epsilon \forall i \in \mathcal{N}$, however this can become restrictive for large \mathcal{N} due to the curse of dimensionality. The first step should be to reduce \mathcal{N} through the use of summary statistics, such as averaging over regions, or stations, or by performing an e.g. Principle Component Analysis (PCA) decomposition.

An alternative is to introduce a tolerance (T) such that only some proportion of ρ_i need to be smaller than ϵ : $\sum_{i=0}^{\mathcal{N}} H(\rho_i - \epsilon) < T$, where H is the Heaviside function. This tolerance can be specified when sampling using the `esem.abc_sampler.ABCSampler`. An efficient implementation of this approach whereby the acceptance is calculated in batches on the GPU can be particularly useful when dealing with high-dimensional outputs `esem.abc_sampler.ABCSampler`. It is recommended however, to choose $T = 0$ as a first approximation and then identify any particular observations which generate a very large implausibilities, since this provides a mechanism for identifying potential structural (or observational) errors.

A useful way of identifying such observations is using the `esem.abc_sampler.ABCSampler.get_implausibility()` method which returns the full implausibility matrix ρ_i . Note this may be very large (N-samples x N-observations) so it is recommended that only a subset of the full sample space be requested. The offending observations can then be removed and noted for further investigation.

Examples of the ABC sampling can be found in [Calibrating_GPs_using_ABC.ipynb](#).

4.2 Using Markov-chain Monte-Carlo

The ABC method described above is simple and powerful, but somewhat inefficient as it repeatedly samples from the same prior. In reality each rejection or acceptance of a set of parameters provides us with extra information about the ‘true’ form of $p(\theta|Y^0)$ so that the sampler could spend more time in plausible regions of the parameter space. This can then allow us to use smaller values of ϵ and hence find better approximations of $p(\theta|Y^0)$.

Given the joint probability distribution described by Eq. 2 and an initial choice of parameters θ' and (emulated) output Y' , the acceptance probability r of a new set of parameters (θ) is given by:

$$r = \frac{p(Y^0|Y') p(\theta'|\theta) p(\theta')}{p(Y^0|Y) p(\theta|\theta') p(\theta)} \quad (4.4)$$

The `esem.sampler.MCMCSampler` class uses the TensorFlow-probability implementation of Hamiltonian Monte-Carlo (HMC) which uses the gradient information automatically calculated by TensorFlow to inform the proposed new parameters θ . For simplicity, we assume that the proposal distribution is symmetric: $p(\theta'|\theta) = p(\theta|\theta')$, which is implemented as a zero log-acceptance correction in the initialisation of the TensorFlow target distribution. The target log probability provided to the TensorFlow HMC algorithm is then:

$$\log(r) = \log(p(Y^0|Y')) + \log(p(\theta')) - \log(p(Y^0|Y)) - \log(p(\theta)) \quad (4.5)$$

Note, that for this implementation the distance metric ρ must be cast as a probability distribution with values $[0, 1]$. We therefore assume that this discrepancy can be approximated as a normal distribution centred about zero, with standard deviation equal to the sum of the squares of the variances as described in Eq. 3:

$$p(Y^0|Y) \approx \frac{1}{\sigma_t \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{Y^0 - Y}{\sigma_t} \right)^2}, \quad \sigma_t = \sqrt{\sigma_E^2 + \sigma_Y^2 + \sigma_R^2 + \sigma_S^2} \quad (4.6)$$

The `esem.sampler.MCMCSampler.sample()` method will then return the requested number of accepted samples as well as reporting the acceptance rate, which provides a useful metric for tuning the algorithm. It should be noted that MCMC algorithms can be sensitive to a number of key parameters, including the number of burn-in steps used (and discarded) before sampling occurs and the step size. Each of these can be controlled via keyword arguments to the `esem.sampler.MCMCSampler.sample()` method.

This approach can provide much more efficient sampling of the emulator and provide improved parameter estimates, especially when used with informative priors which can guide the sampler.

Examples of the MCMC sampling can be found in [Calibrating_GPs_using_MCMC.ipynb](#) and [CMIP6_emulator.ipynb](#).

EXAMPLES

5.1 Emulating using GPs

```
[1]: import os
    ## Ignore my broken HDF5 install...
    os.putenv("HDF5_DISABLE_VERSION_CHECK", '1')
```

```
[2]: import iris
    from utils import get_bc_ppe_data

    from esem import gp_model
    from esem.utils import get_random_params

    import iris.quickplot as qplt
    import matplotlib.pyplot as plt
    %matplotlib inline
```

```
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\h5py\__init__.py:40:
↳ UserWarning: h5py is running against HDF5 1.10.6 when it was built against 1.10.5,
↳ this may cause problems
  '{0}.{1}.{2}'.format(*version.hdf5_built_version_tuple)
```

5.1.1 Read in the parameters and observables

```
[4]: ppe_params, ppe_aaod = get_bc_ppe_data()
```

```
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\__init__.py:249:
↳ IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and will
↳ be removed in a future release. Please remove code that sets this property.
  warn_deprecated(msg.format(name))
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\__init__.py:249:
↳ IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and will
↳ be removed in a future release. Please remove code that sets this property.
  warn_deprecated(msg.format(name))
```

```
[5]: n_test = 5

X_test, X_train = ppe_params[:n_test], ppe_params[n_test:]
Y_test, Y_train = ppe_aaod[:n_test,0], ppe_aaod[n_test:,0]
```

5.1.2 Setup and run the models

```
[9]: model = gp_model(X_train, Y_train)

[15]: model.train()

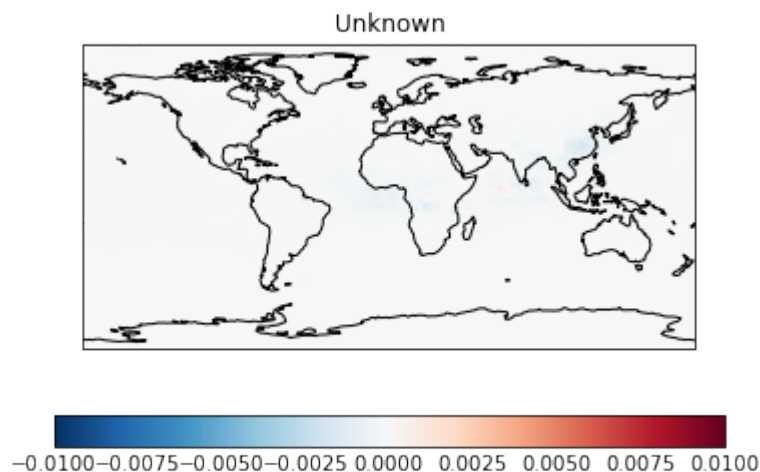
[16]: m, v = model.predict(X_test.values)

[17]: ## validation_plot(Y_test.data.flatten(), m.data.flatten(), v.data.flatten())

[18]: qplt.pcolormesh((m.collapsed('sample', iris.analysis.MEAN)-Y_test.collapsed('job', iris.
    ↳analysis.MEAN)), cmap='RdBu_r', vmin=-0.01, vmax=0.01)
plt.gca().coastlines()

C:\Users\duncan\miniconda3\envs\gcem_dev\lib\site-packages\iris\coords.py:1410:
↳UserWarning: Collapsing a non-contiguous coordinate. Metadata may not be fully
↳descriptive for 'sample'.
    warnings.warn(msg.format(self.name()))
C:\Users\duncan\miniconda3\envs\gcem_dev\lib\site-packages\iris\coords.py:1410:
↳UserWarning: Collapsing a non-contiguous coordinate. Metadata may not be fully
↳descriptive for 'job'.
    warnings.warn(msg.format(self.name()))
C:\Users\duncan\miniconda3\envs\gcem_dev\lib\site-packages\iris\coords.py:1193:
↳UserWarning: Coordinate 'longitude' is not bounded, guessing contiguous bounds.
    'contiguous bounds.'.format(self.name()))
C:\Users\duncan\miniconda3\envs\gcem_dev\lib\site-packages\iris\coords.py:1193:
↳UserWarning: Coordinate 'latitude' is not bounded, guessing contiguous bounds.
    'contiguous bounds.'.format(self.name()))

[18]: <cartopy.mpl.feature_artist.FeatureArtist at 0x19ea2c99388>
```



```
[19]: ## Note the model variance is constant across the outputs
v.data.max()

[19]: 1.175408691451335e-06
```

```
[15]: get_random_params(3, int(1e5)).shape
[15]: (100000, 3)

[16]: m, sd = model.batch_stats(get_random_params(3, int(1e3)))
100%|#####| 1000/1000 [00:09<00:00, 149.00sample/s]

[17]: m, sd = model.batch_stats(get_random_params(3, int(1e4)), batch_size=10)
100%|#####| 10000/10000 [00:07<00:00, 1459.07sample/s]

[19]: m, sd = model.batch_stats(get_random_params(3, int(1e6)), batch_size=10000)
100%|#####| 1000000/1000000 [00:09<00:00, 122569.61sample/s]

[ ]:
```

5.2 Emulating using CNNs

```
[2]: import os
    ## Ignore my broken HDF5 install...
    os.putenv("HDF5_DISABLE_VERSION_CHECK", '1')

[3]: import iris

    from utils import get_bc_ppe_data

    from esim import cnn_model
    from esim.utils import get_random_params

    import iris.quickplot as qplt
    import matplotlib.pyplot as plt
    %matplotlib inline

C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\h5py\__init__.py:40:
↳ UserWarning: h5py is running against HDF5 1.10.6 when it was built against 1.10.5,
↳ this may cause problems
    '{0}.{1}.{2}'.format(*version.hdf5_built_version_tuple)
```

5.2.1 Read in the parameters and data

```
[4]: ppe_params, ppe_aaod = get_bc_ppe_data()

C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\__init__.py:249:
↳ IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and will
↳ be removed in a future release. Please remove code that sets this property.
    warn_deprecated(msg.format(name))
```

(continues on next page)

(continued from previous page)

```
C:\Users\duncan\miniconda3\envs\gcem_dev\lib\site-packages\iris\__init__.py:249:
↳ IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and will
↳ be removed in a future release. Please remove code that sets this property.
warn_deprecated(msg.format(name))
```

```
[5]: ## Ensure the time dimension is last - this is treated as the color 'channel'
ppe_aaod.transpose((0,2,3,1))
```

```
[6]: n_test = 5

X_test, X_train = ppe_params[:n_test], ppe_params[n_test:]
Y_test, Y_train = ppe_aaod[:n_test], ppe_aaod[n_test:]
```

```
[7]: Y_train
```

```
[7]: <iris 'Cube' of Absorption optical thickness - total 550nm / (1) (job: 34; latitude: 96;
↳ longitude: 192; time: 12)>
```

5.2.2 Setup and run the models

```
[8]: model = cnn_model(X_train, Y_train)
```

```
[9]: model.train()

Epoch 1/100
4/4 [=====] - 2s 329ms/step - loss: 1.2647 - val_loss: 0.4618
Epoch 2/100
4/4 [=====] - 1s 172ms/step - loss: 1.1055 - val_loss: 0.4615
Epoch 3/100
4/4 [=====] - 1s 168ms/step - loss: 1.1964 - val_loss: 0.4610
Epoch 4/100
4/4 [=====] - 1s 171ms/step - loss: 1.0877 - val_loss: 0.4589
Epoch 5/100
4/4 [=====] - 1s 172ms/step - loss: 1.3223 - val_loss: 0.4561
Epoch 6/100
4/4 [=====] - 1s 170ms/step - loss: 1.0132 - val_loss: 0.4535
Epoch 7/100
4/4 [=====] - 1s 169ms/step - loss: 1.1708 - val_loss: 0.4544
Epoch 8/100
4/4 [=====] - 1s 177ms/step - loss: 1.2777 - val_loss: 0.4419
Epoch 9/100
4/4 [=====] - 1s 196ms/step - loss: 1.0134 - val_loss: 0.4406
Epoch 10/100
4/4 [=====] - 1s 197ms/step - loss: 1.0496 - val_loss: 0.4189
Epoch 11/100
4/4 [=====] - 1s 179ms/step - loss: 0.8041 - val_loss: 0.4329
Epoch 12/100
4/4 [=====] - 1s 170ms/step - loss: 1.1344 - val_loss: 0.4080
Epoch 13/100
4/4 [=====] - 1s 173ms/step - loss: 0.9210 - val_loss: 0.3832
```

(continues on next page)

(continued from previous page)

```

Epoch 14/100
4/4 [=====] - 1s 171ms/step - loss: 0.8290 - val_loss: 0.3700
Epoch 15/100
4/4 [=====] - 1s 179ms/step - loss: 0.9029 - val_loss: 0.3750
Epoch 16/100
4/4 [=====] - 1s 198ms/step - loss: 0.8040 - val_loss: 0.3580
Epoch 17/100
4/4 [=====] - 1s 169ms/step - loss: 0.9728 - val_loss: 0.3320
Epoch 18/100
4/4 [=====] - 1s 170ms/step - loss: 0.8510 - val_loss: 0.3260
Epoch 19/100
4/4 [=====] - 1s 171ms/step - loss: 0.6017 - val_loss: 0.3425
Epoch 20/100
4/4 [=====] - 1s 170ms/step - loss: 0.5955 - val_loss: 0.3306
Epoch 21/100
4/4 [=====] - 1s 174ms/step - loss: 0.7674 - val_loss: 0.2911
Epoch 22/100
4/4 [=====] - 1s 168ms/step - loss: 0.6860 - val_loss: 0.2662
Epoch 23/100
4/4 [=====] - 1s 171ms/step - loss: 0.7090 - val_loss: 0.2540
Epoch 24/100
4/4 [=====] - 1s 173ms/step - loss: 0.6631 - val_loss: 0.2368
Epoch 25/100
4/4 [=====] - 1s 171ms/step - loss: 0.6363 - val_loss: 0.2402
Epoch 26/100
4/4 [=====] - 1s 170ms/step - loss: 0.4646 - val_loss: 0.2249
Epoch 27/100
4/4 [=====] - 1s 172ms/step - loss: 0.6093 - val_loss: 0.2016
Epoch 28/100
4/4 [=====] - 1s 171ms/step - loss: 0.5923 - val_loss: 0.1915
Epoch 29/100
4/4 [=====] - 1s 168ms/step - loss: 0.4313 - val_loss: 0.1813
Epoch 30/100
4/4 [=====] - 1s 173ms/step - loss: 0.4832 - val_loss: 0.1717
Epoch 31/100
4/4 [=====] - 1s 176ms/step - loss: 0.5329 - val_loss: 0.1609
Epoch 32/100
4/4 [=====] - 1s 172ms/step - loss: 0.4968 - val_loss: 0.1572
Epoch 33/100
4/4 [=====] - 1s 174ms/step - loss: 0.5791 - val_loss: 0.1503
Epoch 34/100
4/4 [=====] - 1s 173ms/step - loss: 0.4020 - val_loss: 0.1417
Epoch 35/100
4/4 [=====] - 1s 170ms/step - loss: 0.5336 - val_loss: 0.1327
Epoch 36/100
4/4 [=====] - 1s 172ms/step - loss: 0.4808 - val_loss: 0.1267
Epoch 37/100
4/4 [=====] - 1s 171ms/step - loss: 0.4852 - val_loss: 0.1332
Epoch 38/100
4/4 [=====] - 1s 176ms/step - loss: 0.4243 - val_loss: 0.1177
Epoch 39/100
4/4 [=====] - 1s 172ms/step - loss: 0.3802 - val_loss: 0.1113

```

(continues on next page)

(continued from previous page)

```

Epoch 40/100
4/4 [=====] - 1s 173ms/step - loss: 0.5062 - val_loss: 0.1103
Epoch 41/100
4/4 [=====] - 1s 172ms/step - loss: 0.4019 - val_loss: 0.1086
Epoch 42/100
4/4 [=====] - 1s 175ms/step - loss: 0.4136 - val_loss: 0.1003
Epoch 43/100
4/4 [=====] - 1s 171ms/step - loss: 0.3220 - val_loss: 0.0984
Epoch 44/100
4/4 [=====] - 1s 172ms/step - loss: 0.3285 - val_loss: 0.0946
Epoch 45/100
4/4 [=====] - 1s 173ms/step - loss: 0.3216 - val_loss: 0.0845
Epoch 46/100
4/4 [=====] - 1s 180ms/step - loss: 0.2456 - val_loss: 0.0815
Epoch 47/100
4/4 [=====] - 1s 172ms/step - loss: 0.3025 - val_loss: 0.0778
Epoch 48/100
4/4 [=====] - 1s 175ms/step - loss: 0.2722 - val_loss: 0.0804
Epoch 49/100
4/4 [=====] - 1s 175ms/step - loss: 0.3915 - val_loss: 0.0739
Epoch 50/100
4/4 [=====] - 1s 182ms/step - loss: 0.4554 - val_loss: 0.0749
Epoch 51/100
4/4 [=====] - 1s 174ms/step - loss: 0.2702 - val_loss: 0.0712
Epoch 52/100
4/4 [=====] - 1s 169ms/step - loss: 0.3218 - val_loss: 0.0702
Epoch 53/100
4/4 [=====] - 1s 169ms/step - loss: 0.3559 - val_loss: 0.0676
Epoch 54/100
4/4 [=====] - 1s 188ms/step - loss: 0.3849 - val_loss: 0.0663
Epoch 55/100
4/4 [=====] - 1s 195ms/step - loss: 0.3528 - val_loss: 0.0663
Epoch 56/100
4/4 [=====] - 1s 169ms/step - loss: 0.4359 - val_loss: 0.0661
Epoch 57/100
4/4 [=====] - 1s 175ms/step - loss: 0.4033 - val_loss: 0.0661
Epoch 58/100
4/4 [=====] - 1s 172ms/step - loss: 0.3583 - val_loss: 0.0648
Epoch 59/100
4/4 [=====] - 1s 189ms/step - loss: 0.2553 - val_loss: 0.0664
Epoch 60/100
4/4 [=====] - 1s 176ms/step - loss: 0.2613 - val_loss: 0.0655
Epoch 61/100
4/4 [=====] - 1s 198ms/step - loss: 0.3042 - val_loss: 0.0647
Epoch 62/100
4/4 [=====] - 1s 198ms/step - loss: 0.3289 - val_loss: 0.0629
Epoch 63/100
4/4 [=====] - 1s 186ms/step - loss: 0.3522 - val_loss: 0.0680
Epoch 64/100
4/4 [=====] - 1s 174ms/step - loss: 0.3994 - val_loss: 0.0612
Epoch 65/100
4/4 [=====] - 1s 187ms/step - loss: 0.3320 - val_loss: 0.0731

```

(continues on next page)

(continued from previous page)

```

Epoch 66/100
4/4 [=====] - 1s 173ms/step - loss: 0.2637 - val_loss: 0.0632
Epoch 67/100
4/4 [=====] - 1s 183ms/step - loss: 0.2442 - val_loss: 0.0664
Epoch 68/100
4/4 [=====] - 1s 216ms/step - loss: 0.4290 - val_loss: 0.0621
Epoch 69/100
4/4 [=====] - 1s 186ms/step - loss: 0.3970 - val_loss: 0.0641
Epoch 70/100
4/4 [=====] - 1s 180ms/step - loss: 0.2437 - val_loss: 0.0613
Epoch 71/100
4/4 [=====] - 1s 182ms/step - loss: 0.3692 - val_loss: 0.0647
Epoch 72/100
4/4 [=====] - 1s 173ms/step - loss: 0.3234 - val_loss: 0.0618
Epoch 73/100
4/4 [=====] - 1s 179ms/step - loss: 0.3490 - val_loss: 0.0672
Epoch 74/100
4/4 [=====] - 1s 178ms/step - loss: 0.3524 - val_loss: 0.0611
Epoch 75/100
4/4 [=====] - 1s 173ms/step - loss: 0.3622 - val_loss: 0.0620
Epoch 76/100
4/4 [=====] - 1s 177ms/step - loss: 0.3948 - val_loss: 0.0658
Epoch 77/100
4/4 [=====] - 1s 170ms/step - loss: 0.2945 - val_loss: 0.0608
Epoch 78/100
4/4 [=====] - 1s 171ms/step - loss: 0.3269 - val_loss: 0.0638
Epoch 79/100
4/4 [=====] - 1s 171ms/step - loss: 0.3167 - val_loss: 0.0647
Epoch 80/100
4/4 [=====] - 1s 181ms/step - loss: 0.3143 - val_loss: 0.0646
Epoch 81/100
4/4 [=====] - 1s 191ms/step - loss: 0.3812 - val_loss: 0.0682
Epoch 82/100
4/4 [=====] - 1s 245ms/step - loss: 0.3734 - val_loss: 0.0680
Epoch 83/100
4/4 [=====] - 1s 246ms/step - loss: 0.3135 - val_loss: 0.0641
Epoch 84/100
4/4 [=====] - 1s 208ms/step - loss: 0.2190 - val_loss: 0.0662
Epoch 85/100
4/4 [=====] - 1s 173ms/step - loss: 0.3904 - val_loss: 0.0680
Epoch 86/100
4/4 [=====] - 1s 211ms/step - loss: 0.3203 - val_loss: 0.0693
Epoch 87/100
4/4 [=====] - 1s 263ms/step - loss: 0.3333 - val_loss: 0.0649
Epoch 88/100
4/4 [=====] - 1s 172ms/step - loss: 0.2221 - val_loss: 0.0668
Epoch 89/100
4/4 [=====] - 1s 214ms/step - loss: 0.3148 - val_loss: 0.0658
Epoch 90/100
4/4 [=====] - 1s 170ms/step - loss: 0.4059 - val_loss: 0.0717
Epoch 91/100
4/4 [=====] - 1s 170ms/step - loss: 0.2097 - val_loss: 0.0657

```

(continues on next page)

(continued from previous page)

```

Epoch 92/100
4/4 [=====] - 1s 170ms/step - loss: 0.2152 - val_loss: 0.0628
Epoch 93/100
4/4 [=====] - 1s 171ms/step - loss: 0.2845 - val_loss: 0.0644
Epoch 94/100
4/4 [=====] - 1s 168ms/step - loss: 0.2331 - val_loss: 0.0658
Epoch 95/100
4/4 [=====] - 1s 178ms/step - loss: 0.3202 - val_loss: 0.0653
Epoch 96/100
4/4 [=====] - 1s 176ms/step - loss: 0.1792 - val_loss: 0.0677
Epoch 97/100
4/4 [=====] - 1s 170ms/step - loss: 0.3031 - val_loss: 0.0691
Epoch 98/100
4/4 [=====] - 1s 174ms/step - loss: 0.2034 - val_loss: 0.0690
Epoch 99/100
4/4 [=====] - 1s 170ms/step - loss: 0.2300 - val_loss: 0.0703
Epoch 100/100
4/4 [=====] - 1s 169ms/step - loss: 0.2708 - val_loss: 0.0712

```

```
[14]: m, v = model.predict(X_test.to_numpy())
```

```

[15]: ## TODO: Tidy this up a bit
plt.figure(figsize=(12, 8))
plt.subplot(2,2,1)
qplt.pcolormesh(m[0].collapsed('time', iris.analysis.MEAN))
plt.gca().set_title('Predicted')
plt.gca().coastlines()

plt.subplot(2,2,2)
qplt.pcolormesh(Y_test[0].collapsed('time', iris.analysis.MEAN))
plt.gca().set_title('Test')
plt.gca().coastlines()

plt.subplot(2,2,3)
qplt.pcolormesh((m.collapsed(['sample', 'time'], iris.analysis.MEAN)-Y_test.collapsed([
    ↪ 'job', 'time'], iris.analysis.MEAN)), cmap='RdBu_r', vmin=-0.01, vmax=0.01)
plt.gca().coastlines()
plt.gca().set_title('Difference')

```

```

C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\coords.py:1193:
↪ UserWarning: Coordinate 'longitude' is not bounded, guessing contiguous bounds.
'contiguous bounds.'.format(self.name()))
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\coords.py:1193:
↪ UserWarning: Coordinate 'latitude' is not bounded, guessing contiguous bounds.
'contiguous bounds.'.format(self.name()))
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\coords.py:1193:
↪ UserWarning: Coordinate 'longitude' is not bounded, guessing contiguous bounds.
'contiguous bounds.'.format(self.name()))
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\coords.py:1193:
↪ UserWarning: Coordinate 'latitude' is not bounded, guessing contiguous bounds.
'contiguous bounds.'.format(self.name()))
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\coords.py:1410:
↪ UserWarning: Collapsing a non-contiguous coordinate. Metadata may not be fully
↪ descriptive for 'sample'.

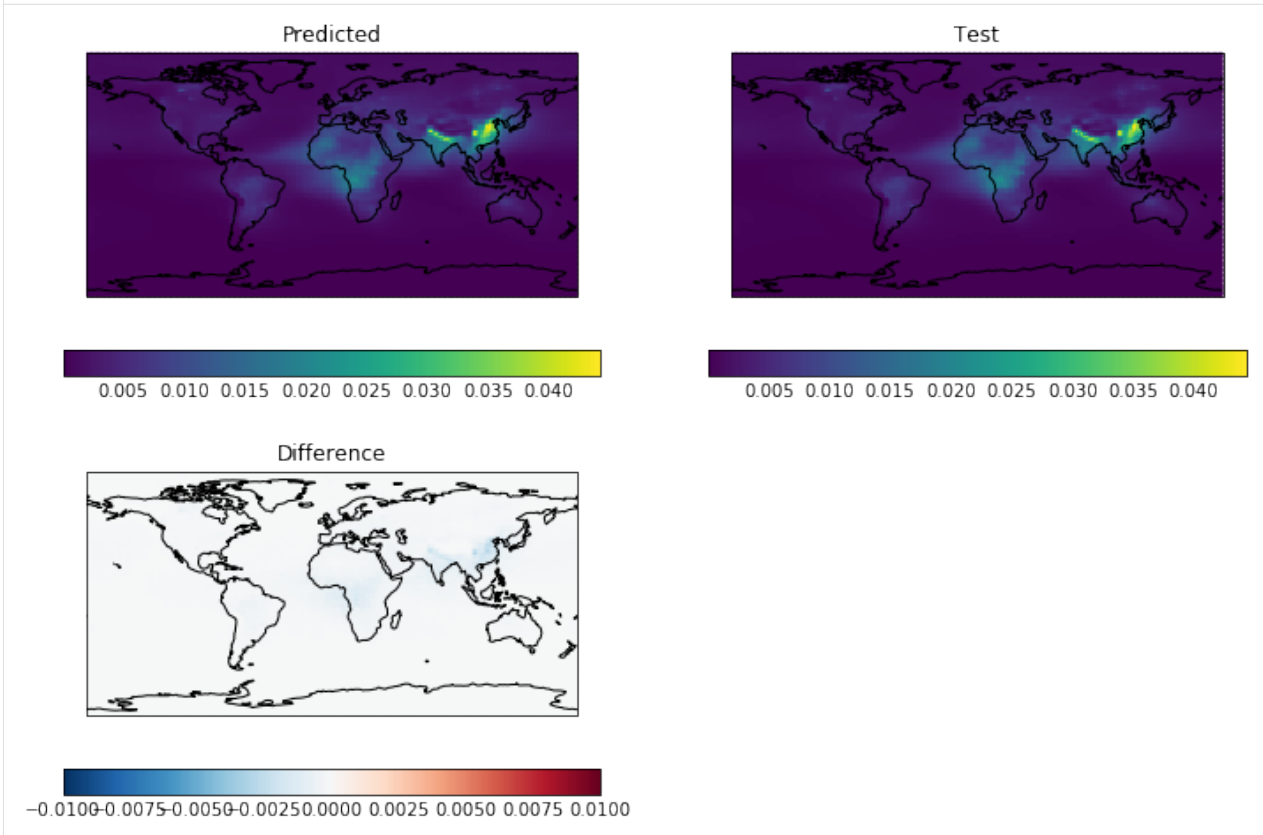
```

(continues on next page)

(continued from previous page)

```
warnings.warn(msg.format(self.name()))
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\coords.py:1410:
↳UserWarning: Collapsing a non-contiguous coordinate. Metadata may not be fully
↳descriptive for 'job'.
warnings.warn(msg.format(self.name()))
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\coords.py:1193:
↳UserWarning: Coordinate 'longitude' is not bounded, guessing contiguous bounds.
'contiguous bounds.'.format(self.name()))
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\coords.py:1193:
↳UserWarning: Coordinate 'latitude' is not bounded, guessing contiguous bounds.
'contiguous bounds.'.format(self.name()))
```

[15]: Text(0.5,1,'Difference')



[16]: m, sd = model.batch_stats(get_random_params(3, int(1e5)), batch_size=1000)

100%|#####| 100000/100000 [26:03<00:00, 44.36sample/s]

[]:

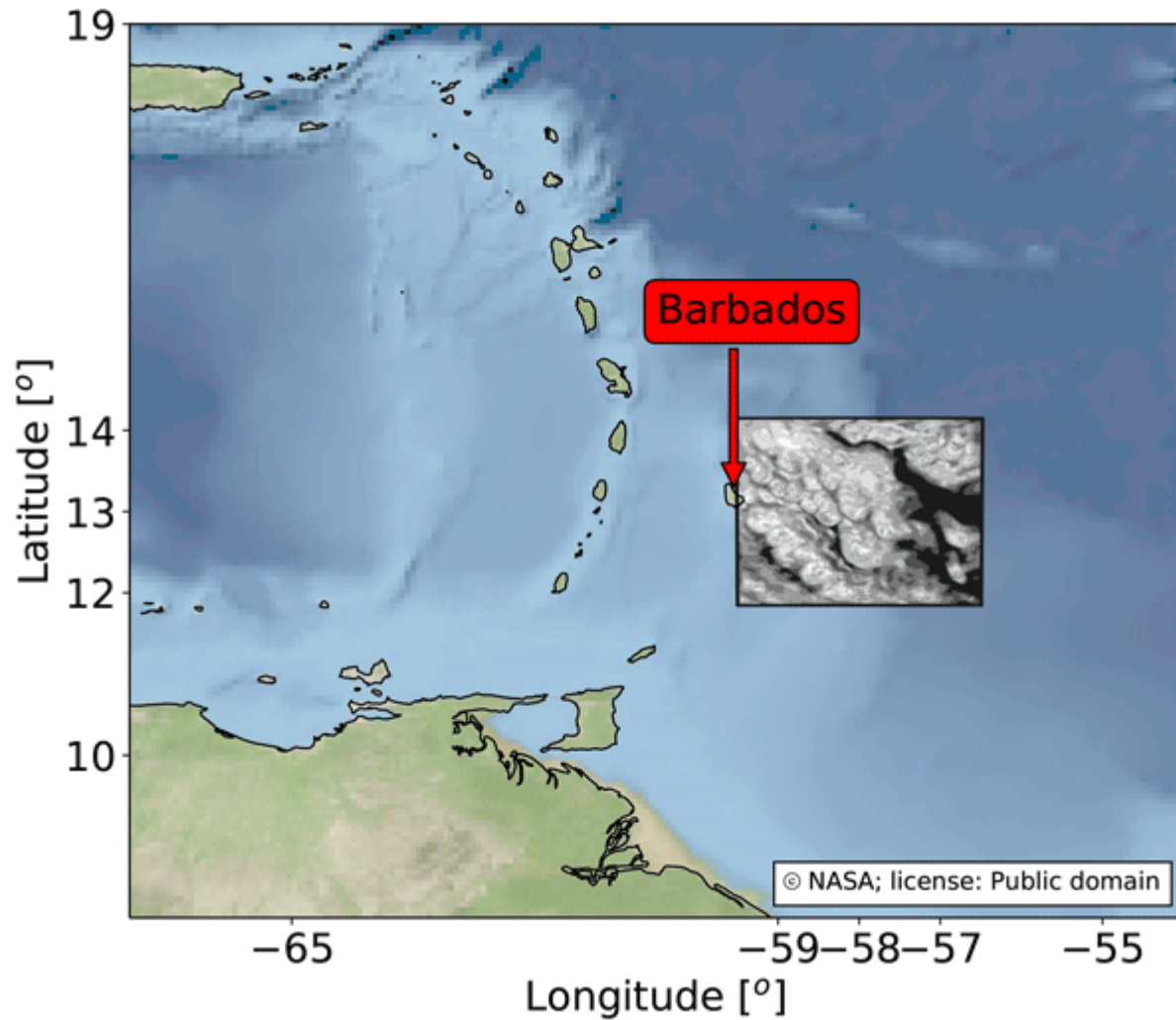
5.3 Random Forest Example: Cloud-resolving model sensitivity

In this example, we will use an ensemble of large-domain simulations of realistic shallow cloud fields to explore the sensitivity of shallow precipitation to local changes in the environment.

The simulation data we use for training the emulator is taken from a recent study Dagan and Stier (2020), where they performed ensemble daily simulations for one month-long period during December 2013 over the ocean to the East of Barbados, such that they sampled variability associated with shallow convection. Each day of the month consisted of two runs, both forced by realistic boundary conditions taken from reanalysis, but with different cloud droplet number concentrations (CDNC) to represent clean and polluted conditions. The altered CDNC was found to have little impact on the precipitation rate in the simulations, and so we simply treat the CDNC change as a perturbation to the initial conditions, and combine the two CDNC runs from each day together to increase the amount of data available for training the emulator. At hourly resolution, this provides us with 1488 data points.

However, given that the amount of precipitation is strongly tied to the local cloud regime, not fully controlling for cloud regime can introduce spurious correlations when training the emulator. As such we also filter out all hours which are not associated with shallow convective clouds. To do this, we consider domain-mean vertical profiles of total cloud water content (liquid + ice), q_t , and filter out all hours where the vertical sum of q_t below 600hPa exceeds 10^{-6} kg/kg. This condition allows us to filter out hours associated with the onset and development of deep convection in the domain, as well as masking out hours with high cirrus layers or hours dominated by transient mesoscale convective activity which is advected in by the boundary conditions. After this, we are left with 850 hourly data point which meet our criteria and can be used to train the emulator.

In this example we emulate hourly precip using domain-mean features: liquid water path (LWP), geopotential height at 700hPa (z_{700}), Estimated Inversion Strength (EIS), sea-surface temperature (SST) and the vertical pressure velocity at 700hPa (w_{700}).



References:

Dagan, G. and Stier, P.: Ensemble daily simulations for elucidating cloud–aerosol interactions under a large spread of realistic environmental conditions, *Atmos. Chem. Phys.*, 20, 6291–6303, <https://doi.org/10.5194/acp-20-6291-2020>, 2020.

```
[1]: import numpy as np
import pandas as pd
import iris

from utils import get_crm_data
from esim.utils import LeaveOneOut, plot_results, prettify_plot, add_121_line

from esim import rf_model

from matplotlib import pyplot as plt
plt.style.use('default')
%matplotlib inline
```

Concatenate 20cdnc and 200cdnc runs into one dataframe

```
[2]: df_crm = get_crm_data()
df_crm
```

```
[2]:      precip    pres_msl      LWP      WS10      lhfl_s      shfl_s  \
0      0.004593  101407.410  0.035898  6.639860 -167.53857  5.745860
1      0.006900  101356.266  0.044468  6.822748 -176.93939  4.438721
2      0.008916  101316.420  0.051559  6.798490 -182.61697  3.649221
3      0.008932  101270.490  0.057509  6.756970 -188.87599  3.033055
4      0.016204  101256.270  0.064226  6.763690 -194.85498  2.826119
..      ...      ...      ...      ...      ...      ...
845    0.063121  101309.750  0.064794  8.253145 -191.23718  12.219704
846    0.064601  101303.110  0.063914  8.326073 -192.57118  11.947702
847    0.046773  101332.234  0.059974  8.404624 -193.80084  12.372276
848    0.056623  101394.280  0.062895  8.385845 -192.18195  13.336615
849    0.064975  101438.690  0.069100  8.429897 -192.28928  13.679647

      LTS      w500      w600      w700      wmax850      wstd850      zg500  \
0    13.180252 -0.014463 -0.012311 -0.010275 -0.000024  0.000947  56627.516
1    13.279678 -0.015064 -0.012710 -0.008676  0.000030  0.000382  56572.645
2    13.333527 -0.014811 -0.012014 -0.006025  0.000642  0.000511  56525.613
3    13.328018 -0.013470 -0.012141 -0.004758  0.001519  0.000476  56471.332
4    13.317032 -0.010917 -0.011119 -0.003158  0.003252  0.000958  56443.758
..      ...      ...      ...      ...      ...      ...      ...
845   10.142059 -0.024480 -0.006400 -0.007968 -0.000044  0.001105  56084.273
846   10.162674 -0.019426  0.000300 -0.003904 -0.000034  0.000588  56071.547
847   10.166580 -0.014384  0.004355 -0.000284 -0.000251  0.000650  56079.492
848   10.149658 -0.016936  0.002702  0.000667  0.000013  0.000509  56117.140
849   10.164474 -0.021005  0.000413  0.000406 -0.000102  0.000838  56146.297

      zg600      zg700      rh850      rh700      u_shear      EIS  \
0    42694.220  30541.566  67.243774  60.067740 -4.662799  0.989443
1    42640.473  30488.172  69.299180  58.453730 -4.322696  1.130803
2    42593.594  30442.703  71.522900  56.912193 -3.925541  1.242463
3    42539.062  30390.662  74.115690  55.652990 -3.556972  1.304206
4    42510.805  30364.852  77.510765  54.434470 -3.319007  1.362710
..      ...      ...      ...      ...      ...      ...
845   42214.140  30222.945  83.696740  77.278465 -5.993636 -2.696190
846   42206.734  30214.715  84.196236  77.536760 -5.848422 -2.673406
847   42225.984  30236.312  84.394960  77.754560 -5.663757 -2.643809
848   42272.740  30286.500  84.437530  78.009740 -5.427930 -2.635981
849   42305.730  30322.684  84.389620  78.030650 -5.215088 -2.612224

      SST
0    301.173248
1    301.173248
2    301.173248
3    301.173248
4    301.173248
..      ...
845   300.126465
846   300.126465
847   300.126465
848   300.126465
```

(continues on next page)

(continued from previous page)

```
849 300.126465
```

```
[850 rows x 20 columns]
```

Extract the precipitation timeseries as target data

```
[3]: precip = df_crm['precip'].to_numpy().reshape(-1,1)
```

5.3.1 Visualize the precipitation landscape

In the ensemble, shallow precipitation is highly correlated with many different physical features. Most obviously there is a high correlation with liquid water path (LWP), 10-metre windspeed (WS10) and geopotential height at 700hPa (z_{700}).

We can use these correlations to create “collapsing spaces” for investigating the relationships between shallow precipitation and the local meteorological environment.

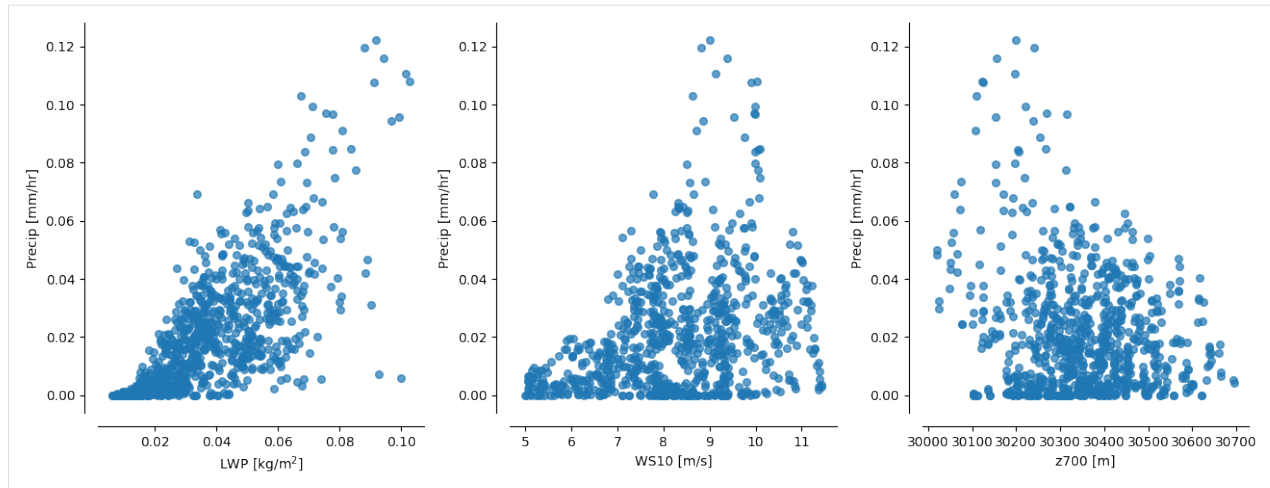
```
[4]: fig, axs = plt.subplots(ncols=3, figsize=(13,5), dpi=100, gridspec_kw={'width_ratios':[1,
    ↪ 1,1]})

axs[0].scatter(df_crm['LWP'], df_crm['precip'], s=30, marker='o', alpha=0.7, label=r'CDNC
    ↪ ${20}$ and CDNC${200}$')
axs[0].set_xlabel('LWP [kg/m^{2}]')
axs[0].set_ylabel('Precip [mm/hr]')
prettify_plot(axs[0])

axs[1].scatter(df_crm['WS10'], df_crm['precip'], s=30, marker='o', alpha=0.7, label=r
    ↪ 'CDNC${20}$ and CDNC${200}$')
axs[1].set_xlabel('WS10 [m/s]')
axs[1].set_ylabel('Precip [mm/hr]')
prettify_plot(axs[1])

axs[2].scatter(df_crm['zg700'], df_crm['precip'], s=30, marker='o', alpha=0.7, label=r
    ↪ 'CDNC${20}$ and CDNC${200}$')
axs[2].set_xlabel('z700 [m]')
axs[2].set_ylabel('Precip [mm/hr]')
prettify_plot(axs[2])

fig.tight_layout()
plt.savefig("Figs/1hr_D13shal_lwp-ws-pr.png", dpi=200)
```



Also, good to note that each of these predictors ``LWP, WS10, z700`` are mutually uncorrelated (see plots below)

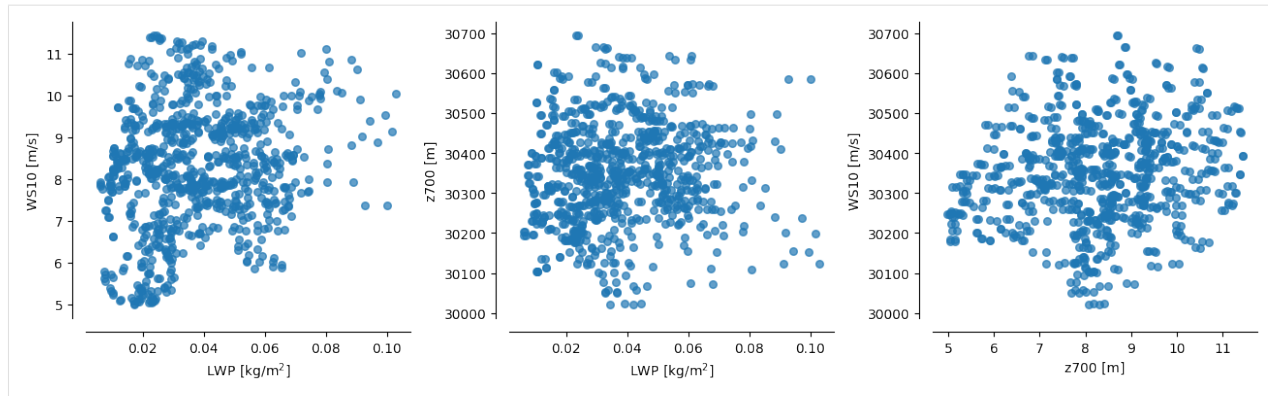
```
[5]: fig, axs = plt.subplots(ncols=3, figsize=(13,4), dpi=100, gridspec_kw={'width_ratios':[1,
    ↪ 1,1]})

axs[0].scatter(df_crm['LWP'], df_crm['WS10'], s=30, marker='o', alpha=0.7, label=r'CDNC$_{20}$ and CDNC$_{200}$')
axs[0].set_xlabel('LWP [kg/m$^{2}$]')
axs[0].set_ylabel('WS10 [m/s]')
prettify_plot(axs[0])

axs[1].scatter(df_crm['LWP'], df_crm['zg700'], s=30, marker='o', alpha=0.7, label=r'CDNC$_{20}$ and CDNC$_{200}$')
axs[1].set_xlabel('LWP [kg/m$^{2}$]')
axs[1].set_ylabel('z700 [m]')
prettify_plot(axs[1])

axs[2].scatter(df_crm['WS10'], df_crm['zg700'], s=30, marker='o', alpha=0.7, label=r'CDNC$_{20}$ and CDNC$_{200}$')
axs[2].set_xlabel('z700 [m]')
axs[2].set_ylabel('WS10 [m/s]')
prettify_plot(axs[2])

fig.tight_layout()
#plt.savefig("Figs/1hr_D13shal_lwp-ws-pr.png", dpi=200)
```



Stratifying precip by pairs of these predictors creates nice “collapsing spaces”

Nice for illustrating how the emulated surface compares to the raw data

```
[6]: fig, axs = plt.subplots(ncols=3, figsize=(13,5), dpi=100, gridspec_kw={'width_ratios':[1,
    ↪1,0.05]}))

sc1 = axs[0].scatter(df_crm['LWP'], df_crm['zg700'], c=df_crm['precip'],
    vmin=0, vmax=0.12, s=30, marker='o', alpha=0.7, label=r'CDNC$_{20}$ and
    ↪CDNC$_{200}$')

axs[0].set_xlabel(r'LWP [kg m$^{-2}$]')
axs[0].set_ylabel(r'z$_{700}$ [m]')

axs[0].set_title("Hourly output: Dec2013, shallow clouds")
axs[0].legend(loc='lower right')

prettify_plot(axs[0])

sc2 = axs[1].scatter(df_crm['LWP'], df_crm['WS10'], c=df_crm['precip'],
    vmin=0, vmax=0.12, s=30, marker='o', alpha=0.7, label=r'CDNC$_{20}$ and
    ↪CDNC$_{200}$')

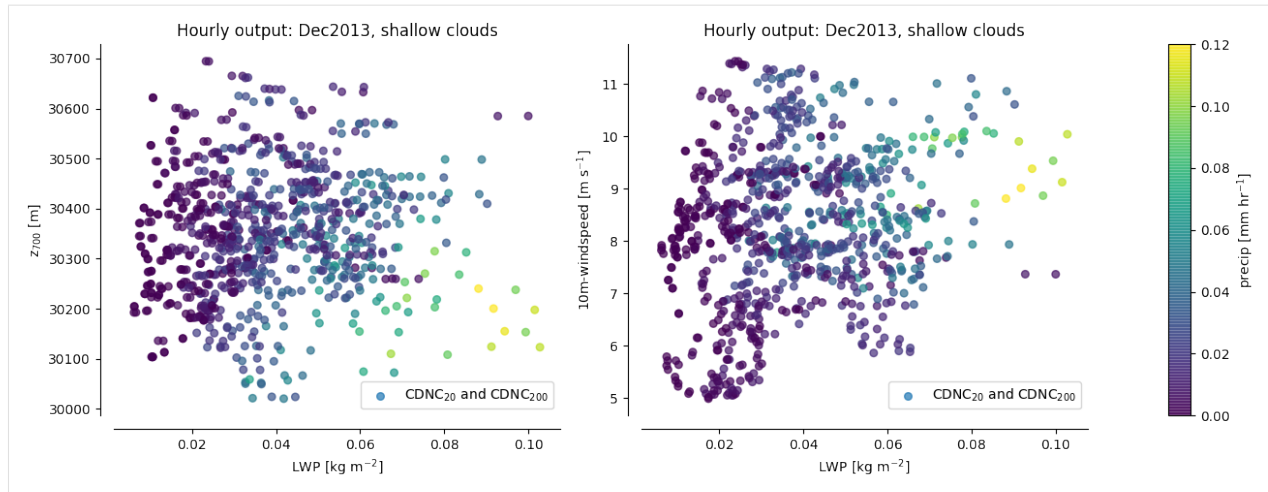
axs[1].set_xlabel(r'LWP [kg m$^{-2}$]')
axs[1].set_ylabel(r'10m-windspeed [m s$^{-1}$]')

axs[1].set_title("Hourly output: Dec2013, shallow clouds")
axs[1].legend(loc='lower right')

prettify_plot(axs[1])

fig.colorbar(sc1, cax=axs[2], label=r'precip [mm hr$^{-1}$]')

fig.tight_layout()
plt.savefig("Figs/1hr_D13shal_lwp-ws-pr.png", dpi=200)
```

5.3.2 Emulation

Our aim is to emulate shallow precipitation as a function of the environmental conditions, and then plot the predictions in LWP–z700 space to compare with the scatter points above.

To do this we choose a set of predictors which are typical “cloud-controlling factors” such as SST, Estimated Inversion Strength, vertical velocity at 700 hPa, LWP and z700. Other variables could also be chosen and it’s worth exploring this just to get a sense for how the model behaves.

After validating the model using Leave-One-Out cross-validation, we then retrain the model using the full dataset, and use this model to predict the precipitation across a wide range of values. Finally, for the purpose of plotting in LWP–z700 space, we reduce the dimensionality of our final prediction by averaging over all features with aren’t LWP or z700. This gives us a smooth field to compare with the scatter points.

```
[4]: params = df_crm.loc[:, ['LWP', 'zg700', 'EIS', 'SST', 'w700']]
```

```
print("The input params are: \n", params, "\n")
```

The input params are:

	LWP	zg700	EIS	SST	w700
0	0.035898	30541.566	0.989443	301.173248	-0.010275
1	0.044468	30488.172	1.130803	301.173248	-0.008676
2	0.051559	30442.703	1.242463	301.173248	-0.006025
3	0.057509	30390.662	1.304206	301.173248	-0.004758
4	0.064226	30364.852	1.362710	301.173248	-0.003158
...
845	0.064794	30222.945	-2.696190	300.126465	-0.007968
846	0.063914	30214.715	-2.673406	300.126465	-0.003904
847	0.059974	30236.312	-2.643809	300.126465	-0.000284
848	0.062895	30286.500	-2.635981	300.126465	0.000667
849	0.069100	30322.684	-2.612224	300.126465	0.000406

[850 rows x 5 columns]

LeaveOneOut cross-validation and plotting

```
[11]: %%time

# Ignore the mountain of warnings
import warnings
from sklearn.exceptions import DataConversionWarning
warnings.filterwarnings(action='ignore', category=DataConversionWarning)

outp = LeaveOneOut(Xdata=params, Ydata=precip, model='RandomForest', n_estimators=50,
↳random_state=0)

truth_RF, pred_RF = outp[:,0], outp[:,1]

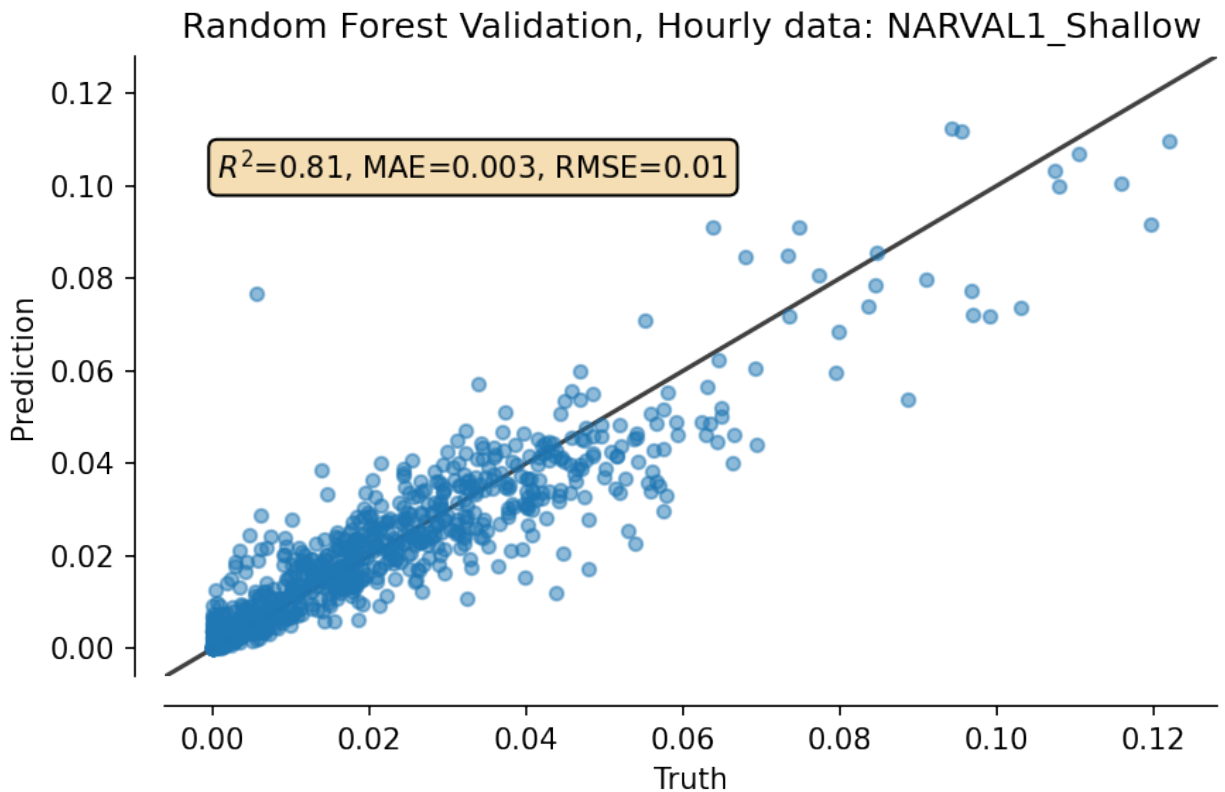
CPU times: user 1min 52s, sys: 1.91 s, total: 1min 53s
Wall time: 1min 55s
```

```
[11]: from sklearn.metrics import mean_squared_error

""" Validation plot """
fig, ax = plt.subplots(dpi=150)

plot_results(ax, truth_RF, pred_RF, title="Random Forest Validation, Hourly data:
↳NARVAL1_Shallow")

fig.tight_layout()
```



Now, retrain model on all data, and extrapolate over whole space

```
[8]: X_train = params.to_numpy()
      Y_train = precip.ravel()
      model = rf_model(X_train, Y_train)
```

```
[9]: model.train()

/Users/andrewwilliams/Desktop/PhD/Projects/GCEm/GCEm/rf_model.py:29:
↳DataConversionWarning: A column-vector y was passed when a 1d array was expected.
↳Please change the shape of y to (n_samples,), for example using ravel().
      self.model.fit(X=self.training_params, y=Y_flat)
```

```
[21]: %%time

      # Now, make grid for plotting RF predictions
      # more n_points means higher resolution, but takes exponentially longer
      n_points = 30

      min_vals = params.min()
      max_vals = params.max()

      # For uniform prediction over full params space
      space=np.linspace(min_vals, max_vals, n_points)

      # Reshape to (N,D)
      reshape_to_ND = np.transpose(space)
      Xs_uniform = np.meshgrid(*reshape_to_ND)
      test = np.array([_.flatten() for _ in Xs_uniform]).T

      # Predict
      predictions,_ = model.predict(test)
      predictions = predictions.reshape(Xs_uniform[0].shape)

      # Now, take mean over all parameters except [LWP, zg700], assumed to be first 2 indices
      predictions_reduced = np.mean(predictions, axis=tuple(range(2, predictions.ndim)))

      CPU times: user 1min 2s, sys: 17.5 s, total: 1min 19s
      Wall time: 1min 27s
```

```
[23]: # Now, make grid for plotting RF predictions

      LWP_grid = np.linspace(min_vals['LWP'], max_vals['LWP'], num=n_points)
      zg700_grid = np.linspace(min_vals['zg700'], max_vals['zg700'], num=n_points)

      lwp, zg = np.meshgrid(LWP_grid, zg700_grid)
```

5.3.3 Plot!

```
[24]: fig, ax = plt.subplots(ncols=3, figsize=(7,4), dpi=250, gridspec_kw={'width_ratios':[1,1,1]
↳0.05, 0.05}})
fig.suptitle("Hourly output: Dec2013, shallow clouds")

cp = ax[0].pcolormesh(LWP_grid, zg700_grid,
                    predictions_reduced, vmin=0, vmax=0.12, alpha=1)

fig.colorbar(cp, cax=ax[1], orientation='vertical', shrink=0.05, label=r'Precip [mm hr$^{\rightarrow{-1}}$]')

"""Overlap errors"""
ax[0].scatter(df_crm['LWP'], df_crm['zg700'], c=df_crm['precip'],
             vmin=0, vmax=0.12, s=30, marker='o', edgecolors="None")

ers = ax[0].scatter(df_crm['LWP'], df_crm['zg700'], c=(truth_RF-pred_RF)/(truth_RF+pred_
↳RF), facecolors="None",
                 vmin=-1, vmax=1, s=30, marker='o', lw=0.7, alpha=0.5, cmap='seismic')
ers.set_facecolors("None")
fig.colorbar(ers, cax=ax[2], label=r'$\frac{\mathrm{Truth-Prediction}}{\mathrm{
↳{Truth+Prediction}}}$')

for idx, _ in enumerate(ax[:1]):
    _ .set_xlabel(r'LWP [kg m$^{\rightarrow{-2}}$]')
    _ .set_ylabel(r'z$_{700}$ [m]')

    if idx==0:
        _ .set_xlim(min_vals['LWP'], max_vals['LWP'])
        _ .set_ylim(min_vals['zg700'], max_vals['zg700'])

fig.tight_layout()
fig.subplots_adjust(top=0.85) # Put this AFTER tight_layout() call !

"""Add validation plot as inset to first axis"""
axins = ax[0].inset_axes([0.79, 0.79, 0.2, 0.2], # x0, y0, width, height
                        transform=ax[0].transAxes)
axins.scatter(truth_RF, pred_RF, s=2, alpha=0.3)
add_121_line(axins)
axins.set_xlabel('Model', position=(0.5,0), fontsize=8, labelpad=-0.01)
axins.set_xticks([0, 0.05, 0.1])
axins.set_xticklabels(labels=[0, 0.05, 0.1], fontdict={'fontsize':6})

axins.set_ylabel('Emulator', position=(0,0.5), fontsize=8, labelpad=-0.01)
axins.set_yticks([0, 0.05, 0.1])
axins.set_yticklabels(labels=[0, 0.05, 0.1], fontdict={'fontsize':6})

axins.tick_params(axis='both', which='major', pad=0.01)

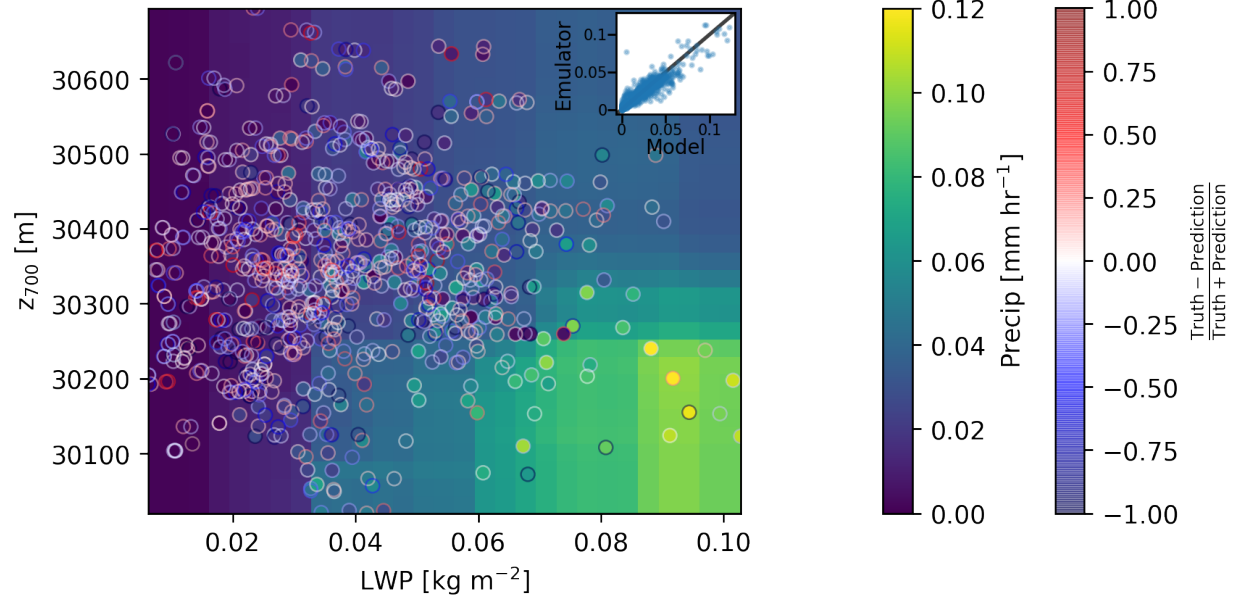
plt.savefig("./Figs/1hr_D13shal_lwp-zg700-pr_w-errorpoints.png", dpi=300)
```

```
<ipython-input-24-ba7a7a607fa9>:5: MatplotlibDeprecationWarning: shading='flat' when X_
↳and Y have the same dimensions as C is deprecated since 3.3. Either specify the_
↳corners of the quadrilaterals with X and Y, or pass shading='auto', 'nearest' or
↳'gouraud', or set rcParams['pcolor.shading']. This will become an error two minor_
↳releases later.
```

(continued from previous page)

```
cp = ax[0].pcolormesh(LWP_grid, zg700_grid,
```

Hourly output: Dec2013, shallow clouds



```
[ ]:
```

5.4 Calibrating GPs using ABC

```
[1]: import os
      # Ignore my broken HDF5 install...
      os.putenv("HDF5_DISABLE_VERSION_CHECK", '1')

[2]: import pandas as pd
      import cis
      import iris

      from utils import get_aeronet_data, get_bc_ppe_data

      from esim.utils import validation_plot, plot_parameter_space, get_random_params,
      ↪ ensemble_collocate
      from esim import gp_model
      from esim.abc_sampler import ABCSampler, constrain

      import iris.quickplot as qplt
      import matplotlib.pyplot as plt
      %matplotlib inline
```

```
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\h5py\__init__.py:40:
↳ UserWarning: h5py is running against HDF5 1.10.6 when it was built against 1.10.5,
↳ this may cause problems
  '{0}.{1}.{2}'.format(*version.hdf5_built_version_tuple)
```

5.4.1 Read in the parameters and observables

```
[3]: aaod = get_aeronet_data()
print(aaod)

Ungridded data: Absorption_AOD440nm / (1)
  Shape = (10098,)

  Total number of points = 10098
  Number of non-masked points = 10098
  Long name = Absorption_AOD440nm
  Standard name = None
  Units = 1
  Missing value = -999.0
  Range = (5.1e-05, 0.47236)
  History =
  Coordinates:
    longitude
      Long name =
      Standard name = longitude
      Units = degrees_east
      Missing value = None
      Range = (-155.576755, 141.3407)
      History =
    latitude
      Long name =
      Standard name = latitude
      Units = degrees_north
      Missing value = None
      Range = (-35.495807, 79.990278)
      History =
    time
      Long name =
      Standard name = time
      Units = days since 1600-01-01 00:00:00
      Missing value = None
      Range = (cftime.DatetimeGregorian(2017, 1, 1, 12, 0, 0, 0), cftime.
↳ DatetimeGregorian(2018, 1, 2, 12, 0, 0, 0))
      History =
```

```
[4]: # Read in the PPE parameters, AAOD and DRE
ppe_params, ppe_aaod, ppe_dre = get_bc_ppe_data(dre=True)
```

```
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\__init__.py:249:
↳ IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and will
↳ be removed in a future release. Please remove code that sets this property.
```

(continues on next page)

(continued from previous page)

```
warn_deprecated(msg.format(name))
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\__init__.py:249:
↳ IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and will
↳ be removed in a future release. Please remove code that sets this property.
warn_deprecated(msg.format(name))
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\__init__.py:249:
↳ IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and will
↳ be removed in a future release. Please remove code that sets this property.
warn_deprecated(msg.format(name))
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\iris\__init__.py:249:
↳ IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and will
↳ be removed in a future release. Please remove code that sets this property.
warn_deprecated(msg.format(name))
```

```
[5]: # Take the annual mean of the DRE
ppe_dre, = ppe_dre.collapsed('time', iris.analysis.MEAN)
```

```
WARNING:root:Creating guessed bounds as none exist in file
WARNING:root:Creating guessed bounds as none exist in file
WARNING:root:Creating guessed bounds as none exist in file
WARNING:root:Creating guessed bounds as none exist in file
```

5.4.2 Collocate the model on to the observations

```
[6]: col_ppe_aaod = ensemble_collocate(ppe_aaod, aaod)
```

```
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳ FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳ use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳ an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳ different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳ FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳ use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳ an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳ different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳ FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳ use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳ an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳ different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳ FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳ use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳ an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳ different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳ FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳ use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳ an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳ different result.
34dout = self.data[indx]
```

(continues on next page)

```
deut = self.data[indx]
```

→ an array index, `arr[np.array(seq)]`, which will result either in an error or a (continues on next page)

(continued from previous page)

```
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
```

(Continued on next page)

→ an array index, `arr[np.array(seq)]`, which will result either in an error or a (continues on next page)

(continued from previous page)

```
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\ma\core.py:3225:
↳FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated;
↳use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as
↳an array index, `arr[np.array(seq)]`, which will result either in an error or a
↳different result.
dout = self.data[indx]
```

```
[7]: n_test = 8
```

```
X_test, X_train = ppe_params[:n_test], ppe_params[n_test:]
Y_test, Y_train = col_ppe_aod[:n_test], col_ppe_aod[n_test:]
```

```
[8]: Y_train
```

```
[8]: <iris 'Cube' of Absorption optical thickness - total 550nm / (1) (job: 31; obs: 10098)>
```

5.4.3 Setup and run the models

Explore different model choices

```
[9]: from esim.utils import leave_one_out, prediction_within_ci
from scipy import stats
import numpy as np

from esim.data_processors import Log
res_l = leave_one_out(X_train, Y_train, model='GaussianProcess', data_
    ↳ processors=[Log(constant=0.1)], kernel=['Linear', 'Exponential', 'Bias'])
r2_values_l = [stats.linregress(x.data.compressed(), y.data[:, ~x.data.mask].
    ↳ flatten())[2]**2 for x,y,_ in res_l]
ci95_values_l = [prediction_within_ci(x.data.flatten(), y.data.flatten(), v.data.
    ↳ flatten())[2].sum()/x.data.count() for x,y,v in res_l]
print("Mean R^2: {:.2f}".format(np.asarray(r2_values_l).mean()))
print("Mean proportion within 95% CI: {:.2f}".format(np.asarray(ci95_values_l).mean()))

res = leave_one_out(X_train, Y_train, model='GaussianProcess', kernel=['Linear', 'Bias'])
r2_values = [stats.linregress(x.data.flatten(), y.data.flatten())[2]**2 for x,y,v in res]
ci95_values = [prediction_within_ci(x.data.flatten(), y.data.flatten(), v.data.
    ↳ flatten())[2].sum()/x.data.count() for x,y,v in res]
print("Mean R^2: {:.2f}".format(np.asarray(r2_values).mean()))
print("Mean proportion within 95% CI: {:.2f}".format(np.asarray(ci95_values).mean()))

# Note that while the Log pre-processing leads to slightly better R^2, the model is
    ↳ under-confident and
# has too large uncertainties which would adversely affect our implausibility metric.

c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid
    ↳ value encountered in log
    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid
    ↳ value encountered in log
    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid
    ↳ value encountered in log
    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid
    ↳ value encountered in log
    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid
    ↳ value encountered in log
    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid
    ↳ value encountered in log
    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid
    ↳ value encountered in log
    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid
    ↳ value encountered in log
    return np.log(data + self.constant)
```

(continues on next page)

(continued from previous page)

[illegible]

(continues on next page)

(continued from previous page)

```

    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid_
↪value encountered in log
    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid_
↪value encountered in log
    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid_
↪value encountered in log
    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid_
↪value encountered in log
    return np.log(data + self.constant)
c:\users\duncan\pycharmprojects\gcm\esem\data_processors.py:67: RuntimeWarning: invalid_
↪value encountered in log
    return np.log(data + self.constant)

```

```

Mean R^2: 1.00
Mean proportion within 95% CI: 1.00
Mean R^2: 0.99
Mean proportion within 95% CI: 0.94

```

Build final model

```
[10]: model = gp_model(X_train, Y_train, kernel=['Linear', 'Bias'])
```

```
[11]: model.train()
```

```
[12]: m, v = model.predict(X_test.values)
```

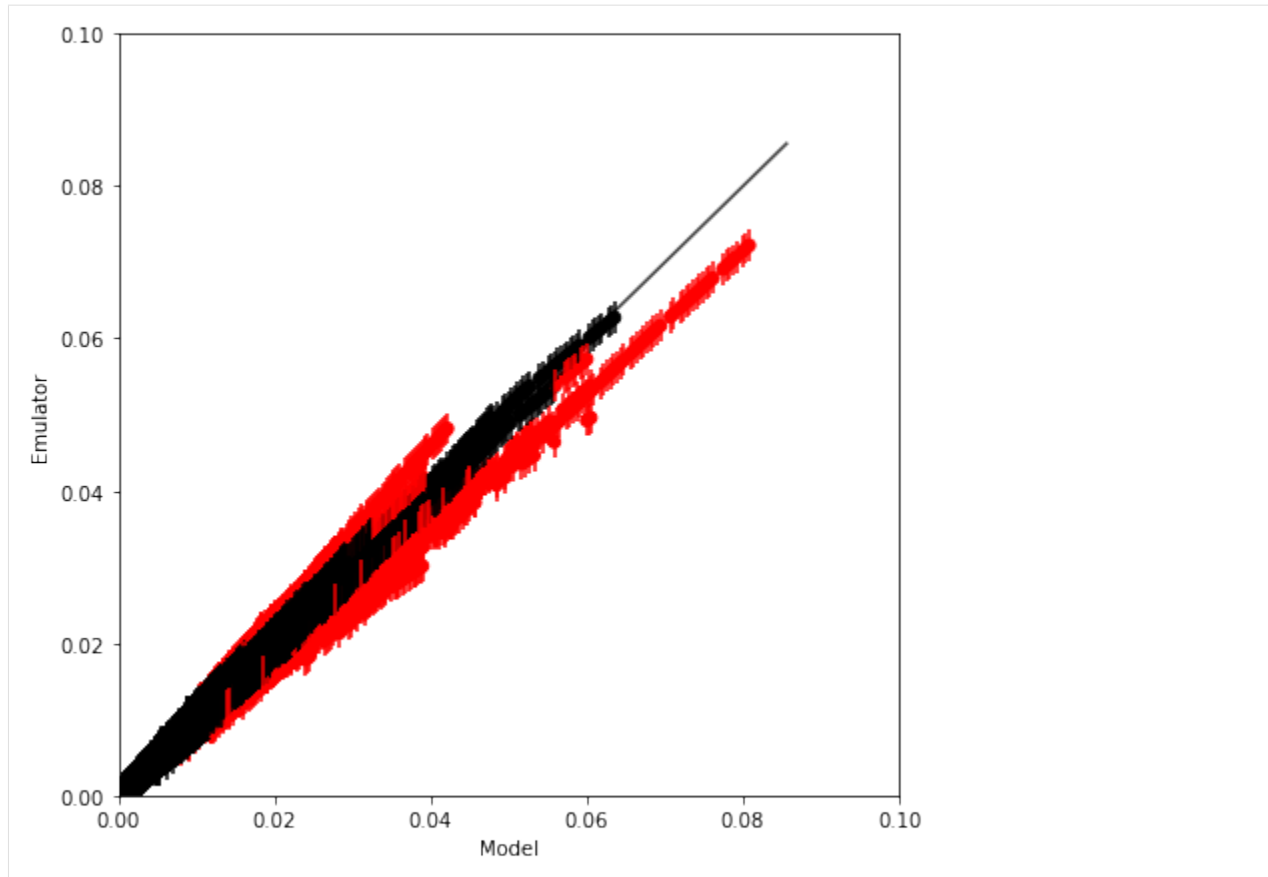
```
[13]: validation_plot(Y_test.data.flatten(), m.data.flatten(), v.data.flatten(),
                    minx=0, maxx=0.1, miny=0., maxy=0.1)
```

```
Proportion of 'Bad' estimates : 5.52%
```

```

C:\Users\duncan\miniconda3\envs\gcm_dev\lib\site-packages\numpy\core\_asarray.py:83:
↪UserWarning: Warning: converting a masked element to nan.
    return array(a, dtype, copy=False, order=order)

```



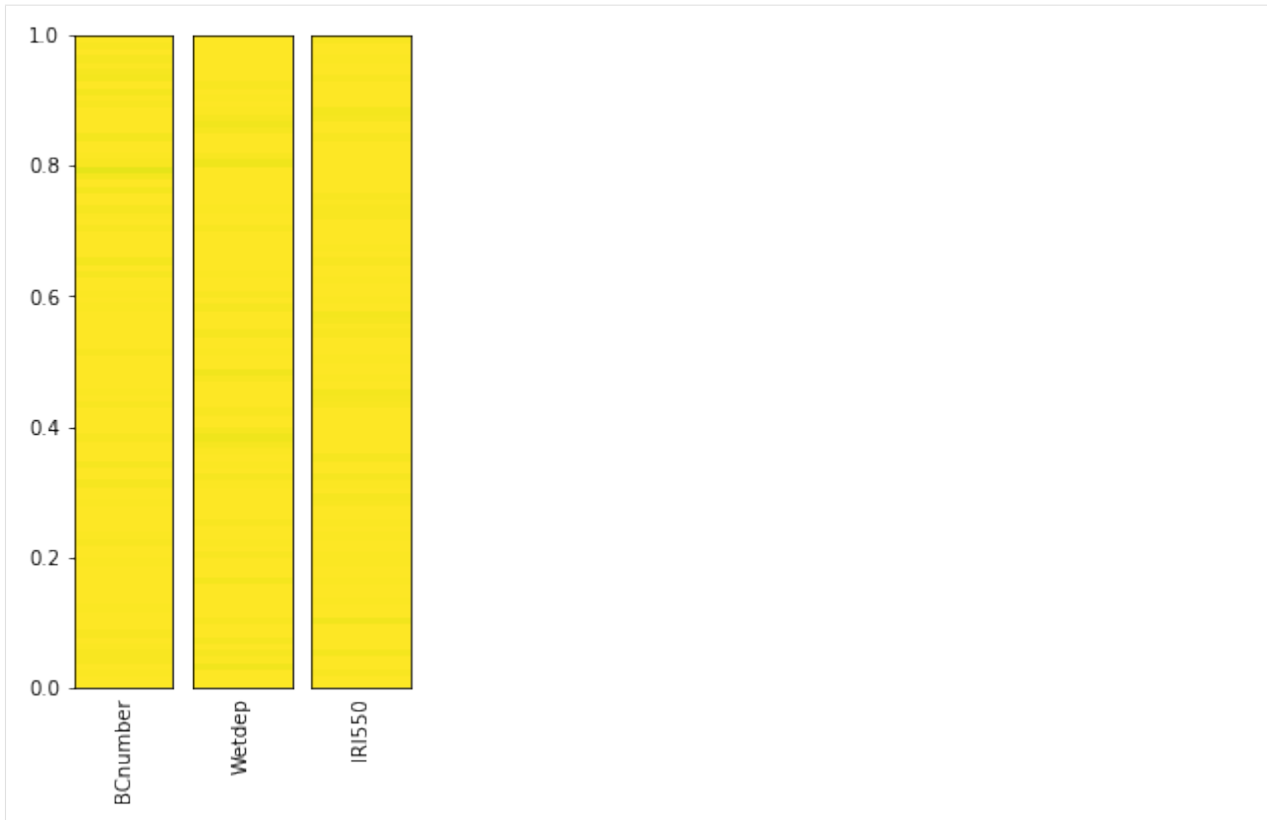
5.4.4 Sample and constrain the models

Emulating $1e6$ sample points directly would require 673 Gb of memory so we can either run $1e6$ samples for each point, or run the constraint everywhere, but in batches. Here we do the latter, optionally on the GPU, using the ‘naive’ algorithm for calculating the running mean and variance of the various properties.

The rejection sampling happens in a similar manner so that only as much memory as is used for one batch is ever used.

```
[14]: # In this case
sample_points = pd.DataFrame(data=get_random_params(3, int(1e6)), columns=X_train.
    ↪ columns)

[15]: # Note that smoothing the parameter distribution can be slow for large numbers of points
plot_parameter_space(sample_points, fig_size=(3,6), smooth=False)
```



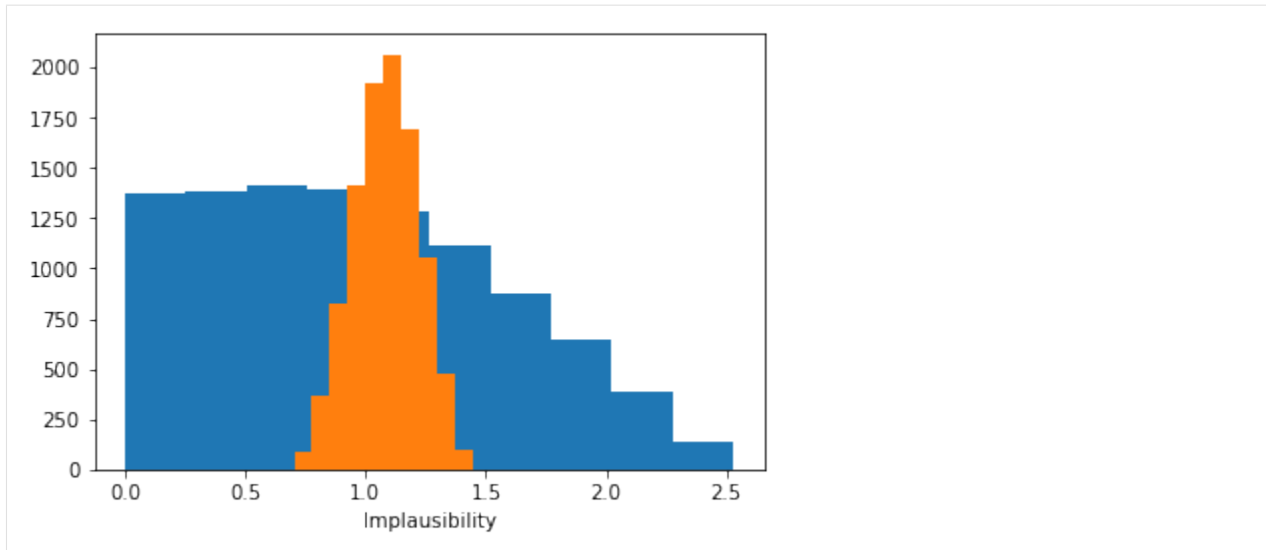
```
[16]: # Setup the sampler to compare against our AeroNet data
sampler = ABCSampler(model, aaod, obs_uncertainty=0.5, repres_uncertainty=0.5)

[17]: # Calculate the implausibility for each sample against each observation - note this can
      ↪ be very large so we only sample a fraction!
implaus = sampler.get_implausibility(sample_points[:,100], batch_size=1000)

# The implausibility distributions for different observations can be very different.
_ = plt.hist(implaus.data[:, 1400])
_ = plt.hist(implaus.data[:, 14])
plt.gca().set(xlabel='Implausibility')

100%|#####| 10000/10000 [00:05<00:00, 1859.15sample/s]

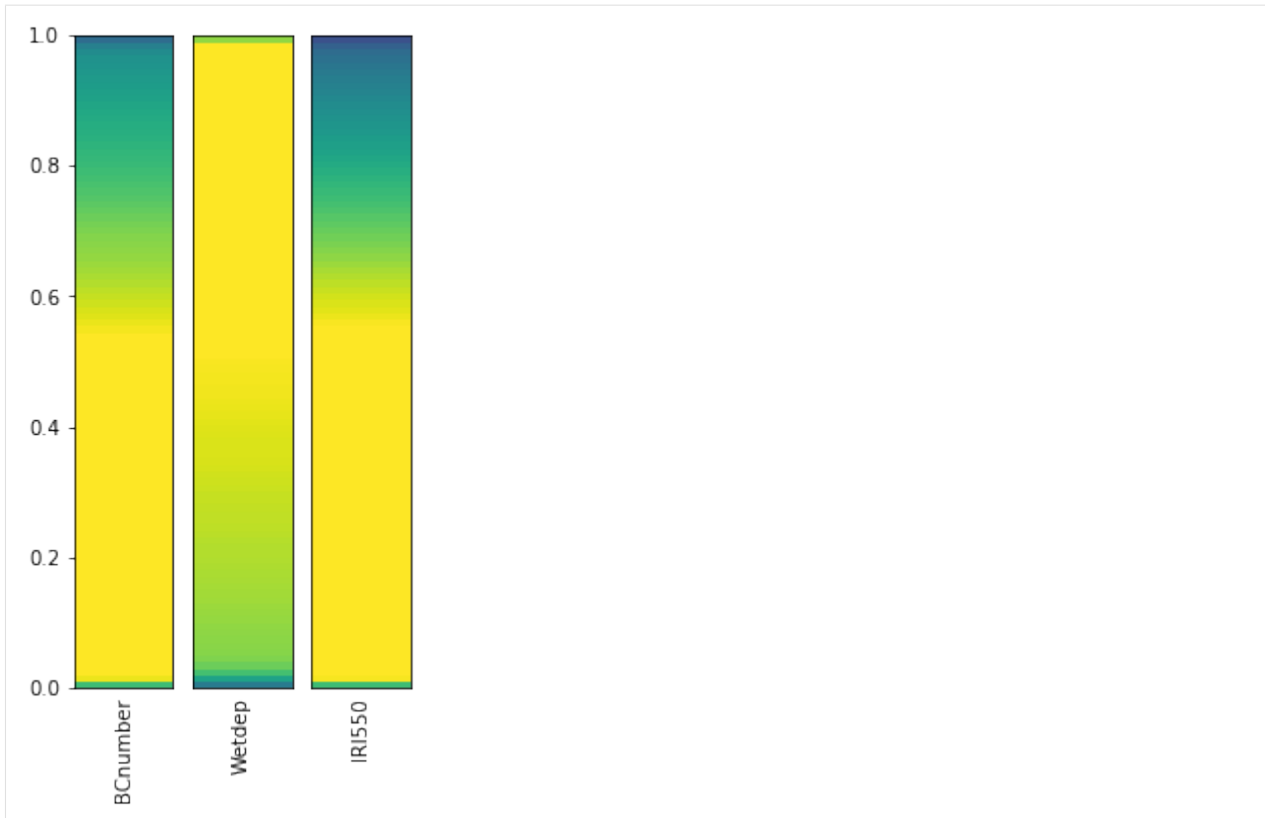
[17]: [Text(0.5,0,'Implausibility')]
```



```
[18]: # Find the valid samples in our full 1million samples by comparing against a given
      ↪ tolerance and threshold
      valid_samples = sampler.batch_constrain(sample_points, batch_size=1000, tolerance=.1)
      print("Remaining points: {}".format(valid_samples.sum()))

      100%|#####| 1000000/1000000 [05:23<00:00, 2972.65sample/s]Remaining points: 729607
```

```
[19]: # Plot the reduced parameter distribution
      constrained_sample = sample_points[valid_samples]
      plot_parameter_space(constrained_sample, fig_size=(3,6))
```

```
[20]: # We can also easily plot the joint distributions
# Only plot every one in 100 points as scatter plots with large numbers of points are
# slow...
import matplotlib

# Mimic Seaborn scaling without requiring the whole package
scale = 1.5
matplotlib.rcParams['font.size'] = 12 * scale
matplotlib.rcParams['axes.labelsize'] = 12 * scale
matplotlib.rcParams['axes.titlesize'] = 12 * scale
matplotlib.rcParams['xtick.labelsize'] = 11 * scale
matplotlib.rcParams['ytick.labelsize'] = 11 * scale
matplotlib.rcParams['lines.linewidth'] = 1.5 * scale
matplotlib.rcParams['lines.markersize'] = 6 * scale
#

m, _ = model.predict(constrained_sample[:, :100].values)
Zs = m.data
# Plot the emulated AAOD value (averaged over observation locations) for each point
grr = pd.plotting.scatter_matrix(constrained_sample[:, :100], c=Zs.mean(axis=1),
# figsize=(12, 10), marker='o',
# hist_kwds={'bins': 20,}, s=20, alpha=.8, vmin=1e-3,
# vmax=1e-2, range_padding=0.,
# density_kwds={'range': [[0., 1.], [0., 1.]], 'colormap':
# 'viridis'},
# )
```

(continues on next page)

(continued from previous page)

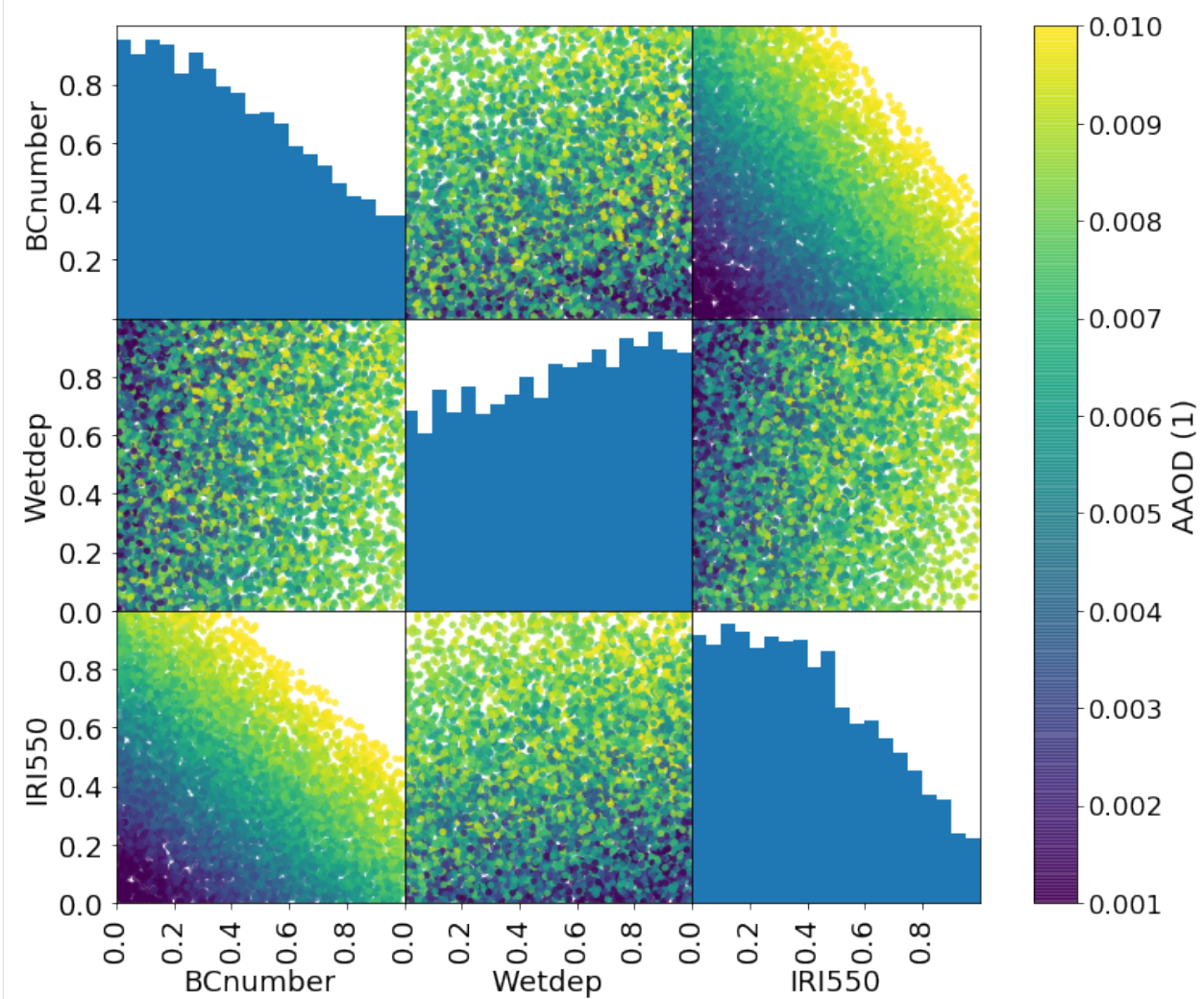
```

# Matplotlib dragons...
grr[0][0].set_yticklabels([0.2, 0.4, 0.6, 0.8], fontsize=12 * scale)
for i in range(2):
    grr[i+1][0].set_yticklabels([0.0, 0.2, 0.4, 0.6, 0.8], fontsize=12 * scale)
for i in range(3):
    grr[2][i].set_xticks([0.0, 0.2, 0.4, 0.6, 0.8])
    grr[2][i].set_xticklabels([0.0, 0.2, 0.4, 0.6, 0.8], fontsize=12 * scale)

plt.colorbar(grr[0][1].collections[0], ax=grr, use_gridspec=True, label='AAOD (1)')

plt.savefig('BCPPE_constrained_params_paper.png', transparent=True)

```



5.4.5 Explore the uncertainty in Direct Radiative Effect of Aerosol in constrained sample-space

```
[21]: dre_test, dre_train = ppe_dre[:n_test], ppe_dre[n_test:]
```

```
ari_model = gp_model(X_train, dre_train, name="ARI", kernel=['Linear', 'Bias'])
ari_model.train()
```

```
[22]: # Calculate the mean and std-dev DRE over each set of sample points
```

```
unconstrained_mean_ari, unconstrained_sd_ari = ari_model.batch_stats(sample_points,
↳ batch_size=1000)
constrained_mean_ari, constrained_sd_ari = ari_model.batch_stats(constrained_sample,
↳ batch_size=1000)
```

```
-----
TypeError                                Traceback (most recent call last)
<ipython-input-22-1ce0d8f67b9a> in <module>
      1 # Calculate the mean and std-dev DRE over each set of sample points
      2
----> 3 unconstrained_mean_ari, unconstrained_sd_ari = ari_model.batch_stats(sample_
↳ points, batch_size=1000)
      4 constrained_mean_ari, constrained_sd_ari = ari_model.batch_stats(constrained_
↳ sample, batch_size=1000)

c:\users\duncan\pycharmprojects\gcm\esem\emulator.py in batch_stats(self, sample_points,
↳ batch_size)
    112         # TODO: Make sample points optional and just sample from a uniform_
↳ distribution if not provided
    113         mean, sd = _tf_stats(self, tf.constant(sample_points),
--> 114                             tf.constant(batch_size, dtype=tf.int64),
    115                             pbar=tf_tqdm(batch_size=batch_size,
    116                                         total=sample_points.shape[0]))

~\miniconda3\envs\gcm_dev\lib\site-packages\tensorflow\python\eager\def_function.py in _
↳ __call__(self, *args, **kwargs)
    826         tracing_count = self.experimental_get_tracing_count()
    827         with trace.Trace(self._name) as tm:
--> 828             result = self._call(*args, **kwargs)
    829             compiler = "xla" if self._experimental_compile else "nonXla"
    830             new_tracing_count = self.experimental_get_tracing_count()

~\miniconda3\envs\gcm_dev\lib\site-packages\tensorflow\python\eager\def_function.py in _
↳ call(self, *args, **kwargs)
    869         # This is the first call of __call__, so we have to initialize.
    870         initializers = []
--> 871         self._initialize(args, kwargs, add_initializers_to=initializers)
    872         finally:
    873             # At this point we know that the initialization is complete (or less

~\miniconda3\envs\gcm_dev\lib\site-packages\tensorflow\python\eager\def_function.py in _
↳ initialize(self, args, kwargs, add_initializers_to)
    724         self._concrete_stateful_fn = (
```

(continues on next page)

(continued from previous page)

```

725         self._stateful_fn._get_concrete_function_internal_garbage_collected( #_
↳ pylint: disable=protected-access
--> 726             *args, **kwargs))
727
728     def invalid_creator_scope(*unused_args, **unused_kwargs):

~\miniconda3\envs\gcm_dev\lib\site-packages\tensorflow\python\eager\function.py in _get_
↳ concrete_function_internal_garbage_collected(self, *args, **kwargs)
2967         args, kwargs = None, None
2968         with self._lock:
-> 2969             graph_function, _ = self._maybe_define_function(args, kwargs)
2970         return graph_function
2971

~\miniconda3\envs\gcm_dev\lib\site-packages\tensorflow\python\eager\function.py in _
↳ maybe_define_function(self, args, kwargs)
3359
3360         self._function_cache.missed.add(call_context_key)
-> 3361         graph_function = self._create_graph_function(args, kwargs)
3362         self._function_cache.primary[cache_key] = graph_function
3363

~\miniconda3\envs\gcm_dev\lib\site-packages\tensorflow\python\eager\function.py in _
↳ create_graph_function(self, args, kwargs, override_flat_arg_shapes)
3204         arg_names=arg_names,
3205         override_flat_arg_shapes=override_flat_arg_shapes,
-> 3206         capture_by_value=self._capture_by_value),
3207         self._function_attributes,
3208         function_spec=self.function_spec,

~\miniconda3\envs\gcm_dev\lib\site-packages\tensorflow\python\framework\func_graph.py_
↳ in func_graph_from_py_func(name, python_func, args, kwargs, signature, func_graph,
↳ autograph, autograph_options, add_control_dependencies, arg_names, op_return_value,
↳ collections, capture_by_value, override_flat_arg_shapes)
988         _, original_func = tf_decorator.unwrap(python_func)
989
--> 990         func_outputs = python_func(*func_args, **func_kwargs)
991
992         # invariant: `func_outputs` contains only Tensors, CompositeTensors,

~\miniconda3\envs\gcm_dev\lib\site-packages\tensorflow\python\eager\def_function.py in_
↳ wrapped_fn(*args, **kwargs)
632         xla_context.Exit()
633     else:
--> 634         out = weak_wrapped_fn().__wrapped__(*args, **kwargs)
635         return out
636

~\miniconda3\envs\gcm_dev\lib\site-packages\tensorflow\python\framework\func_graph.py_
↳ in wrapper(*args, **kwargs)
975         except Exception as e: # pylint:disable=broad-except
976             if hasattr(e, "ag_error_metadata"):

```

(continues on next page)

(continued from previous page)

```
--> 977         raise e.ag_error_metadata.to_exception(e)
    978     else:
    979         raise
```

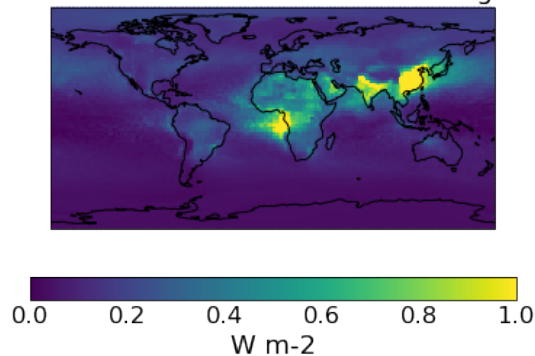
TypeError: in user code:

TypeError: tf_batch_stats() takes from 2 to 3 positional arguments but 4 were given

```
[23]: # The original (unconstrained DRE)
      qplt.pcolormesh(unconstrained_sd_ari, vmin=0., vmax=1)
      plt.gca().coastlines()
```

[23]: <cartopy.mpl.feature_artist.FeatureArtist at 0x2114ad520c8>

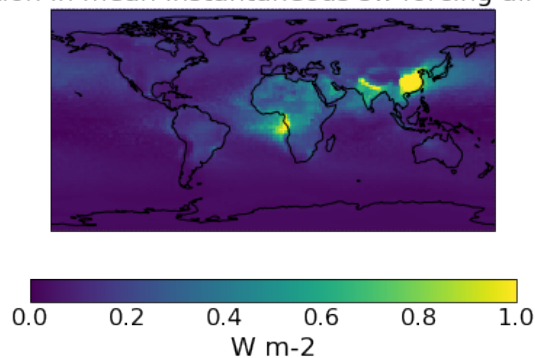
Ensemble standard deviation in mean instantaneous sw forcing all sky top of the atmosphere



```
[24]: # The constrained DRE
      qplt.pcolormesh(constrained_sd_ari, vmin=0., vmax=1)
      plt.gca().coastlines()
```

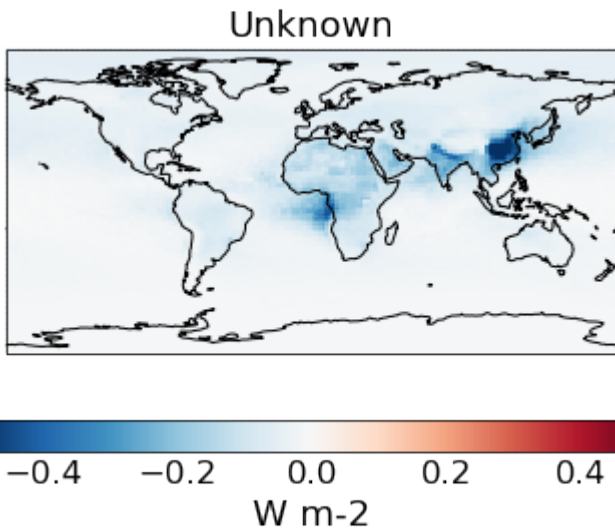
[24]: <cartopy.mpl.feature_artist.FeatureArtist at 0x2114a4e35c8>

Ensemble standard deviation in mean instantaneous sw forcing all sky top of the atmosphere



```
[25]: # The change in spread after the constraint is applied
      qplt.pcolormesh((constrained_sd_ari-unconstrained_sd_ari), cmap='RdBu_r', vmin=-5e-1,
      ↪vmax=5e-1)
      plt.gca().coastlines()
```

```
[25]: <cartopy.mpl.feature_artist.FeatureArtist at 0x21132beb148>
```



```
[ ]:
```

5.5 Calibrating GPs using MCMC

```
[1]: import os
      # Ignore my broken HDF5 install...
      os.putenv("HDF5_DISABLE_VERSION_CHECK", '1')

[2]: import pandas as pd
      import numpy as np
      import iris

      from utils import get_aeronet_data, get_bc_ppe_data
      from esim import gp_model
      from esim.sampler import MCMCSampler

      import os

      import matplotlib.pyplot as plt
      %matplotlib inline

      # GPU = "1"

      # os.environ["CUDA_VISIBLE_DEVICES"] = GPU
```

5.5.1 Read in the parameters and observables

```
[3]: ppe_params, ppe_aod = get_bc_ppe_data()
```

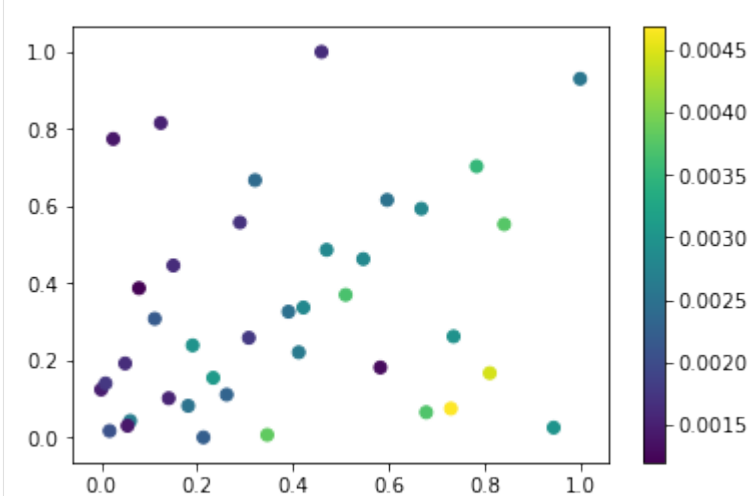
```
/Users/watson-parris/miniconda3/envs/gcem/lib/python3.8/site-packages/iris/__init__.py:
↳249: IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and
↳will be removed in a future release. Please remove code that sets this property.
warn_deprecated(msg.format(name))
/Users/watson-parris/miniconda3/envs/gcem/lib/python3.8/site-packages/iris/__init__.py:
↳249: IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and
↳will be removed in a future release. Please remove code that sets this property.
warn_deprecated(msg.format(name))
```

```
[4]: # Calculate the global, annual mean AAOD (CIS will automatically apply the weights)
mean_aaod, = ppe_aod.collapsed(['latitude', 'longitude', 'time'], 'mean')
```

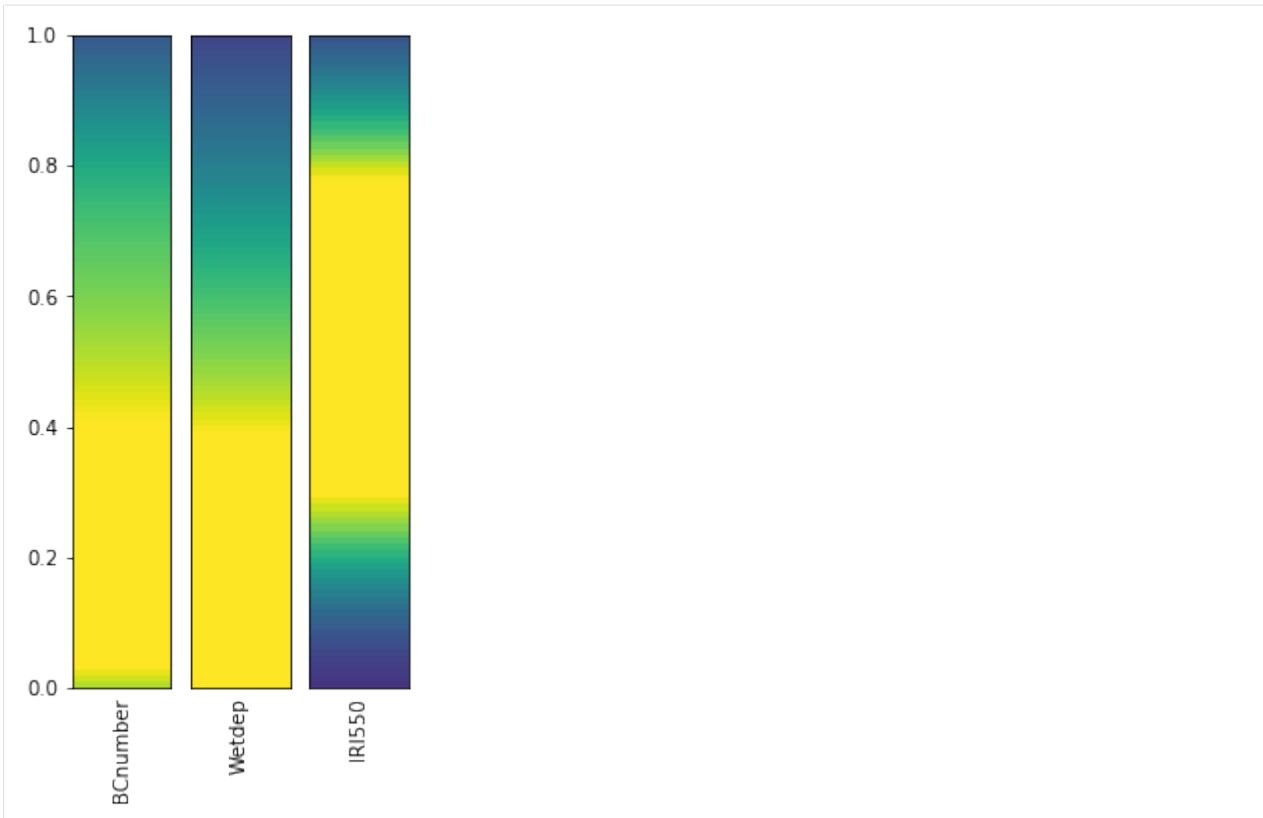
```
WARNING:root:Creating guessed bounds as none exist in file
WARNING:root:Creating guessed bounds as none exist in file
WARNING:root:Creating guessed bounds as none exist in file
/Users/watson-parris/miniconda3/envs/gcem/lib/python3.8/site-packages/iris/analysis/
↳cartography.py:394: UserWarning: Using DEFAULT_SPHERICAL_EARTH_RADIUS.
warnings.warn("Using DEFAULT_SPHERICAL_EARTH_RADIUS.")
```

```
[5]: plt.scatter(ppe_params.BCnumber, ppe_params.Wetdep, c=mean_aaod.data)
plt.colorbar()
```

```
[5]: <matplotlib.colorbar.Colorbar at 0x7fc03dea0580>
```



```
[6]: from esim.utils import plot_parameter_space
plot_parameter_space(ppe_params, fig_size=(3,6))
```



```
[7]: n_test = 8
```

```
X_test, X_train = ppe_params[:n_test], ppe_params[n_test:]
Y_test, Y_train = mean_aaod[:n_test], mean_aaod[n_test:]
```

5.5.2 Setup and run the models

```
[8]: model = gp_model(X_train, Y_train)
```

WARNING: Using default kernel - be sure you understand the assumptions this implies.
 ↳ Consult e.g. <http://www.cs.toronto.edu/~duvenaud/cookbook/> for an excellent
 ↳ description of different kernel choices.

```
[9]: model.train()
```

```
[10]: m, v = model.predict(X_test.values)
```

```
[11]: Y_test.data
```

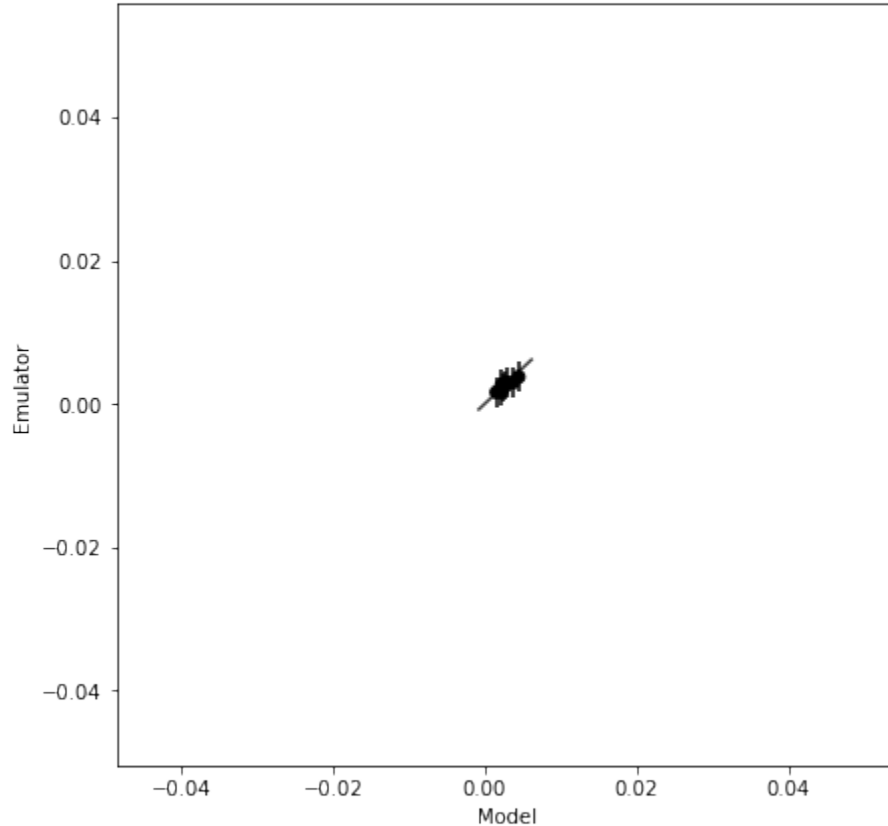
```
[11]: masked_array(data=[0.002822510314728863, 0.002615034735649063,
                        0.002186747212596059, 0.0015112621889264352,
                        0.004456776854431466, 0.0025203727123839113,
                        0.0022378865435589133, 0.003730134076583199],
                    mask=[False, False, False, False, False, False, False, False],
                    fill_value=1e+20)
```



```
[12]: from esim.utils import validation_plot
```

```
validation_plot(Y_test.data.flatten(), m.data.flatten(), v.data.flatten())
```

```
Proportion of 'Bad' estimates : 0.00%
```



```
[13]: # Set the objective as one of the test datasets
```

```
sampler = MCMCSampler(model, Y_test[0])
```

```
[14]: samples = sampler.sample(n_samples=8000, mcmc_kwargs=dict(num_burnin_steps=1000) )
```

```
Acceptance rate: 0.9026122150748647
```

```
[15]: new_samples = pd.DataFrame(data=samples, columns=ppe_params.columns)
```

```
m, _ = model.predict(new_samples.values)
```

```
Zs = m.data
```

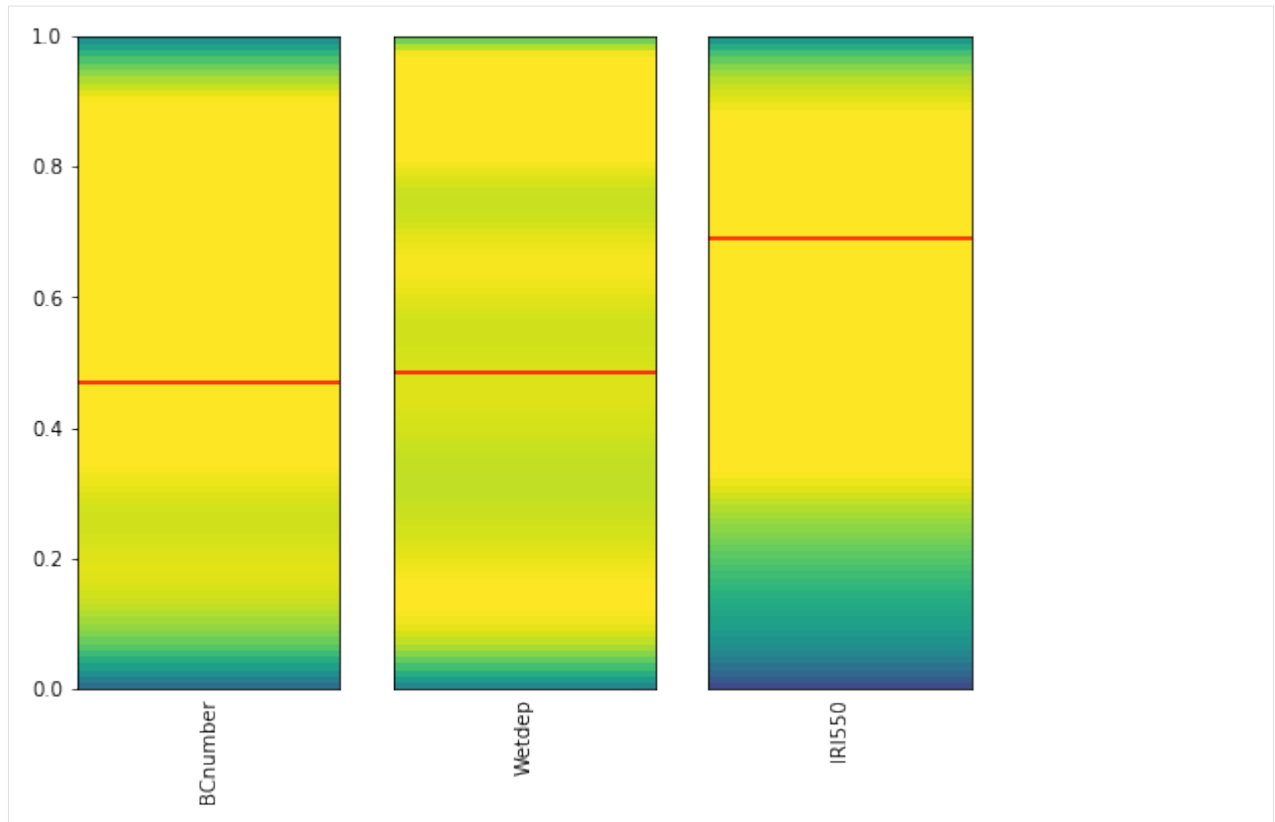
```
[16]: print("Sample mean: {}".format(Zs.mean()))
```

```
print("Sample std dev: {}".format(Zs.std()))
```

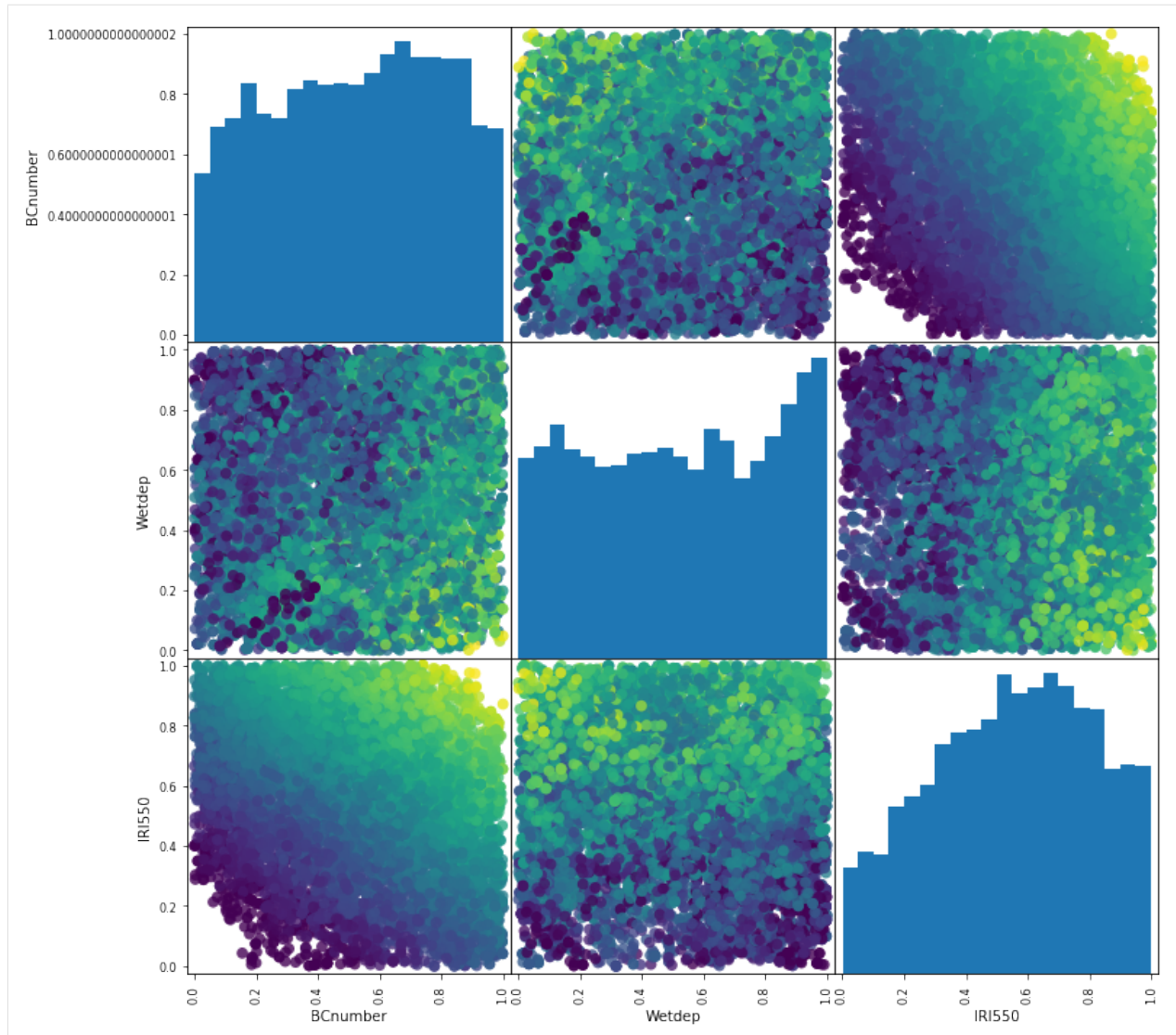
```
Sample mean: 0.0026673169026350933
```

```
Sample std dev: 0.000807943471868746
```

```
[17]: plot_parameter_space(new_samples, target_df=X_test.iloc[0])
```



```
[18]: grr = pd.plotting.scatter_matrix(new_samples, c=Zs, figsize=(12, 12), marker='o',  
                                     hist_kws={'bins': 20}, s=60, alpha=.8, vmin=1e-3,  
                                     ↪vmax=5e-3)
```



```
[19]: from esem.abc_sampler import ABCSampler, constrain
      from esem.utils import get_random_params
```

```
[20]: sampler = ABCSampler(model, Y_test[0])

      samples = sampler.sample(n_samples=2000, threshold=0.5)
      valid_points = pd.DataFrame(data=samples, columns=ppe_params.columns)

      Acceptance rate: 0.34405642525374164
```

```
[21]: m, _ = model.predict(valid_points.values)
      Zs = m.data
```

```
[22]: print("Sample mean: {}".format(Zs.mean()))
      print("Sample std dev: {}".format(Zs.std()))

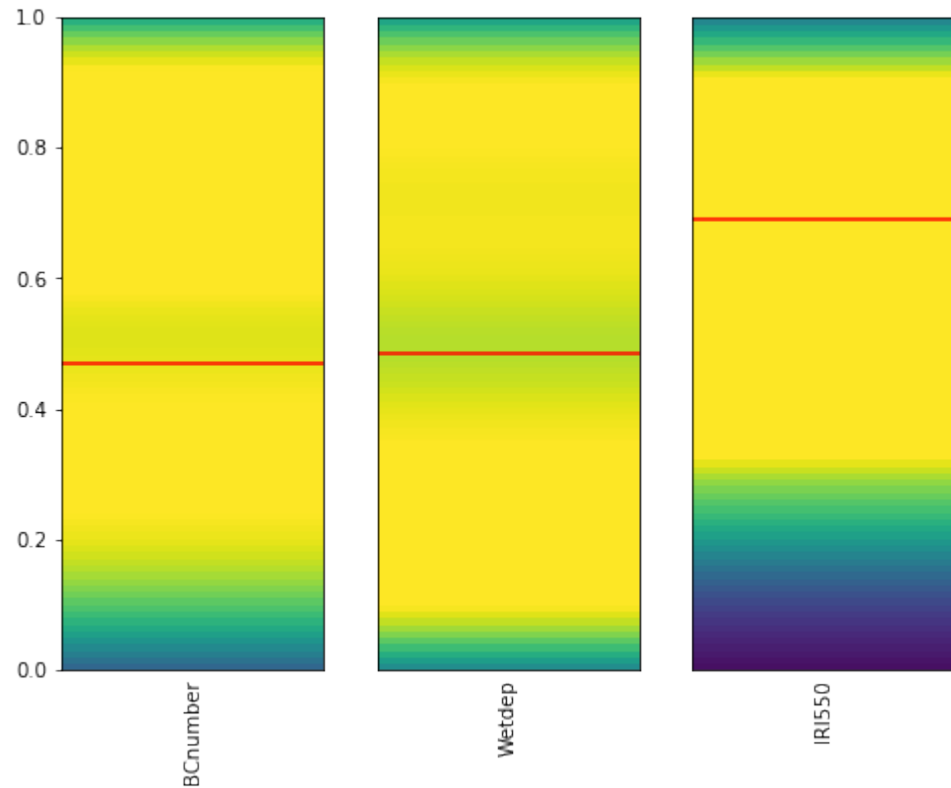
      Sample mean: 0.0028015036302732293
```

(continues on next page)

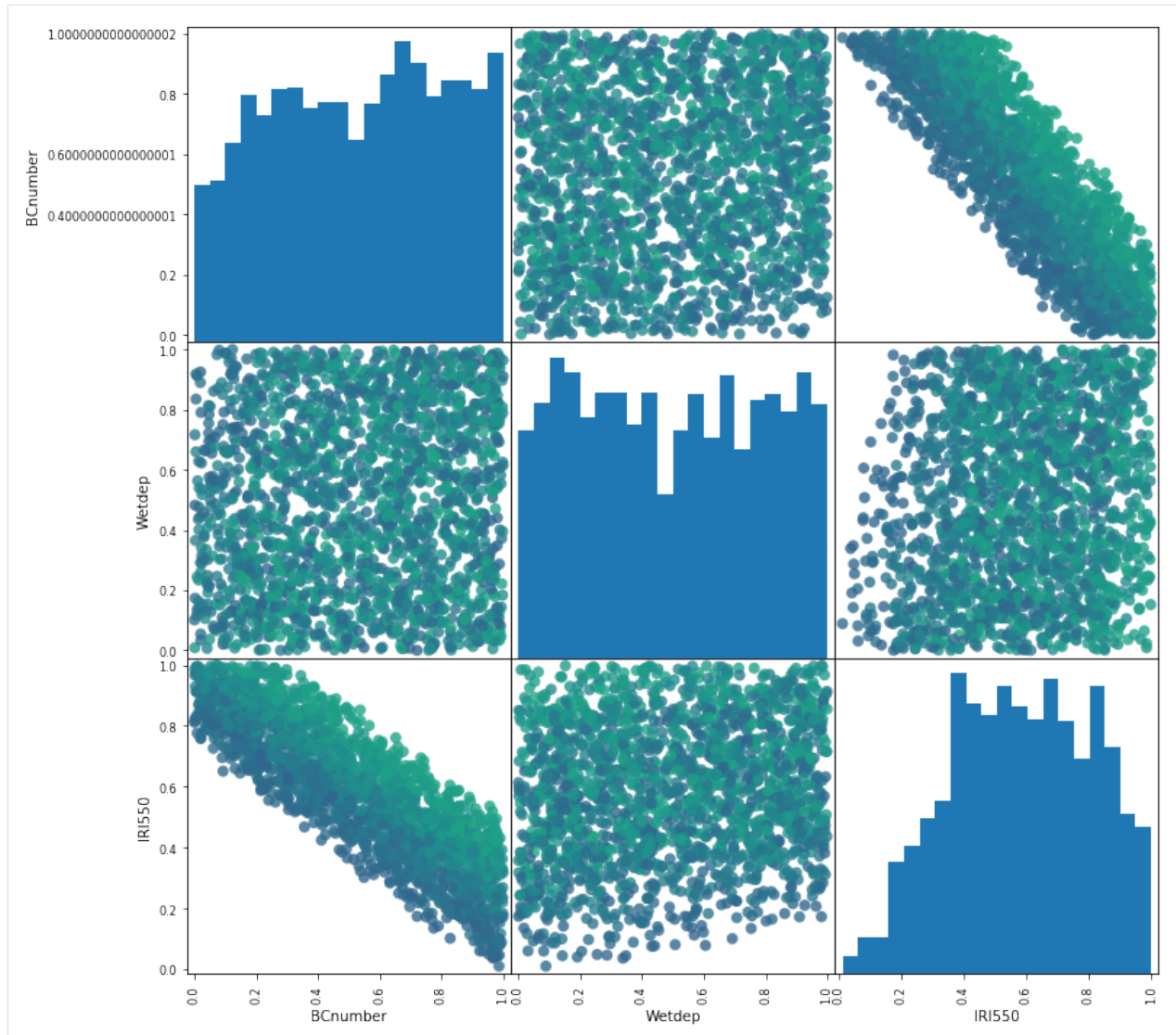
(continued from previous page)

Sample std dev: 0.00030028439418992273

```
[23]: plot_parameter_space(valid_points, target_df=X_test.iloc[0])
```



```
[24]: grr = pd.plotting.scatter_matrix(valid_points, c=Zs, figsize=(12, 12), marker='o',
                                         hist_kws={'bins': 20}, s=60, alpha=.8, vmin=1e-3,
                                         ↪vmax=5e-3)
```



```
[ ]:
```

5.6 CMIP6 Emulation

```
[1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from esim import gp_model
from esim.utils import validation_plot
%matplotlib inline
```

```
[2]: df = pd.read_csv('CMIP6_scenarios.csv', index_col=0).dropna()
```

```
[3]: # These are the models included
```

```
df.model.unique()
```

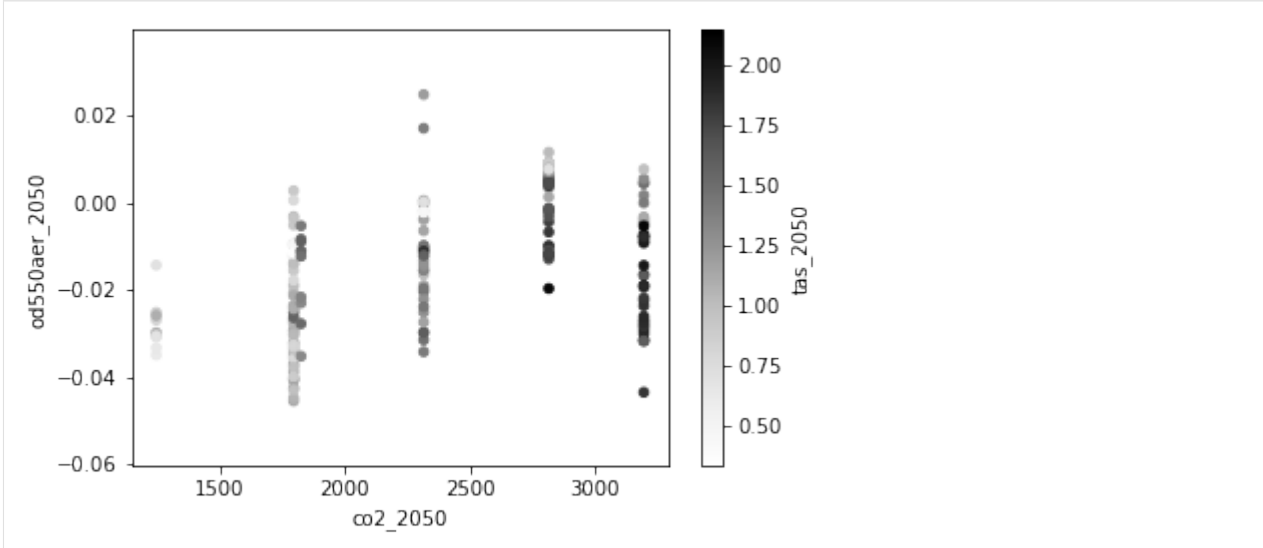
```
[3]: array(['CanESM5', 'ACCESS-ESM1-5', 'ACCESS-CM2', 'MPI-ESM1-2-HR',  
        'MIROC-ES2L', 'HadGEM3-GC31-LL', 'UKESM1-0-LL', 'MPI-ESM1-2-LR',  
        'CESM2', 'CESM2-WACCM', 'NorESM2-LM'], dtype=object)
```

```
[4]: # And these scenarios
```

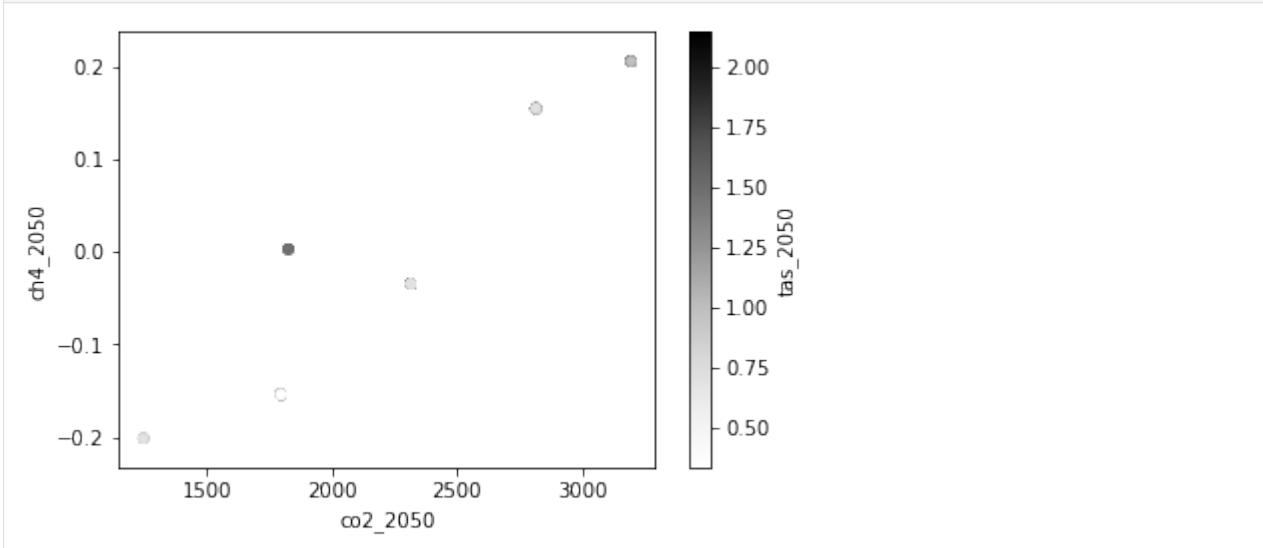
```
df.scenario.unique()
```

```
[4]: array(['ssp126', 'ssp119', 'ssp245', 'ssp370', 'ssp585', 'ssp434'],  
        dtype=object)
```

```
[5]: ax = df.plot.scatter(x='co2_2050', y='od550aer_2050', c='tas_2050')
```



```
[6]: ax = df.plot.scatter(x='co2_2050', y='ch4_2050', c='tas_2050')
```



```
[7]: # Collapse ensemble members
```

```
df = df.groupby(['model', 'scenario']).mean()
df
```

```
[7]:
```

		tas_2050	od550aer_2050	tas_2100	od550aer_2100	\
model	scenario					
ACCESS-CM2	ssp126	1.000582	-0.025191	1.407198	-0.038552	
	ssp245	1.187015	-0.012222	2.489487	-0.027589	
	ssp370	1.103402	0.006077	3.791109	0.001188	
	ssp585	1.478602	-0.008842	5.016310	-0.016805	
ACCESS-ESM1-5	ssp126	0.819904	-0.013947	0.967981	-0.015895	
	ssp245	1.078853	-0.004727	2.052867	-0.009796	
	ssp370	1.047218	0.003862	3.422980	0.000495	
	ssp585	1.412020	-0.002210	4.096812	-0.001356	
CESM2	ssp126	0.974787	-0.003323	1.134402	-0.010735	
	ssp245	1.119462	0.000371	2.303586	0.000121	
	ssp370	1.154434	0.007554	3.501974	0.016237	
	ssp585	1.453782	0.002178	4.980513	0.016340	
CESM2-WACCM	ssp126	0.770853	0.000561	1.053462	-0.006278	
	ssp245	1.118329	0.000520	2.252389	0.006330	
	ssp370	1.022903	0.011563	3.526784	0.035894	
	ssp585	1.311378	0.005453	5.010599	0.038768	
CanESM5	ssp119	0.641977	-0.030857	0.333723	-0.036624	
	ssp126	0.989167	-0.029769	0.953098	-0.041451	
	ssp245	1.318495	-0.016138	2.454336	-0.033573	
	ssp370	1.724648	-0.004288	4.744037	-0.021374	
	ssp434	1.304015	-0.023741	1.863815	-0.040477	
HadGEM3-GC31-LL	ssp585	1.877571	-0.025367	5.984793	-0.040852	
	ssp126	0.768553	-0.026508	1.314000	-0.037692	
	ssp245	1.262740	-0.014397	2.651430	-0.024341	
	ssp585	1.606117	-0.007902	5.547793	-0.015396	
MIROC-ES2L	ssp119	0.703733	-0.014334	0.557255	-0.020002	
	ssp126	0.785037	-0.017767	0.572429	-0.028761	
	ssp245	0.789146	-0.010066	1.589075	-0.019236	
	ssp370	1.131194	0.001314	2.623115	0.006749	
	ssp585	1.117971	-0.003265	3.748485	-0.003236	
MPI-ESM1-2-HR	ssp126	0.370221	-0.009666	0.351424	-0.014088	
	ssp245	0.735819	-0.001871	1.352010	-0.008150	
	ssp370	0.819827	0.008950	2.616101	0.003430	
	ssp585	0.946618	0.004603	3.153594	-0.004407	
MPI-ESM1-2-LR	ssp126	0.518146	-0.009654	0.358184	-0.014067	
	ssp370	0.877299	0.008940	2.610222	0.003437	
	ssp585	1.029724	0.004613	3.278956	-0.004377	
NorESM2-LM	ssp126	0.331449	-0.011588	0.376486	-0.017219	
	ssp245	0.533042	-0.001380	1.295670	-0.005419	
	ssp370	0.704050	0.007691	2.552758	0.022834	
	ssp585	0.983765	0.007667	2.957163	0.009652	
UKESM1-0-LL	ssp119	0.918189	-0.027682	0.970074	-0.035182	
	ssp126	1.284645	-0.024166	1.563342	-0.036149	
	ssp245	1.604553	-0.010934	3.040434	-0.022142	
	ssp370	1.749074	0.004618	5.046036	0.007180	
	ssp434	1.511831	-0.009184	2.371976	-0.027848	
	ssp585	2.028151	-0.007564	6.038747	-0.006429	

(continues on next page)

(continued from previous page)

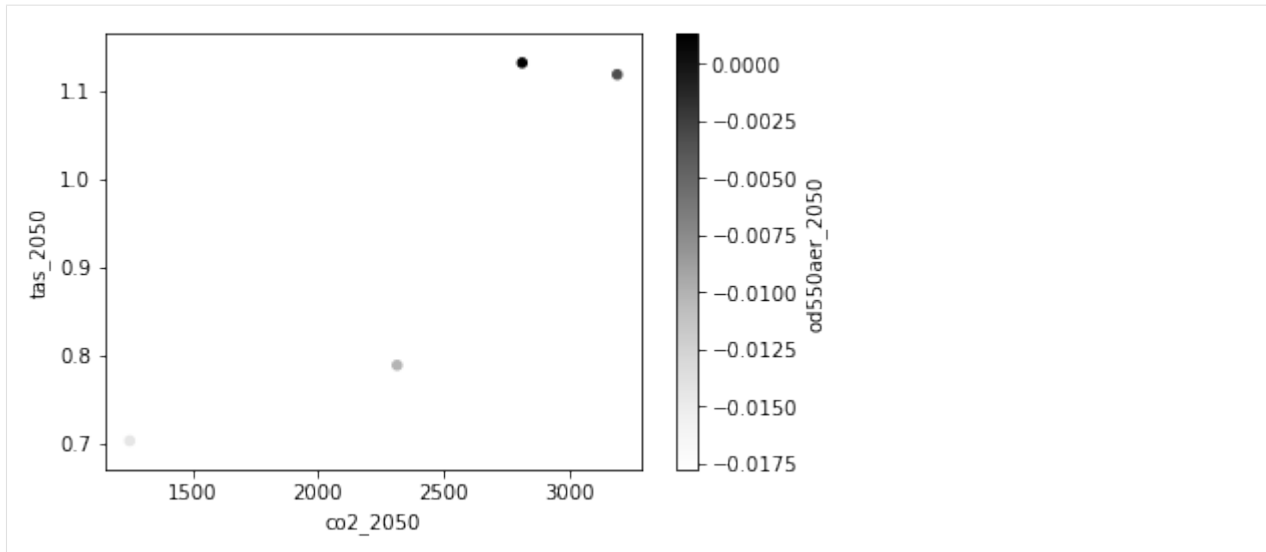
model	scenario	co2_2050	co2_2100	so2_2050	so2_2100	\
ACCESS-CM2	ssp126	1795.710867	1848.864201	-0.064020	-0.082066	
	ssp245	2314.385253	3932.717046	-0.035997	-0.059775	
	ssp370	2813.146604	6912.965613	0.004142	-0.020179	
	ssp585	3192.373467	10283.292188	-0.028185	-0.059536	
ACCESS-ESM1-5	ssp126	1795.710867	1848.864201	-0.064020	-0.082066	
	ssp245	2314.385253	3932.717046	-0.035997	-0.059775	
	ssp370	2813.146604	6912.965613	0.004142	-0.020179	
	ssp585	3192.373467	10283.292188	-0.028185	-0.059536	
CESM2	ssp126	1795.710867	1848.864201	-0.064020	-0.082066	
	ssp245	2314.385253	3932.717046	-0.035997	-0.059775	
	ssp370	2813.146604	6912.965613	0.004142	-0.020179	
	ssp585	3192.373467	10283.292188	-0.028185	-0.059536	
CESM2-WACCM	ssp126	1795.710867	1848.864201	-0.064020	-0.082066	
	ssp245	2314.385253	3932.717046	-0.035997	-0.059775	
	ssp370	2813.146604	6912.965613	0.004142	-0.020179	
	ssp585	3192.373467	10283.292188	-0.028185	-0.059536	
CanESM5	ssp119	1247.788346	905.867767	-0.069425	-0.080790	
	ssp126	1795.710867	1848.864201	-0.064020	-0.082066	
	ssp245	2314.385253	3932.717046	-0.035997	-0.059775	
	ssp370	2813.146604	6912.965613	0.004142	-0.020179	
	ssp434	1825.355018	1612.733274	-0.055372	-0.077382	
	ssp585	3192.373467	10283.292188	-0.028185	-0.059536	
HadGEM3-GC31-LL	ssp126	1795.710867	1848.864201	-0.064020	-0.082066	
	ssp245	2314.385253	3932.717046	-0.035997	-0.059775	
	ssp585	3192.373467	10283.292188	-0.028185	-0.059536	
MIROC-ES2L	ssp119	1247.788346	905.867767	-0.069425	-0.080790	
	ssp126	1795.710867	1848.864201	-0.064020	-0.082066	
	ssp245	2314.385253	3932.717046	-0.035997	-0.059775	
	ssp370	2813.146604	6912.965613	0.004142	-0.020179	
	ssp585	3192.373467	10283.292188	-0.028185	-0.059536	
MPI-ESM1-2-HR	ssp126	1795.710867	1848.864201	-0.064020	-0.082066	
	ssp245	2314.385253	3932.717046	-0.035997	-0.059775	
	ssp370	2813.146604	6912.965613	0.004142	-0.020179	
	ssp585	3192.373467	10283.292188	-0.028185	-0.059536	
MPI-ESM1-2-LR	ssp126	1795.710867	1848.864201	-0.064020	-0.082066	
	ssp370	2813.146604	6912.965613	0.004142	-0.020179	
	ssp585	3192.373467	10283.292188	-0.028185	-0.059536	
NorESM2-LM	ssp126	1795.710867	1848.864201	-0.064020	-0.082066	
	ssp245	2314.385253	3932.717046	-0.035997	-0.059775	
	ssp370	2813.146604	6912.965613	0.004142	-0.020179	
	ssp585	3192.373467	10283.292188	-0.028185	-0.059536	
UKESM1-0-LL	ssp119	1247.788346	905.867767	-0.069425	-0.080790	
	ssp126	1795.710867	1848.864201	-0.064020	-0.082066	
	ssp245	2314.385253	3932.717046	-0.035997	-0.059775	
	ssp370	2813.146604	6912.965613	0.004142	-0.020179	
	ssp434	1825.355018	1612.733274	-0.055372	-0.077382	
	ssp585	3192.373467	10283.292188	-0.028185	-0.059536	
		ch4_2050	ch4_2100			

(continues on next page)

(continued from previous page)

model	scenario		
ACCESS-CM2	ssp126	-0.153999	-0.241390
	ssp245	-0.034673	-0.092845
	ssp370	0.154481	0.368859
	ssp585	0.205365	0.104140
ACCESS-ESM1-5	ssp126	-0.153999	-0.241390
	ssp245	-0.034673	-0.092845
	ssp370	0.154481	0.368859
	ssp585	0.205365	0.104140
CESM2	ssp126	-0.153999	-0.241390
	ssp245	-0.034673	-0.092845
	ssp370	0.154481	0.368859
	ssp585	0.205365	0.104140
CESM2-WACCM	ssp126	-0.153999	-0.241390
	ssp245	-0.034673	-0.092845
	ssp370	0.154481	0.368859
	ssp585	0.205365	0.104140
CanESM5	ssp119	-0.201275	-0.257405
	ssp126	-0.153999	-0.241390
	ssp245	-0.034673	-0.092845
	ssp370	0.154481	0.368859
	ssp434	0.002907	-0.032744
	ssp585	0.205365	0.104140
HadGEM3-GC31-LL	ssp126	-0.153999	-0.241390
	ssp245	-0.034673	-0.092845
	ssp585	0.205365	0.104140
MIROC-ES2L	ssp119	-0.201275	-0.257405
	ssp126	-0.153999	-0.241390
	ssp245	-0.034673	-0.092845
	ssp370	0.154481	0.368859
	ssp585	0.205365	0.104140
MPI-ESM1-2-HR	ssp126	-0.153999	-0.241390
	ssp245	-0.034673	-0.092845
	ssp370	0.154481	0.368859
	ssp585	0.205365	0.104140
MPI-ESM1-2-LR	ssp126	-0.153999	-0.241390
	ssp370	0.154481	0.368859
	ssp585	0.205365	0.104140
NorESM2-LM	ssp126	-0.153999	-0.241390
	ssp245	-0.034673	-0.092845
	ssp370	0.154481	0.368859
	ssp585	0.205365	0.104140
UKESM1-0-LL	ssp119	-0.201275	-0.257405
	ssp126	-0.153999	-0.241390
	ssp245	-0.034673	-0.092845
	ssp370	0.154481	0.368859
	ssp434	0.002907	-0.032744
	ssp585	0.205365	0.104140

```
[8]: ax = df.query("model == 'MIROC-ES2L'").plot.scatter(x='co2_2050', y='tas_2050', c=
      ↪ 'od550aer_2050')
```



```
[9]: from utils import normalize
      # Merge the year columns in to a long df
      df=pd.wide_to_long(df.reset_index(), ["tas", "od550aer", "co2", "ch4", "so2"], i=['model', 'scenario'], j="year", suffix='_(\d+)')
      # Choose only the 2050 data since the aerosol signal is pretty non-existent by 2100
      df = df[df.index.isin(["_2050"], level=2)]
```

```
[10]: df.describe()
```

```
[10]:
```

	tas	od550aer	co2	ch4	so2
count	47.000000	47.000000	47.000000	47.000000	47.000000
mean	1.085538	-0.006851	2415.708964	0.024789	-0.035145
std	0.381735	0.011775	613.880074	0.152382	0.025320
min	0.331449	-0.030857	1247.788346	-0.201275	-0.069425
25%	0.804487	-0.014140	1795.710867	-0.153999	-0.064020
50%	1.047218	-0.004727	2314.385253	-0.034673	-0.035997
75%	1.307697	0.003020	2813.146604	0.154481	-0.028185
max	2.028151	0.011563	3192.373467	0.205365	0.004142

```
[11]: # Do a 20/80 split of the data for test and training
      msk = np.random.rand(len(df)) < 0.8
      train, test = df[msk], df[~msk]
```

5.6.1 Try a few different models

```
[12]: from esim.utils import leave_one_out, prediction_within_ci
      from scipy import stats

      # Try just modelling the temperature based on cumulative CO2
      res = leave_one_out(df[['co2']], df[['tas']].values, model='GaussianProcess', kernel=[
          'Linear'])

      r2_values = stats.linregress(*np.squeeze(np.asarray(res, dtype=float)).T[0:2])[2]**2
```

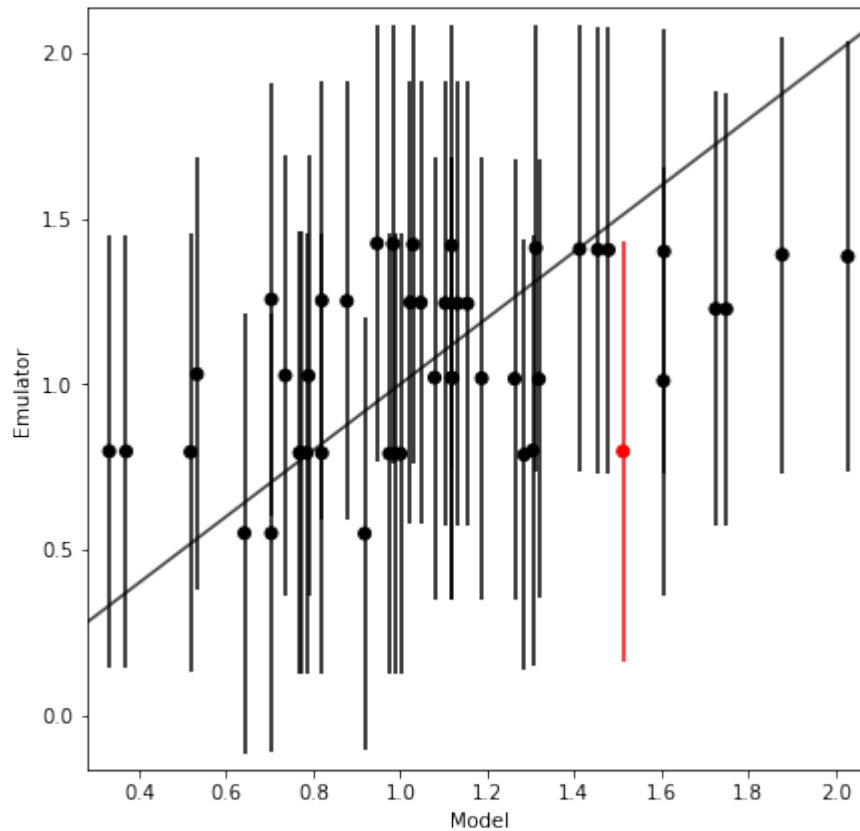
(continues on next page)

(continued from previous page)

```
print("R^2: {:.2f}".format(r2_values))
validation_plot(*np.squeeze(np.asarray(res, dtype=float)).T)
```

R^2: 0.25

Proportion of 'Bad' estimates : 2.13%

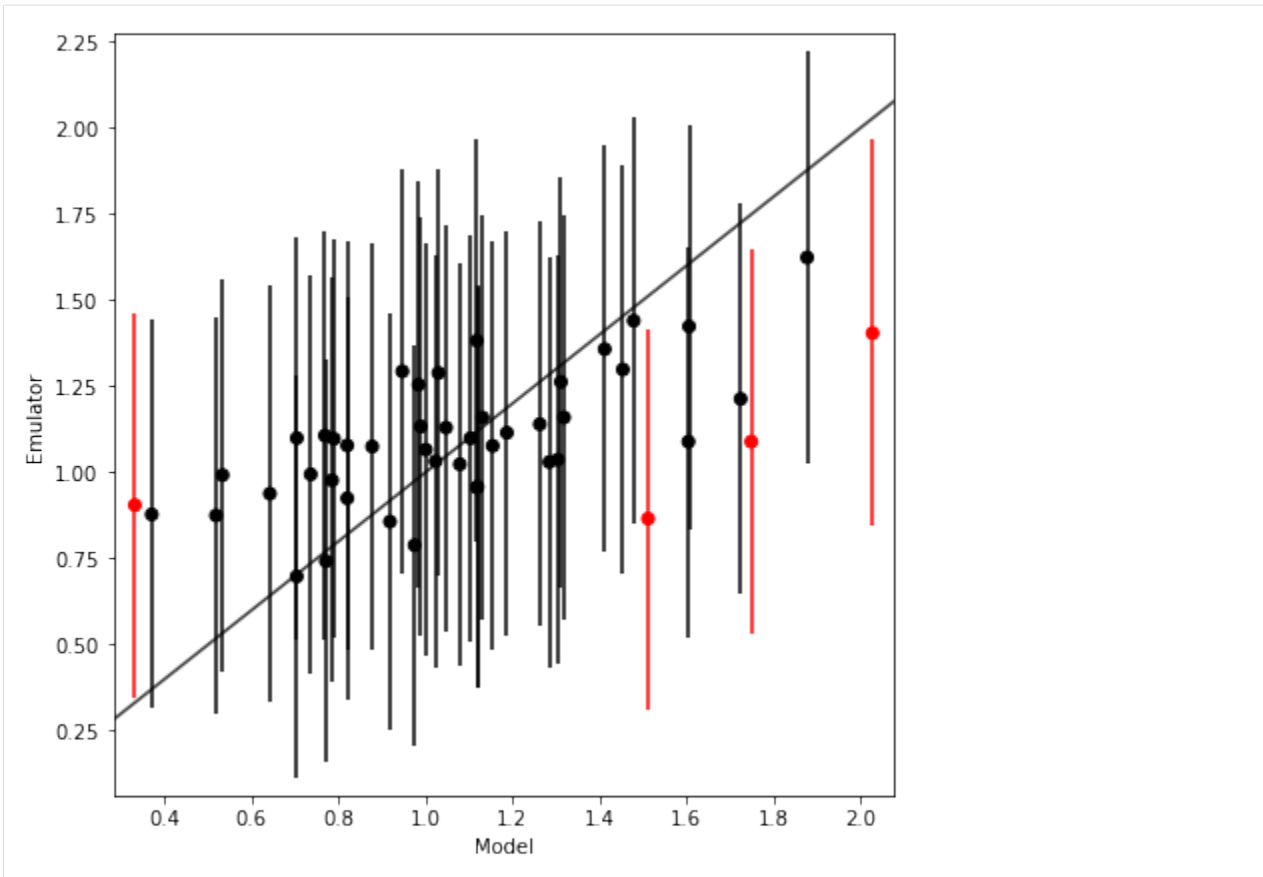


```
[13]: # This model still doesn't do brilliantly, but it's better than just CO2
res = leave_one_out(df[['co2', 'od550aer']], df[['tas']].values, model='GaussianProcess',
    ↪ kernel=['Linear'])
```

```
r2_values = stats.linregress(*np.squeeze(np.asarray(res, dtype=float)).T[0:2])[2]**2
print("R^2: {:.2f}".format(r2_values))
validation_plot(*np.squeeze(np.asarray(res, dtype=float)).T)
```

R^2: 0.40

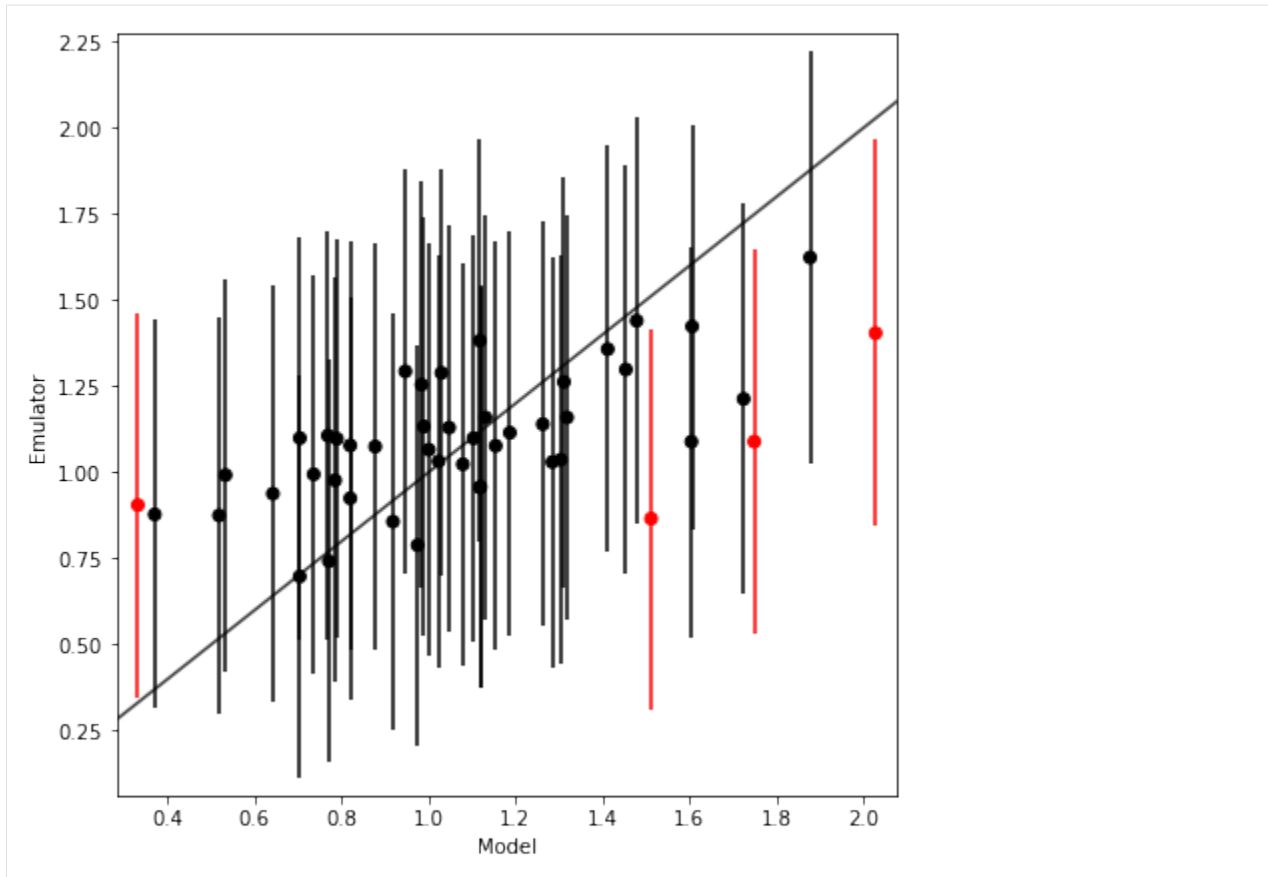
Proportion of 'Bad' estimates : 8.51%



```
[14]: # Adding Methane doesn't seem to improve the picture
res = leave_one_out(df[['co2', 'od550aer', 'ch4']], df[['tas']].values, model=
↳ 'GaussianProcess', kernel=['Linear', 'Bias'])

r2_values = stats.linregress(*np.squeeze(np.asarray(res, dtype=float)).T[0:2])[2]**2
print("R^2: {:.2f}".format(r2_values))
validation_plot(*np.squeeze(np.asarray(res, dtype=float)).T)

R^2: 0.40
Proportion of 'Bad' estimates : 8.51%
```



5.6.2 Plot the best

```
[15]: m = gp_model(df[['co2', 'od550aer']], df[['tas']].values, kernel=['Linear'])
      m.train()
```

```
[16]: # Sample a large AOD/CO2 space using the emulator
      xx, yy = np.meshgrid(np.linspace(0, 4000, 25), np.linspace(-.05, 0.05, 20))
      X_new = np.stack([xx.flat, yy.flat], axis=1)
      Y_new, Y_new_sigma = m.predict(X_new)
```

```
[17]: # Calculate the scenario mean values for comparison
      scn_mean = train.groupby(['scenario']).mean()
```

```
[18]: import matplotlib

      scale = 1.5
      matplotlib.rcParams['font.size'] = 12 * scale
      matplotlib.rcParams['lines.linewidth'] = 1.5 * scale
      matplotlib.rcParams['lines.markersize'] = 6 * scale

      plt.figure(figsize=(12, 6))
```

(continues on next page)

(continued from previous page)

```

norm = matplotlib.colors.Normalize(vmin=-2.5,vmax=2.5)
p = plt.contourf(xx, yy, Y_new.reshape(xx.shape), norm=norm, levels=30, cmap='RdBu_r')

plt.scatter(train.co2, train.od550aer, c=train.tas, norm=norm, edgecolors='k', cmap=
    ↪ 'RdBu_r', marker='x')
plt.scatter(scen_mean.co2, scen_mean.od550aer, c=scen_mean.tas, norm=norm, edgecolors='k',
    ↪ cmap='RdBu_r', marker='s')

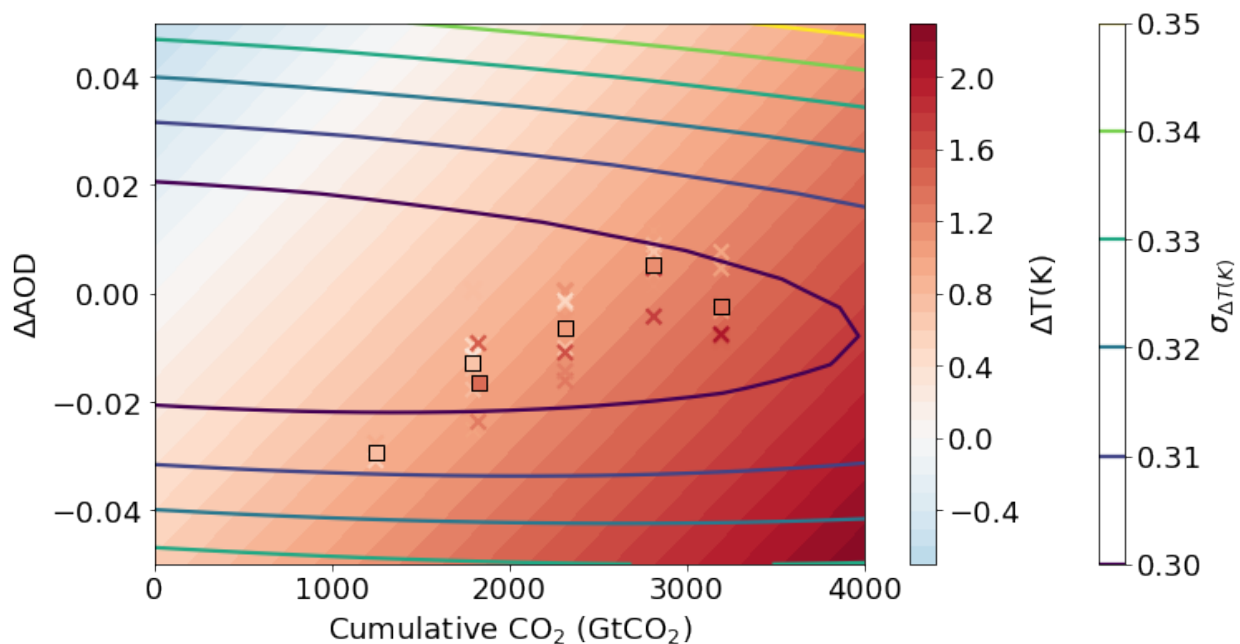
c = plt.contour(xx, yy, np.sqrt(Y_new_sigma.reshape(xx.shape)), cmap='viridis', levels=6)

plt.setp(plt.gca(), xlabel='Cumulative CO2$ (GtCO2$)', ylabel='$\Delta AOD$')

plt.colorbar(c, label='$\sigma_{\Delta T(K)}$')
plt.colorbar(p, label='$\Delta T(K)$')

# Cumulative CO2, delta T and delta AOD all relative to a 2015-2020 average. Each point
    ↪ represents a single model integration for different scenarios in the CMIP6 archive.
plt.savefig('CMIP6_emulator_paper_v1.1.png', transparent=True)

```



5.6.3 Sample emissions for a particular temperature target

```

[19]: from esim.sampler import MCMCSampler

# The MCMC algorithm works much better with a normalised parameter range, so recreate
    ↪ the model
m = gp_model(pd.concat([df[['co2']]/4000, (df[['od550aer']]+0.05)/0.1], axis=1), df[['tas
    ↪ ']].values, kernel=['Linear'])
m.train()

```

(continues on next page)

(continued from previous page)

```
# Target 1.2 degrees above present day (roughly 2 degrees above pre-industrial)
sampler = MCMCSampler(m, np.asarray([1.2], dtype=np.float64))
samples = sampler.sample(n_samples=8000, mcmc_kwargs=dict(num_burnin_steps=1000) )
```

```
Acceptance rate: 0.9614173964786951
```

```
[20]: # Get the emulated temperatures for these samples
new_samples = pd.DataFrame(data=samples, columns=['co2', 'od550aer'])
Z, _ = m.predict(new_samples.values)
```

```
[21]: fig = plt.figure(figsize=(9, 6))

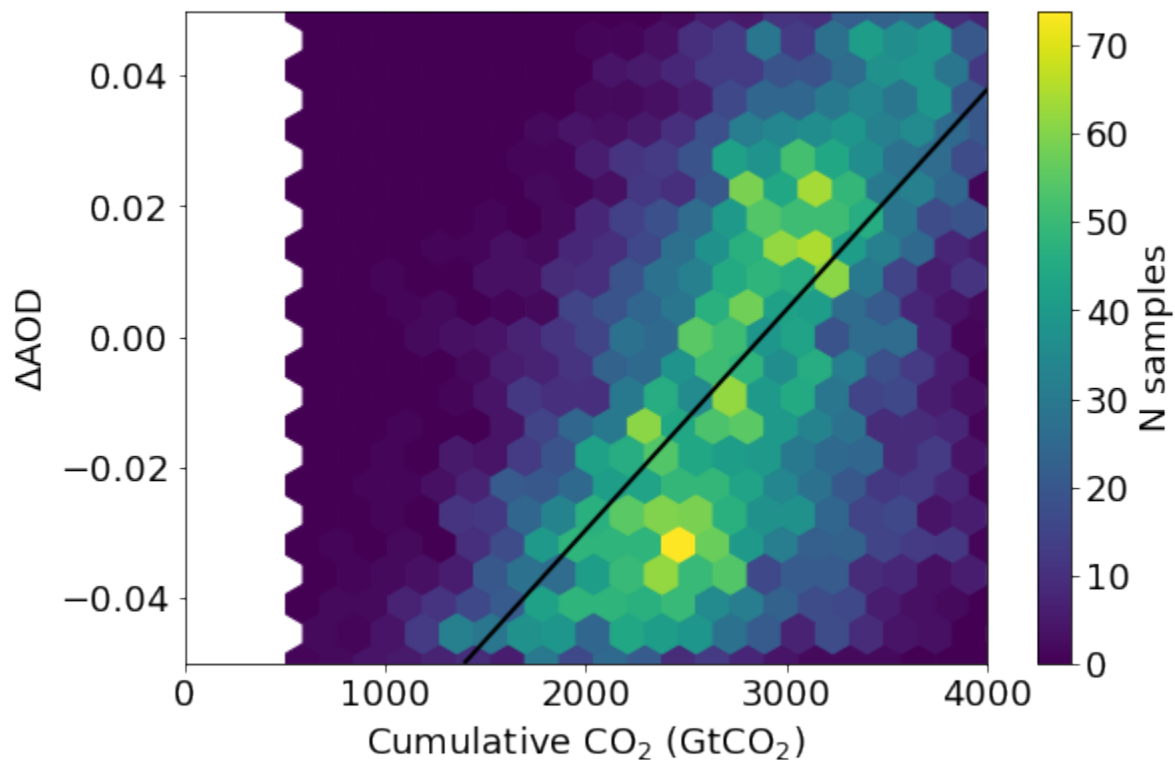
cl = plt.contour(xx, yy, Y_new.reshape(xx.shape), levels = [1.2],
                colors=('k',),linestyles=('-',),linewidths=(2,))

cl=plt.hexbin(new_samples.co2*4000, new_samples.od550aer*0.1-0.05, gridsize=20)

plt.setp(plt.gca(), xlabel='Cumulative CO2$ (GtCO2$)', ylabel='$\Delta$AOD')

plt.colorbar(cl, label='N samples')
plt.setp(plt.gca(), ylim=[-0.05, 0.05], xlim=[0, 4000])

plt.savefig('CMIP6_emulator_sampled.png', transparent=True)
```



```
[ ]:
```

5.7 Create paper emulation figure

```
[1]: import os
    ## Ignore my broken HDF5 install...
    os.putenv("HDF5_DISABLE_VERSION_CHECK", '1')
```

```
[2]: import iris

    from utils import get_bc_ppe_data

    from GCEm import cnn_model, gp_model
    from GCEm.utils import get_random_params

    import iris.quickplot as qplt
    import iris.analysis.maths as imath
    import matplotlib.pyplot as plt
    %matplotlib inline
```

5.7.1 Read in the parameters and data

```
[3]: ppe_params, ppe_aaod = get_bc_ppe_data()

/Users/watson-parris/miniconda3/envs/gcem/lib/python3.8/site-packages/iris/__init__.py:
↳249: IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and
↳will be removed in a future release. Please remove code that sets this property.
    warn_deprecated(msg.format(name))
/Users/watson-parris/miniconda3/envs/gcem/lib/python3.8/site-packages/iris/__init__.py:
↳249: IrisDeprecation: setting the 'Future' property 'netcdf_promote' is deprecated and
↳will be removed in a future release. Please remove code that sets this property.
    warn_deprecated(msg.format(name))
```

```
[4]: ## Ensure the dependent time dimension is last - this is treated as the color 'channel'
    ## ppe_aaod.transpose((0,2,3,1))
    ppe_aaod = ppe_aaod.collapsed('time')[0]
```

```
WARNING:root:Creating guessed bounds as none exist in file
WARNING:root:Creating guessed bounds as none exist in file
WARNING:root:Creating guessed bounds as none exist in file
```

```
[5]: n_test = 5

    X_test, X_train = ppe_params[:n_test], ppe_params[n_test:]
    Y_test, Y_train = ppe_aaod[:n_test], ppe_aaod[n_test:]
```

```
[6]: Y_train
```

```
[6]: <iris 'Cube' of Absorption optical thickness - total 550nm / (1) (job: 34; latitude: 96;
↳longitude: 192)>
```


5.7.2 Setup and run the models

```
[7]: nn_model = cnn_model(X_train, Y_train)
```

```
[8]: nn_model.model.model.summary()
```

Model: "decoder"

Layer (type)	Output Shape	Param #
=====		
input_1 (InputLayer)	[(None, 3)]	0

dense (Dense)	(None, 221184)	884736

reshape (Reshape)	(None, 96, 192, 12)	0

conv2d_transpose (Conv2DTran	(None, 96, 192, 12)	2172

conv2d_transpose_1 (Conv2DTr	(None, 96, 192, 1)	181
=====		

Total params: 887,089
 Trainable params: 887,089
 Non-trainable params: 0

```
[9]: nn_model.train()
```

```
Epoch 1/100
4/4 [=====] - 1s 284ms/step - loss: 0.8766 - val_loss: 0.4800
Epoch 2/100
4/4 [=====] - 0s 115ms/step - loss: 1.0869 - val_loss: 0.4778
Epoch 3/100
4/4 [=====] - 0s 100ms/step - loss: 1.1900 - val_loss: 0.4759
Epoch 4/100
4/4 [=====] - 0s 106ms/step - loss: 1.0656 - val_loss: 0.4741
Epoch 5/100
4/4 [=====] - 0s 99ms/step - loss: 1.2058 - val_loss: 0.4699
Epoch 6/100
4/4 [=====] - 0s 112ms/step - loss: 1.2211 - val_loss: 0.4672
Epoch 7/100
4/4 [=====] - 0s 105ms/step - loss: 1.0138 - val_loss: 0.4624
Epoch 8/100
4/4 [=====] - 0s 101ms/step - loss: 0.9261 - val_loss: 0.4645
Epoch 9/100
4/4 [=====] - 0s 112ms/step - loss: 0.9354 - val_loss: 0.4569
Epoch 10/100
4/4 [=====] - 0s 100ms/step - loss: 1.1767 - val_loss: 0.4494
Epoch 11/100
4/4 [=====] - 0s 101ms/step - loss: 1.1077 - val_loss: 0.4310
Epoch 12/100
4/4 [=====] - 0s 101ms/step - loss: 0.8575 - val_loss: 0.4343
Epoch 13/100
4/4 [=====] - 0s 101ms/step - loss: 0.8162 - val_loss: 0.4476
```

(continues on next page)

(continued from previous page)

```

Epoch 14/100
4/4 [=====] - 0s 102ms/step - loss: 0.8983 - val_loss: 0.4181
Epoch 15/100
4/4 [=====] - 1s 104ms/step - loss: 0.8968 - val_loss: 0.4003
Epoch 16/100
4/4 [=====] - 0s 100ms/step - loss: 0.8450 - val_loss: 0.3989
Epoch 17/100
4/4 [=====] - 0s 103ms/step - loss: 0.7290 - val_loss: 0.4107
Epoch 18/100
4/4 [=====] - 0s 100ms/step - loss: 0.8785 - val_loss: 0.3785
Epoch 19/100
4/4 [=====] - 0s 106ms/step - loss: 0.7236 - val_loss: 0.3918
Epoch 20/100
4/4 [=====] - 0s 103ms/step - loss: 0.8572 - val_loss: 0.3857
Epoch 21/100
4/4 [=====] - 0s 102ms/step - loss: 0.8973 - val_loss: 0.3543
Epoch 22/100
4/4 [=====] - 0s 101ms/step - loss: 0.9261 - val_loss: 0.3282
Epoch 23/100
4/4 [=====] - 0s 102ms/step - loss: 0.5686 - val_loss: 0.3123
Epoch 24/100
4/4 [=====] - 0s 100ms/step - loss: 0.9107 - val_loss: 0.3041
Epoch 25/100
4/4 [=====] - 0s 101ms/step - loss: 0.7340 - val_loss: 0.2931
Epoch 26/100
4/4 [=====] - 0s 100ms/step - loss: 0.7675 - val_loss: 0.2825
Epoch 27/100
4/4 [=====] - 0s 105ms/step - loss: 0.5598 - val_loss: 0.2795
Epoch 28/100
4/4 [=====] - 0s 101ms/step - loss: 0.6441 - val_loss: 0.2640
Epoch 29/100
4/4 [=====] - 0s 100ms/step - loss: 0.8158 - val_loss: 0.2579
Epoch 30/100
4/4 [=====] - 0s 102ms/step - loss: 0.7065 - val_loss: 0.2447
Epoch 31/100
4/4 [=====] - 0s 105ms/step - loss: 0.6948 - val_loss: 0.2351
Epoch 32/100
4/4 [=====] - 0s 125ms/step - loss: 0.7187 - val_loss: 0.2226
Epoch 33/100
4/4 [=====] - 0s 122ms/step - loss: 0.5605 - val_loss: 0.2154
Epoch 34/100
4/4 [=====] - 0s 104ms/step - loss: 0.5716 - val_loss: 0.2117
Epoch 35/100
4/4 [=====] - 0s 107ms/step - loss: 0.6035 - val_loss: 0.1991
Epoch 36/100
4/4 [=====] - 0s 107ms/step - loss: 0.6084 - val_loss: 0.1892
Epoch 37/100
4/4 [=====] - 0s 103ms/step - loss: 0.6416 - val_loss: 0.1786
Epoch 38/100
4/4 [=====] - 0s 100ms/step - loss: 0.4608 - val_loss: 0.1749
Epoch 39/100
4/4 [=====] - 0s 100ms/step - loss: 0.3582 - val_loss: 0.1674

```

(continues on next page)

(continued from previous page)

```

Epoch 40/100
4/4 [=====] - 0s 113ms/step - loss: 0.4616 - val_loss: 0.1607
Epoch 41/100
4/4 [=====] - 0s 117ms/step - loss: 0.5972 - val_loss: 0.1558
Epoch 42/100
4/4 [=====] - 0s 121ms/step - loss: 0.3961 - val_loss: 0.1471
Epoch 43/100
4/4 [=====] - 0s 111ms/step - loss: 0.5399 - val_loss: 0.1404
Epoch 44/100
4/4 [=====] - 0s 100ms/step - loss: 0.4118 - val_loss: 0.1385
Epoch 45/100
4/4 [=====] - 0s 126ms/step - loss: 0.3920 - val_loss: 0.1316
Epoch 46/100
4/4 [=====] - 0s 121ms/step - loss: 0.4126 - val_loss: 0.1237
Epoch 47/100
4/4 [=====] - 1s 247ms/step - loss: 0.3554 - val_loss: 0.1202
Epoch 48/100
4/4 [=====] - 0s 129ms/step - loss: 0.3429 - val_loss: 0.1145
Epoch 49/100
4/4 [=====] - 1s 128ms/step - loss: 0.3157 - val_loss: 0.1075
Epoch 50/100
4/4 [=====] - 1s 117ms/step - loss: 0.4217 - val_loss: 0.1051
Epoch 51/100
4/4 [=====] - 0s 112ms/step - loss: 0.4561 - val_loss: 0.1078
Epoch 52/100
4/4 [=====] - 0s 121ms/step - loss: 0.4475 - val_loss: 0.0932
Epoch 53/100
4/4 [=====] - 0s 123ms/step - loss: 0.2861 - val_loss: 0.0896
Epoch 54/100
4/4 [=====] - 0s 115ms/step - loss: 0.3469 - val_loss: 0.0886
Epoch 55/100
4/4 [=====] - 0s 126ms/step - loss: 0.3332 - val_loss: 0.0920
Epoch 56/100
4/4 [=====] - 0s 117ms/step - loss: 0.3290 - val_loss: 0.0793
Epoch 57/100
4/4 [=====] - 1s 136ms/step - loss: 0.3292 - val_loss: 0.0792
Epoch 58/100
4/4 [=====] - 0s 104ms/step - loss: 0.2757 - val_loss: 0.0760
Epoch 59/100
4/4 [=====] - 0s 110ms/step - loss: 0.2637 - val_loss: 0.0720
Epoch 60/100
4/4 [=====] - 0s 123ms/step - loss: 0.4001 - val_loss: 0.0729
Epoch 61/100
4/4 [=====] - 0s 110ms/step - loss: 0.4507 - val_loss: 0.0666
Epoch 62/100
4/4 [=====] - 1s 131ms/step - loss: 0.3879 - val_loss: 0.0695
Epoch 63/100
4/4 [=====] - 0s 106ms/step - loss: 0.2597 - val_loss: 0.0610
Epoch 64/100
4/4 [=====] - 0s 111ms/step - loss: 0.3288 - val_loss: 0.0725
Epoch 65/100
4/4 [=====] - 0s 102ms/step - loss: 0.2841 - val_loss: 0.0575

```

(continues on next page)

(continued from previous page)

```

Epoch 66/100
4/4 [=====] - 0s 101ms/step - loss: 0.2730 - val_loss: 0.0544
Epoch 67/100
4/4 [=====] - 0s 101ms/step - loss: 0.2694 - val_loss: 0.0517
Epoch 68/100
4/4 [=====] - 0s 110ms/step - loss: 0.3846 - val_loss: 0.0573
Epoch 69/100
4/4 [=====] - 0s 100ms/step - loss: 0.3429 - val_loss: 0.0528
Epoch 70/100
4/4 [=====] - 0s 106ms/step - loss: 0.2360 - val_loss: 0.0478
Epoch 71/100
4/4 [=====] - 0s 100ms/step - loss: 0.3651 - val_loss: 0.0486
Epoch 72/100
4/4 [=====] - 0s 107ms/step - loss: 0.2351 - val_loss: 0.0469
Epoch 73/100
4/4 [=====] - 0s 101ms/step - loss: 0.2993 - val_loss: 0.0476
Epoch 74/100
4/4 [=====] - 0s 105ms/step - loss: 0.2739 - val_loss: 0.0523
Epoch 75/100
4/4 [=====] - 0s 103ms/step - loss: 0.3909 - val_loss: 0.0438
Epoch 76/100
4/4 [=====] - 0s 109ms/step - loss: 0.2526 - val_loss: 0.0423
Epoch 77/100
4/4 [=====] - 0s 103ms/step - loss: 0.2566 - val_loss: 0.0485
Epoch 78/100
4/4 [=====] - 0s 101ms/step - loss: 0.2969 - val_loss: 0.0411
Epoch 79/100
4/4 [=====] - 0s 99ms/step - loss: 0.2602 - val_loss: 0.0413
Epoch 80/100
4/4 [=====] - 0s 126ms/step - loss: 0.1863 - val_loss: 0.0382
Epoch 81/100
4/4 [=====] - 0s 101ms/step - loss: 0.3251 - val_loss: 0.0394
Epoch 82/100
4/4 [=====] - 0s 109ms/step - loss: 0.3038 - val_loss: 0.0376
Epoch 83/100
4/4 [=====] - 0s 105ms/step - loss: 0.3729 - val_loss: 0.0405
Epoch 84/100
4/4 [=====] - 0s 107ms/step - loss: 0.2542 - val_loss: 0.0389
Epoch 85/100
4/4 [=====] - 0s 109ms/step - loss: 0.2379 - val_loss: 0.0403
Epoch 86/100
4/4 [=====] - 0s 117ms/step - loss: 0.2620 - val_loss: 0.0350
Epoch 87/100
4/4 [=====] - 0s 108ms/step - loss: 0.2459 - val_loss: 0.0356
Epoch 88/100
4/4 [=====] - 0s 107ms/step - loss: 0.1902 - val_loss: 0.0368
Epoch 89/100
4/4 [=====] - 1s 134ms/step - loss: 0.3153 - val_loss: 0.0362
Epoch 90/100
4/4 [=====] - 0s 116ms/step - loss: 0.3270 - val_loss: 0.0371
Epoch 91/100
4/4 [=====] - 1s 123ms/step - loss: 0.3058 - val_loss: 0.0393

```

(continues on next page)

(continued from previous page)

```

Epoch 92/100
4/4 [=====] - 0s 115ms/step - loss: 0.3535 - val_loss: 0.0359
Epoch 93/100
4/4 [=====] - 0s 102ms/step - loss: 0.3056 - val_loss: 0.0379
Epoch 94/100
4/4 [=====] - 0s 106ms/step - loss: 0.3331 - val_loss: 0.0354
Epoch 95/100
4/4 [=====] - 0s 99ms/step - loss: 0.2074 - val_loss: 0.0364
Epoch 96/100
4/4 [=====] - 0s 99ms/step - loss: 0.1920 - val_loss: 0.0364
Epoch 97/100
4/4 [=====] - 0s 101ms/step - loss: 0.3333 - val_loss: 0.0329
Epoch 98/100
4/4 [=====] - 0s 99ms/step - loss: 0.2105 - val_loss: 0.0340
Epoch 99/100
4/4 [=====] - 0s 105ms/step - loss: 0.3683 - val_loss: 0.0355
Epoch 100/100
4/4 [=====] - 0s 100ms/step - loss: 0.3073 - val_loss: 0.0361

```

```
[10]: ## Linear model: 0.3566 - val_loss: 0.0867
```

```
[11]: nn_prediction, _ = nn_model.predict(X_test.values)
```

```
[12]: gp_model_ = gp_model(X_train, Y_train, kernel=['Bias', 'Linear'])
gp_model_.train()
```

```
[13]: gp_prediction, _ = gp_model_.predict(X_test.values)
```

```
[14]: import matplotlib
import cartopy.crs as ccrs
import iris.plot as iplt

plt.figure(figsize=(30, 10))
matplotlib.rcParams['font.size'] = 24

plt.subplot(2,3,1, projection=ccrs.Mollweide())
plt.annotate("(a)", (0.,1.), xycoords='axes fraction')
iplt.pcolormesh(imath.log10(Y_test[0]), vmin=-4, vmax=-1)
plt.gca().set_title('Truth')
plt.gca().coastlines()

plt.subplot(2,3,2, projection=ccrs.Mollweide())
plt.annotate("(b)", (0.,1.), xycoords='axes fraction')
iplt.pcolormesh(imath.log10(gp_prediction[0]), vmin=-4, vmax=-1)
plt.gca().set_title('GP')
plt.gca().coastlines()

plt.subplot(2,3,3, projection=ccrs.Mollweide())
plt.annotate("(c)", (0.,1.), xycoords='axes fraction')
im=iplt.pcolormesh(imath.log10(nn_prediction[0]), vmin=-4, vmax=-1)

```

(continues on next page)

(continued from previous page)

```

plt.gca().set_title('CNN')
plt.colorbar(im, fraction=0.046, pad=0.04, label='log(AAOD)')
plt.gca().coastlines()

plt.subplot(2,3,5, projection=ccrs.Mollweide())
plt.annotate("(d)", (0.,1.), xycoords='axes fraction')
iplt.pcolormesh((gp_prediction.collapsed(['sample'], iris.analysis.MEAN)-Y_test.
↳collapsed(['job'], iris.analysis.MEAN)), cmap='RdBu_r', vmin=-0.001, vmax=0.001)
plt.gca().coastlines()
plt.gca().set_title('Difference')

plt.subplot(2,3,6, projection=ccrs.Mollweide())
plt.annotate("(e)", (0.,1.), xycoords='axes fraction')
im=iplt.pcolormesh((nn_prediction.collapsed(['sample'], iris.analysis.MEAN)-Y_test.
↳collapsed(['job'], iris.analysis.MEAN)), cmap='RdBu_r', vmin=-1e-3, vmax=1e-3)
cb = plt.colorbar(im, fraction=0.046, pad=0.04)
cb.ax.set_yticklabels(["{: .2e} ".format(i) for i in cb.get_ticks()]) ## set ticks of your
↳format
plt.gca().coastlines()
plt.gca().set_title('Difference')

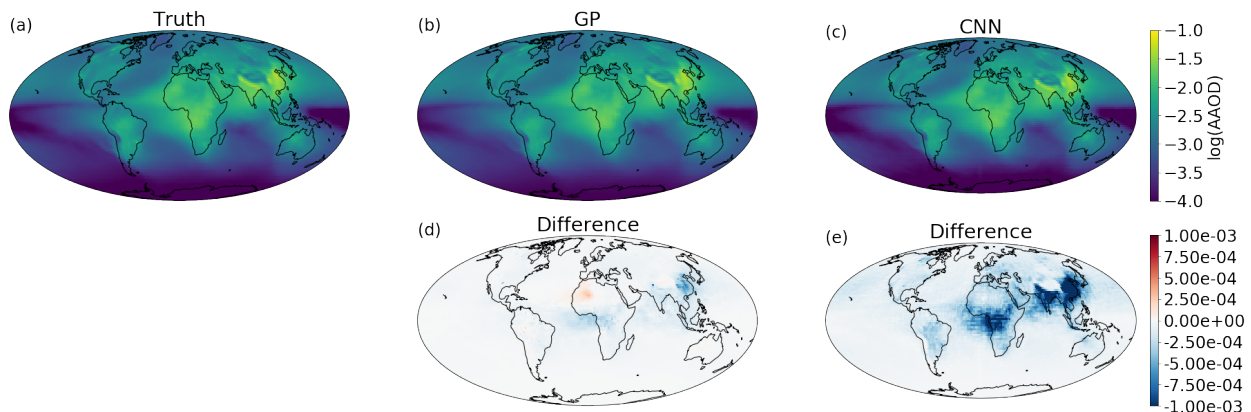
plt.savefig('BCPPE_emulator_paper.png', transparent=True)

```

```

/Users/watson-parris/miniconda3/envs/gcem/lib/python3.8/site-packages/iris/coords.py:
↳1410: UserWarning: Collapsing a non-contiguous coordinate. Metadata may not be fully
↳descriptive for 'sample'.
warnings.warn(msg.format(self.name()))
/Users/watson-parris/miniconda3/envs/gcem/lib/python3.8/site-packages/iris/coords.py:
↳1410: UserWarning: Collapsing a non-contiguous coordinate. Metadata may not be fully
↳descriptive for 'sample'.
warnings.warn(msg.format(self.name()))

```



[15]:

```

COLOR = 'white'
matplotlib.rcParams['text.color'] = COLOR
matplotlib.rcParams['axes.labelcolor'] = COLOR
matplotlib.rcParams['xtick.color'] = COLOR
matplotlib.rcParams['ytick.color'] = COLOR
matplotlib.rcParams['font.size'] = 20

```

(continues on next page)

(continued from previous page)

```

plt.figure(figsize=(30, 10))

plt.subplot(2,3,1, projection=ccrs.Mollweide())
plt.annotate("(a)", (0.,1.), xycoords='axes fraction')
iplt.pcolormesh(imath.log10(Y_test[0]), vmin=-4, vmax=-1)
plt.gca().set_title('Truth')
plt.gca().coastlines()

plt.subplot(2,3,2, projection=ccrs.Mollweide())
plt.annotate("(b)", (0.,1.), xycoords='axes fraction')
iplt.pcolormesh(imath.log10(gp_prediction[0]), vmin=-4, vmax=-1)
plt.gca().set_title('GP')
plt.gca().coastlines()

plt.subplot(2,3,3, projection=ccrs.Mollweide())
plt.annotate("(c)", (0.,1.), xycoords='axes fraction')
im=iplt.pcolormesh(imath.log10(nn_prediction[0]), vmin=-4, vmax=-1)
plt.gca().set_title('CNN')
plt.colorbar(im, fraction=0.046, pad=0.04, label='log(AAOD)')
plt.gca().coastlines()

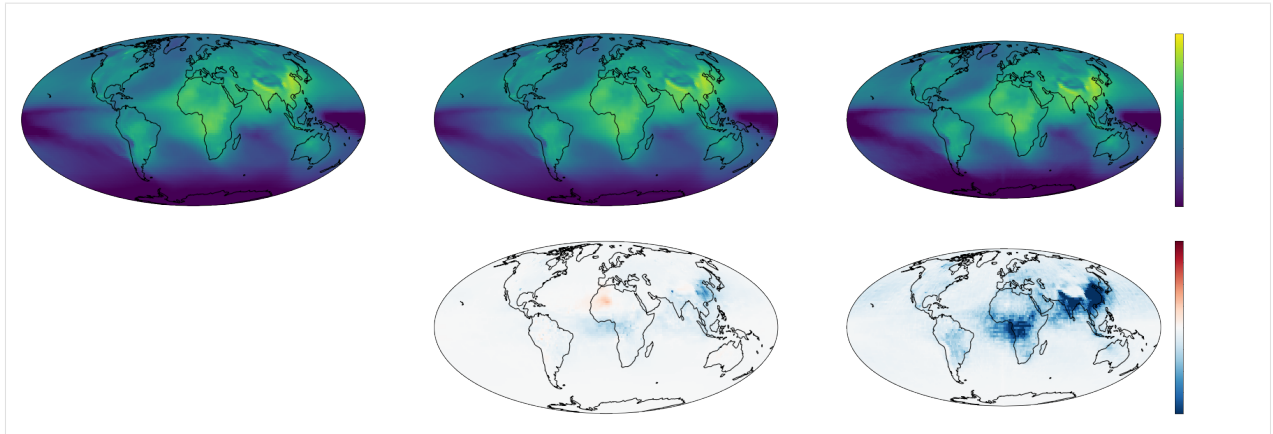
plt.subplot(2,3,5, projection=ccrs.Mollweide())
plt.annotate("(d)", (0.,1.), xycoords='axes fraction')
iplt.pcolormesh((gp_prediction.collapsed(['sample'], iris.analysis.MEAN)-Y_test.
↳collapsed(['job'], iris.analysis.MEAN)), cmap='RdBu_r', vmin=-0.001, vmax=0.001)
plt.gca().coastlines()
plt.gca().set_title('Difference')

plt.subplot(2,3,6, projection=ccrs.Mollweide())
plt.annotate("(e)", (0.,1.), xycoords='axes fraction')
im=iplt.pcolormesh((nn_prediction.collapsed(['sample'], iris.analysis.MEAN)-Y_test.
↳collapsed(['job'], iris.analysis.MEAN)), cmap='RdBu_r', vmin=-1e-3, vmax=1e-3)
cb=plt.colorbar(im, fraction=0.046, pad=0.04)
cb.ax.set_yticklabels(["{:.2e}".format(i) for i in cb.get_ticks()]) ## set ticks of your
↳format
plt.gca().coastlines()
plt.gca().set_title('Difference')

plt.savefig('BCPPE_emulator_talk.png', transparent=True)

/Users/watson-parris/miniconda3/envs/gcem/lib/python3.8/site-packages/iris/coords.py:
↳1410: UserWarning: Collapsing a non-contiguous coordinate. Metadata may not be fully
↳descriptive for 'sample'.
warnings.warn(msg.format(self.name()))
/Users/watson-parris/miniconda3/envs/gcem/lib/python3.8/site-packages/iris/coords.py:
↳1410: UserWarning: Collapsing a non-contiguous coordinate. Metadata may not be fully
↳descriptive for 'sample'.
warnings.warn(msg.format(self.name()))

```



[]:

API REFERENCE

This page provides an auto-generated summary of xarray's API. For more details and examples, refer to the relevant chapters in the main part of the documentation.

6.1 Top-level functions

This provides the main interface for ESEm and should be the starting point for most users.

<i>gp_model</i>	Create a Gaussian process (GP) based emulator with provided <i>training_params</i> (X) and <i>training_data</i> (Y) which assumes independent inputs (and outputs).
<i>cnn_model</i>	Create a simple two layer Convolutional Neural Network Emulator using Keras.
<i>rf_model</i>	Create a simple Random Forest Emulator using sklearn.

6.1.1 esem.gp_model

`esem.gp_model(training_params, training_data, data_processors=None, kernel=None, kernel_op='add', active_dims=None, noise_variance=1.0, name="", gpu=0)`

Create a Gaussian process (GP) based emulator with provided *training_params* (X) and *training_data* (Y) which assumes independent inputs (and outputs).

The *kernel* is a key parameter in GP emulation and care should be taken in choosing it.

Parameters

- **training_params** (DataFrame) – The training parameters
- **training_data** (xarray.DataArray or iris.Cube or *array_like*) – The training data - the leading dimension should represent training samples
- **data_processors** (list of *esem.data_processors.DataProcessor*) – A list of 'DataProcessor' to apply to the data transparently before training. Model output will be untransformed before being returned from the Emulator.
- **kernel** (gpflow.kernels.Kernel or list of str or None) – The GP kernel to use. A GPFlow kernel can be specified directly, or a list of kernel names can be provided which will be initialised using the default values (of the correct shape) and combined using *kernel_op*. Alternatively no kernel can be specified and a default will be used.
- **kernel_op** ({'add', 'mul'}) – The operation to perform in order to combine the specified *kernel*'s. Only used if *kernel* is a list of strings.

- **noise_variance** (float) – The noise variance to initialise the GP regression model
- **active_dims** (list of int or slice or None) – The dimensions to train the GP over (by default all of the dimensions are used)
- **name** (str) – An optional name for the emulator
- **gpu** (int) – The GPU to use (only applicable for multi-GPU) machines

Returns Emulator – An esim emulator object which can be trained and sampled from

6.1.2 esim.cnn_model

`esim.cnn_model(training_params, training_data, data_processors=None, filters=12, learning_rate=0.001, decay=0.01, kernel_size=(3, 5), loss='mean_squared_error', activation='tanh', optimizer='RMSprop', name='', gpu=0)`

Create a simple two layer Convolutional Neural Network Emulator using Keras.

Note that X should include both the train and validation data

Parameters

- **training_params** (pd.DataFrame) – The training parameters
- **training_data** (xarray.DataArray or iris.cube.Cube or *array_like*) – The training data - the leading dimension should represent training samples
- **data_processors** (list of `esim.data_processors.DataProcessor`) – A list of *DataProcessor* to apply to the data transparently before training. Model output will be un-transformed before being returned from the Emulator.
- **filters** (int) – The dimensionality of the first convolutional layer output space
- **learning_rate** (float) – The learning rate to use with the chosen optimizer
- **decay** (float) – Any decay to apply to the learning rate
- **kernel_size** (tuple of int) – The convolutional kernel size
- **loss** (str) – The loss function to train against (see <https://keras.io/api/losses/>)
- **activation** (str) – The activation function to use in the final CNN layer (see <https://keras.io/api/layers/activations/>)
- **optimizer** ({'RMSprop', 'Adam'}) – The optimizer to train the model with
- **name** (str) – An optional name for the emulator
- **gpu** (int) – The GPU to use (only applicable for multi-GPU) machines

Returns Emulator – An esim emulator object which can be trained and sampled from

Notes

The Keras model is compiled but not trained until *train* is called on the returned *Emulator* object.

6.1.3 `esem.rf_model`

`esem.rf_model(training_params, training_data, data_processors=None, name='', gpu=0, *args, **kwargs)`

Create a simple Random Forest Emulator using sklearn.

Note that because a Random Forest is just a recursive binary partition over the training data, there is no need to normalize/standardize the inputs.

i.e. At least in theory, Random Forests are invariant to monotonic transformations of the independent variables

Parameters

- **training_params** (pd.DataFrame) – The training parameters
- **training_data** (xarray.DataArray or iris.cube.Cube or *array_like*) – The training data - the leading dimension should represent training samples
- **data_processors** (list of `esem.data_processors.DataProcessor`) – A list of *DataProcessor* to apply to the data transparently before training. Model output will be untransformed before being returned from the Emulator.
- **name** (str) – An optional name for the emulator
- **gpu** (int) – The GPU to use (only applicable for multi-GPU) machines
- **args** (list) – List of optional arguments for `sklearn.ensemble.RandomForestRegressor`
- **kwargs** (dict) – Dict of optional keyword arguments for `sklearn.ensemble.RandomForestRegressor`

Returns *Emulator* – An *esem* emulator object which can be trained and sampled from

6.2 Emulator

<i>Emulator</i>	A class wrapping a statistical emulator
<i>Emulator.train</i>	Train on the training data
<i>Emulator.predict</i>	Make a prediction using a trained emulator
<i>Emulator._predict</i>	The (internal) predict interface used by e.g., a sampler.
<i>Emulator.batch_stats</i>	Return mean and standard deviation in model predictions over samples, without storing the intermediate predictions in memory to enable evaluating large models over more samples than could fit in memory

6.2.1 `esem.emulator.Emulator`

class `esem.emulator.Emulator`(*model*, *training_params*, *training_data*, *name*="", *gpu*=0)

A class wrapping a statistical emulator

training_data

A wrapped representation of the training data

Type `esem.wrappers.DataWrapper`

model

The underlying model which performs the emulation

Type `ModelAdaptor`

name

A human-readable name for the model

Type `str`

__init__(*model*, *training_params*, *training_data*, *name*="", *gpu*=0)

Parameters

- **model** (`ModelAdaptor`) – The (compiled but not trained) model to be wrapped
- **training_params** (`pd.DataFrame` or *array-like*) – The training parameters (X)
- **training_data** (`esem.wrappers.DataWrapper` or `xarray.DataArray` or `iris.Cube` or *array-like*) – The training data - the leading dimension should represent training samples (Y)
- **name** (`str`) – Human readable name for the model
- **gpu** (`int`) – The machine GPU to assign this model to

Methods

__init__(*model*, *training_params*, *training_data*)

Parameters

- **model** (`ModelAdaptor`) – The (compiled but not trained) model to be wrapped

batch_stats(*sample_points*[, *batch_size*])

Return mean and standard deviation in model predictions over samples, without storing the intermediate predictions in memory to enable evaluating large models over more samples than could fit in memory

predict(*x*, **args*, ***kwargs*)

Make a prediction using a trained emulator

train([*verbose*])

Train on the training data

6.2.2 `esem.emulator.Emulator.train`

`Emulator.train(verbose=False, **kwargs)`

Train on the training data

Parameters `verbose` (bool) – Print verbose training output to screen

6.2.3 `esem.emulator.Emulator.predict`

`Emulator.predict(x, *args, **kwargs)`

Make a prediction using a trained emulator

Parameters

- `x` (pd.DataFrame or *array-like*) – The points at which to make predictions from the model
- `args` – The specific arguments needed for prediction with this model
- `kwargs` – Any keyword arguments that might need to be passed through to the model

Returns Emulated prediction and variance with the same type as ``self.training_data``

6.2.4 `esem.emulator.Emulator._predict`

`Emulator._predict(x, *args, **kwargs)`

The (internal) predict interface used by e.g., a sampler. It is still in tf but has been post-processed to allow comparison with obs.

Parameters

- `x` (*array-like*) – The points at which to make predictions from the model
- `args` – The specific arguments needed for prediction with this model
- `kwargs` – Any keyword arguments that might need to be passed through to the model

Returns Emulated prediction and variance as either `np.ndarray` or `tf.Tensor`

6.2.5 `esem.emulator.Emulator.batch_stats`

`Emulator.batch_stats(sample_points, batch_size=1)`

Return mean and standard deviation in model predictions over samples, without storing the intermediate predictions in memory to enable evaluating large models over more samples than could fit in memory

Parameters

- `sample_points` (pd.DataFrame or *array-like*) – The parameter values at which to sample the emulator
- `batch_size` (int) – The number of samples to calculate in each batch. This can be optimised to fill the available (GPU) memory

Returns The batch mean and standard deviation with the same type as ``self.training_data``

6.3 Sampler

This class defines the sampling interface currently used by the ABC and MCMC sampling implementations.

<i>Sampler</i>	A class that efficiently samples a Model object for posterior inference
<i>Sampler.sample</i>	This is the call that does the actual inference.

6.3.1 `esem.sampler.Sampler`

```
class esim.sampler.Sampler(model, obs, obs_uncertainty=0.0, interann_uncertainty=0.0,
                           repres_uncertainty=0.0, struct_uncertainty=0.0, abs_obs_uncertainty=0.0,
                           abs_interann_uncertainty=0.0, abs_repres_uncertainty=0.0,
                           abs_struct_uncertainty=0.0)
```

A class that efficiently samples a Model object for posterior inference

```
__init__(model, obs, obs_uncertainty=0.0, interann_uncertainty=0.0, repres_uncertainty=0.0,
          struct_uncertainty=0.0, abs_obs_uncertainty=0.0, abs_interann_uncertainty=0.0,
          abs_repres_uncertainty=0.0, abs_struct_uncertainty=0.0)
```

Parameters

- **model** (*esem.emulator.Emulator*)
- **obs** (*iris.cube.Cube* or *array-like*) – The objective
- **obs_uncertainty** (float) – Fractional, relative (1 sigma) uncertainty in observations
- **repres_uncertainty** (float) – Fractional, relative (1 sigma) uncertainty due to the spatial and temporal representitiveness of the observations
- **interann_uncertainty** (float) – Fractional, relative (1 sigma) uncertainty introduced when using a model run for a year other than that the observations were measured in.
- **struct_uncertainty** (float) – Fractional, relative (1 sigma) uncertainty in the model itself.
- **abs_obs_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty in observations
- **abs_repres_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty due to the spatial and temporal representitiveness of the observations
- **abs_interann_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty introduced when using a model run for a year other than that the observations were measured in.
- **abs_struct_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty in the model itself.

Methods

`__init__(model, obs[, obs_uncertainty, ...])`

Parameters

- **model** (*esem.emulator.Emulator*)

`sample([prior_x, n_samples])`

This is the call that does the actual inference.

6.3.2 `esem.sampler.Sampler.sample`

`Sampler.sample(prior_x=None, n_samples=1)`

This is the call that does the actual inference.

It should call `model.sample` over the prior, compare with the objective, and then output samples from the posterior distribution

Parameters

- **prior_x** (`tensorflow_probability.distribution`) – The distribution to sample parameters from. By default it will uniformly sample the unit N-D hypercube
- **n_samples** (`int`) – The number of samples to draw

Returns `np.array` – Array of samples

6.3.3 `MCMCSampler`

`MCMCSampler`

Sample from the posterior using the TensorFlow Markov-Chain Monte-Carlo (MCMC) sampling tools.

`MCMCSampler.sample`

This is the call that does the actual inference.

`esem.sampler.MCMCSampler`

class `esem.sampler.MCMCSampler(model, obs, **kwargs)`

Sample from the posterior using the TensorFlow Markov-Chain Monte-Carlo (MCMC) sampling tools. It uses a HamiltonianMonteCarlo kernel.

Notes

Note that NaN observations will create ill-defined likelihoods.

`__init__(model, obs, **kwargs)`

Parameters

- **model** (*esem.emulator.Emulator*)
- **obs** (`iris.cube.Cube` or *array-like*) – The objective
- **obs_uncertainty** (`float`) – Fractional, relative (1 sigma) uncertainty in observations

- **repres_uncertainty** (float) – Fractional, relative (1 sigma) uncertainty due to the spatial and temporal representitiveness of the observations
- **interann_uncertainty** (float) – Fractional, relative (1 sigma) uncertainty introduced when using a model run for a year other than that the observations were measured in.
- **struct_uncertainty** (float) – Fractional, relative (1 sigma) uncertainty in the model itself.
- **abs_obs_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty in observations
- **abs_repres_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty due to the spatial and temporal representitiveness of the observations
- **abs_interann_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty introduced when using a model run for a year other than that the observations were measured in.
- **abs_struct_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty in the model itself.

Methods

`__init__(model, obs, **kwargs)`

Parameters

- **model** (*esem.emulator.Emulator*)

`sample([prior_x, n_samples, kernel_kwargs, ...])` This is the call that does the actual inference.

`esem.sampler.MCMCSampler.sample`

`MCMCSampler.sample(prior_x=None, n_samples=1, kernel_kwargs=None, mcmc_kwargs=None)`

This is the call that does the actual inference.

It should call `model.sample` over the prior, compare with the objective, and then output a posterior distribution

Parameters

- **prior_x** (*tensorflow_probability.distribution*) – The distribution to sample parameters from. By default it will uniformly sample the unit N-D hypercube
- **n_samples** (int) – The number of samples to draw
- **kernel_kwargs** (dict) – kwargs for the MCMC kernel
- **mcmc_kwargs** (dict) – kwargs for the MCMC sampler

Returns `np.array` – Array of samples

6.3.4 ABCSampler

<code>ABCSampler</code>	Sample from the posterior using Approximate Bayesian Computation (ABC).
<code>ABCSampler.sample</code>	Sample the emulator over <i>prior_x</i> and compare with the observations, returning <i>n_samples</i> of the posterior distribution (those points for which the model is compatible with the observations).
<code>ABCSampler.get_implausibility</code>	Calculate the implausibility of the provided sample points, optionally in batches.
<code>ABCSampler.batch_constrain</code>	Constrain the supplied sample points based on the tolerance threshold, optionally in batches.

esem.abc_sampler.ABCSampler

```
class esim.abc_sampler.ABCSampler(model, obs, obs_uncertainty=0.0, interann_uncertainty=0.0,
                                repres_uncertainty=0.0, struct_uncertainty=0.0,
                                abs_obs_uncertainty=0.0, abs_interann_uncertainty=0.0,
                                abs_repres_uncertainty=0.0, abs_struct_uncertainty=0.0)
```

Sample from the posterior using Approximate Bayesian Computation (ABC). This is a style of rejection sampling.

Notes

Note that emulator samples compared to NaN observations are always treated as ‘plausible’.

```
__init__(model, obs, obs_uncertainty=0.0, interann_uncertainty=0.0, repres_uncertainty=0.0,
          struct_uncertainty=0.0, abs_obs_uncertainty=0.0, abs_interann_uncertainty=0.0,
          abs_repres_uncertainty=0.0, abs_struct_uncertainty=0.0)
```

Parameters

- **model** (`esem.emulator.Emulator`)
- **obs** (`iris.cube.Cube` or *array-like*) – The objective
- **obs_uncertainty** (float) – Fractional, relative (1 sigma) uncertainty in observations
- **repres_uncertainty** (float) – Fractional, relative (1 sigma) uncertainty due to the spatial and temporal representitiveness of the observations
- **interann_uncertainty** (float) – Fractional, relative (1 sigma) uncertainty introduced when using a model run for a year other than that the observations were measured in.
- **struct_uncertainty** (float) – Fractional, relative (1 sigma) uncertainty in the model itself.
- **abs_obs_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty in observations
- **abs_repres_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty due to the spatial and temporal representitiveness of the observations
- **abs_interann_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty introduced when using a model run for a year other than that the observations were measured in.
- **abs_struct_uncertainty** (float) – Fractional, absolute (1 sigma) uncertainty in the model itself.

Methods

<code>__init__(model, obs[, obs_uncertainty, ...])</code>	<p>Parameters</p> <ul style="list-style-type: none"> • model (<i>esem.emulator.Emulator</i>)
<code>batch_constrain(sample_points[, tolerance, ...])</code>	Constrain the supplied sample points based on the tolerance threshold, optionally in batches.
<code>get_implausibility(sample_points[, batch_size])</code>	Calculate the implausibility of the provided sample points, optionally in batches.
<code>sample([prior_x, n_samples, tolerance, ...])</code>	Sample the emulator over <i>prior_x</i> and compare with the observations, returning <i>n_samples</i> of the posterior distribution (those points for which the model is compatible with the observations).

esem.abc_sampler.ABCSampler.sample

ABCSampler.**sample**(*prior_x=None, n_samples=1, tolerance=0.0, threshold=3.0*)

Sample the emulator over *prior_x* and compare with the observations, returning *n_samples* of the posterior distribution (those points for which the model is compatible with the observations).

Parameters

- **prior_x** (*tensorflow_probability.distribution* or *None*) – The distribution to sample parameters from. By default it will uniformly sample the unit N-D hypercube
- **n_samples** (*int*) – The number of samples to draw
- **tolerance** (*float*) – The fraction of samples which are allowed to be over the threshold
- **threshold** (*float*) – The number of standard deviations a sample is allowed to be away from the obs

Returns *ndarray[n_samples]* – Array of samples conforming to the specified tolerance and threshold

esem.abc_sampler.ABCSampler.get_implausibility

ABCSampler.**get_implausibility**(*sample_points, batch_size=1*)

Calculate the implausibility of the provided sample points, optionally in batches.

Note this calculates an array of shape (n_sample_points, n_obs) and so can easily exceed available memory if not used carefully.

Parameters

- **sample_points** (*ndarray* or *DataFrame*) – The sample points to calculate the implausibility for
- **batch_size** (*int*) – The size of the batches in which to calculate the implausibility (useful for large samples)

Returns *Cube* – A cube of the implausibility of each sample against each observation

esem.abc_sampler.ABCSampler.batch_constrain

ABCSampler.batch_constrain(*sample_points*, *tolerance*=0.0, *threshold*=3.0, *batch_size*=1)

Constrain the supplied sample points based on the tolerance threshold, optionally in batches.

Return a boolean array indicating if each sample meets the implausibility criteria:

$$I < T$$

Return True (for a sample) if the number of implausibility measures greater than the threshold is less than or equal to the tolerance

Parameters

- **sample_points** (ndarray) – An array of sample points which are to be emulated and compared with the observations
- **tolerance** (float) – The fraction of samples which are allowed to be over the threshold
- **threshold** (float) – The number of standard deviations a sample is allowed to be away from the obs
- **batch_size** (int) – The size of the batches in which to perform the constraining (useful for large samples)

Returns ndarray – A boolean array which is true where the (emulated) samples are compatible with the observations and false otherwise

6.4 Wrappers

<i>ProcessWrapper</i>	This class handles applying any data pre- and post-processing by any provided DataProcessor
<i>DataWrapper</i>	Provide a unified interface for numpy arrays, Iris Cube's and xarray DataArrays.
<i>DataWrapper.name</i>	
<i>DataWrapper.data</i>	
<i>DataWrapper.dtype</i>	
<i>DataWrapper.wrap</i>	Wrap back in a cube if one was provided
<i>CubeWrapper</i>	
<i>DataArrayWrapper</i>	

6.4.1 `esem.wrappers.ProcessWrapper`

class `esem.wrappers.ProcessWrapper`(*data*, *data_processors=None*)

This class handles applying any data pre- and post-processing by any provided `DataProcessor`

`__init__`(*data*, *data_processors=None*)

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__</code> (<i>data</i> [, <i>data_processors</i>])	Initialize self.
<code>post_process</code> (<i>mean</i> , <i>variance</i>)	Any necessary reshaping or un-weightings are performed here
<code>pre_process</code> (<i>data</i>)	Any necessary rescaling or weightings are performed here

Attributes

<code>data</code>

6.4.2 `esem.wrappers.DataWrapper`

class `esem.wrappers.DataWrapper`(*data*, *data_processors=None*)

Provide a unified interface for numpy arrays, Iris Cube's and xarray DataArrays. Emulation outputs will be provided based on the provided input type, preserving appropriate metadata.

`__init__`(*data*, *data_processors=None*)

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__</code> (<i>data</i> [, <i>data_processors</i>])	Initialize self.
<code>name</code> ()	
<code>wrap</code> (<i>data</i> [, <i>name_prefix</i>])	Wrap back in a cube if one was provided

Attributes

<code>data</code>
<code>dtype</code>

6.4.3 `esem.wrappers.DataWrapper.name`

`DataWrapper.name()`

6.4.4 `esem.wrappers.DataWrapper.data`

property `DataWrapper.data`

6.4.5 `esem.wrappers.DataWrapper.dtype`

property `DataWrapper.dtype`

6.4.6 `esem.wrappers.DataWrapper.wrap`

`DataWrapper.wrap(data, name_prefix='Emulated ')`

Wrap back in a cube if one was provided

Parameters

- **data** (*np.array*) – Model output to wrap
- **name_prefix** (*str*) –

Returns

6.4.7 `esem.wrappers.CubeWrapper`

class `esem.wrappers.CubeWrapper(cube, data_processors=None)`

`__init__` (*cube, data_processors=None*)

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__</code> (cube[, data_processors])	Initialize self.
<code>name</code> ()	
<code>wrap</code> (data[, name_prefix])	Wrap back in a cube if one was provided

Attributes

data
dtype

6.4.8 `esem.wrappers.DataArrayWrapper`

`class` `esem.wrappers.DataArrayWrapper`(*dataarray*, *data_processors=None*)

`__init__`(*dataarray*, *data_processors=None*)
Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__</code> (<i>dataarray</i> [, <i>data_processors</i>])	Initialize self.
<code>name</code> ()	
<code>wrap</code> (<i>data</i> [, <i>name_prefix</i>])	Wrap back in a <code>xr.DataArray</code> if one was provided

Attributes

data
dtype

6.5 `ModelAdaptor`

<i>ModelAdaptor</i>	Provides a unified interface for all emulation engines within ESEm.
<i>SKLearnModel</i>	A wrapper around scikit-learn models.
<i>KerasModel</i>	A wrapper around Keras models
<i>GPFlowModel</i>	A wrapper around GPFlow regression models

6.5.1 `esem.model_adaptor.ModelAdaptor`

class `esem.model_adaptor.ModelAdaptor(model)`

Provides a unified interface for all emulation engines within ESEm. Concrete classes must implement both `train()` and `predict()` methods.

See the [API documentation](#) for a list of concrete classes implementing this interface.

`__init__(model)`

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__(model)</code>	Initialize self.
<code>predict(*args, **kwargs)</code>	This is either the tf model which can be called directly, or a generator over the model.predict (in tf, so it's quick).
<code>train(training_params, training_data[, verbose])</code>	Train on the training data :return:

6.5.2 `esem.model_adaptor.SKLearnModel`

class `esem.model_adaptor.SKLearnModel(model)`

A wrapper around [scikit-learn](#) models.

`__init__(model)`

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__(model)</code>	Initialize self.
<code>predict(*args, **kwargs)</code>	This is either the tf model which can be called directly, or a generator over the model.predict (in tf, so it's quick).
<code>train(training_params, training_data[, verbose])</code>	Train the RF model.

6.5.3 `esem.model_adaptor.KerasModel`

class `esem.model_adaptor.KerasModel(model)`

A wrapper around [Keras](#) models

`__init__(model)`

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__(model)</code>	Initialize self.
<code>predict(*args, **kwargs)</code>	This is either the tf model which can be called directly, or a generator over the model.predict (in tf, so it's quick).
<code>train(training_params, training_data[, ...])</code>	Train the Keras model.

6.5.4 `esem.model_adaptor.GPFlowModel`

class `esem.model_adaptor.GPFlowModel(model)`

A wrapper around `GPFlow` regression models

`__init__(model)`

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__(model)</code>	Initialize self.
<code>predict(*args, **kwargs)</code>	This is either the tf model which can be called directly, or a generator over the model.predict (in tf, so it's quick).
<code>train(training_params, training_data[, ...])</code>	Train on the training data :return:

6.6 DataProcessor

<i>DataProcessor</i>	A utility class for transparently processing (transforming) numpy arrays and un-processing TensorFlow Tensors to aid in emulation.
<i>Log</i>	Return $\log(x + c)$ where c can be specified.
<i>Whiten</i>	Scale the data to have zero mean and unit variance
<i>Normalise</i>	Linearly scale the data to lie between [0, 1]
<i>Flatten</i>	Flatten all dimensions except the leading one
<i>Reshape</i>	Ensure the training data is the right shape for the ConvNet
<i>Recast</i>	Cast the data to a given type

6.6.1 `esem.data_processors.DataProcessor`

class `esem.data_processors.DataProcessor`

A utility class for transparently processing (transforming) numpy arrays and un-processing TensorFlow Tensors to aid in emulation.

See the [API documentation](#) for a list of concrete classes implementing this interface.

`__init__()`

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__()</code>	Initialize self.
<code>process(data)</code>	
<code>unprocess(mean, variance)</code>	

6.6.2 `esem.data_processors.Log`

class `esem.data_processors.Log`(*constant=0.0*)

Return $\log(x + c)$ where c can be specified.

`__init__`(*constant=0.0*)

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__</code> ([constant])	Initialize self.
<code>process(data)</code>	
<code>unprocess(mean, variance)</code>	

6.6.3 `esem.data_processors.Whten`

class `esem.data_processors.Whten`

Scale the data to have zero mean and unit variance

`__init__`()

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__</code> ()	Initialize self.
<code>process(data)</code>	
<code>unprocess(mean, variance)</code>	

6.6.4 `esem.data_processors.Normalise`

class `esem.data_processors.Normalise`

Linearly scale the data to lie between [0, 1]

`__init__()`

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__()</code>	Initialize self.
<code>process(data)</code>	
<code>unprocess(mean, variance)</code>	

6.6.5 `esem.data_processors.Flatten`

class `esem.data_processors.Flatten`

Flatten all dimensions except the leading one

`__init__()`

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__()</code>	Initialize self.
<code>process(data)</code>	
<code>unprocess(mean, variance)</code>	

6.6.6 `esem.data_processors.Reshape`

class `esem.data_processors.Reshape`

Ensure the training data is the right shape for the ConvNet

`__init__()`

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__()</code>	Initialize self.
<code>process(data)</code>	
<code>unprocess(mean, variance)</code>	

6.6.7 `esem.data_processors.Recast`

class `esem.data_processors.Recast(new_type)`

Cast the data to a given type

`__init__(new_type)`

Initialize self. See `help(type(self))` for accurate signature.

Methods

<code>__init__(<i>new_type</i>)</code>	Initialize self.
<code>process(data)</code>	
<code>unprocess(mean, variance)</code>	

6.7 Utilities

A collection of associated utilities which might be of use when performing typical ESEm workflows.

<code>plot_results</code>	Validation plot for LeaveOneOut
<code>validation_plot</code>	
<code>plot_parameter_space</code>	
<code>get_uniform_params</code>	Slightly convoluted method for getting a flat set of points evenly
<code>get_random_params</code>	Get points randomly sampling a (unit) N-dimensional space
<code>ensemble_collocate</code>	Efficiently collocate (interpolate) many ensemble members on to a set of (un-gridded) observations
<code>leave_one_out</code>	Function to perform LeaveOneOut cross-validation with different models.
<code>get_param_mask</code>	Determine the most relevant parameters in the input space using a regularised linear model and either the Aikake or Bayesian Information Criterion.

6.7.1 `esem.utils.plot_results`

`esem.utils.plot_results(ax, truth, pred, title)`
 Validation plot for LeaveOneOut

6.7.2 `esem.utils.validation_plot`

`esem.utils.validation_plot(test_mean, pred_mean, pred_var, figsize=(7, 7), minx=None, miny=None, maxx=None, maxy=None)`

6.7.3 `esem.utils.plot_parameter_space`

`esem.utils.plot_parameter_space(df, nbins=100, target_df=None, smooth=True, xmins=None, xmaxs=None, fig_size=(8, 6))`

6.7.4 `esem.utils.get_uniform_params`

`esem.utils.get_uniform_params(n_params, n_samples=5)`

Slightly convoluted method for getting a flat set of points evenly sampling a (unit) N-dimensional space

Parameters

- **n_params** (int) – The number of parameters (dimensions) to sample from
- **n_samples** (int) – The number of uniformly spaced samples (in each dimension)

Returns ndarray – n_samples*n_params parameters uniformly sampled

6.7.5 `esem.utils.get_random_params`

`esem.utils.get_random_params(n_params, n_samples=5)`

Get points randomly sampling a (unit) N-dimensional space

Parameters

- **n_params** (int) – The number of parameters (dimensions) to sample from
- **n_samples** (int) – The number of parameters to (randomly) sample

6.7.6 `esem.utils.ensemble_collocate`

`esem.utils.ensemble_collocate(ensemble, observations, member_dimension='job')`

Efficiently collocate (interpolate) many ensemble members on to a set of (un-gridded) observations

Note: This function requires both Iris and CIS to be installed

Parameters

- **ensemble** (GriddedData) – The ensemble of (model) samples to interpolate on to the observations
- **observations** (UngriddedData) – The observations on to which the observations will be sampled
- **member_dimension** (str) – The name of the dimension which represents the ensemble members in *ensemble*

Returns **col_ensemble** (iris.cube.Cube) – The ensemble values interpolated on to the observation locations, with the ensemble members along the leading dimension.

6.7.7 esim.utils.leave_one_out

`esim.utils.leave_one_out(Xdata, Ydata, model='RandomForest', **model_kwargs)`

Function to perform LeaveOneOut cross-validation with different models.

Parameters

- **Xdata** (*array-like* of shape (n_samples, n_features)) – Parameter values
- **Ydata** (*array-like* of shape (n_samples,)) – Target values.
- **model** ({'RandomForest', 'GaussianProcess', 'NeuralNet'}, *default* 'RandomForest')
- **model_kwargs** (dict) – More arguments to pass to the model.

Returns **output** (list of n_samples (truth, prediction, variance) tuples) – which can then be passed to `esim.utils.validation_plot()`

6.7.8 esim.utils.get_param_mask

`esim.utils.get_param_mask(X, y, criterion='bic', **kwargs)`

Determine the most relevant parameters in the input space using a regularised linear model and either the Aikake or Baysian Information Criterion.

Parameters

- **X** (*array-like* of shape (n_samples, n_features)) – Parameter values
- **y** (*array-like* of shape (n_samples,)) – target values.
- **criterion** ({'bic', 'aic'}, *default* 'bic') – The information criteria to apply for parameter selection. Either Aikake or Baysian Information Criterion.
- **kwargs** (dict) – Further arguments for `sklearn.feature_selection.SelectFromModel`

Returns **mask** (ndarray) – A boolean array of shape [# input features], in which an element is True iff its corresponding feature is selected for retention.

ESEM DESIGN

Here we provide a brief description of the main architectural decisions behind the design for ESEm in order to hopefully make it easier for contributors and users alike to understand the various components and how they fit together.

7.1 Emulation

We try to provide a seamless interface for users whether they provide iris Cube's, xarray DataArray's or numpy ndarrays. This is done using the `esem.wrappers.DataWrapper` and associated subclasses, which keep a copy of the provided object but only exposes the underlying numpy array to the emulation engines. When the data is requested from this wrapper using the `esem.wrappers.DataWrapper.wrap()` method then it will return a copy of the input object (Cube or DataArray) with the data replaced by the emulated data.

This layer will also ensure the underlying (numpy) data is wrapped in a `esem.wrappers.ProcessWrapper`. This class transparently applies any requested `esem.data_processors.DataProcessor` in sequence.

The user can then create an `esem.emulator.Emulator` object by providing a concrete `esem.model_adaptor.ModelAdaptor` such as a `esem.model_adaptor.KerasModel`. There are two layers of abstraction here: The first to deal with different interfaces to different emulation libraries; and the second to apply the pre- and post-processing and allow a single `esem.emulator.Emulator.batch_stats()` method. The `esem.emulator.Emulator._predict()` provides an important internal interface to the underlying model which reverts any data-processing but leaves the emulator output as a TensorFlow Tensor to allow optimal sampling.

The top-level functions `esem.gp_model()`, `esem.cnn_model()` and `esem.rf_model()` provide a simple interface for constructing these emulators and should be sufficient for most users.

7.2 Calibration

We try and keep this interface very simple; a `esem.sampler.Sampler` should be initialised with an `esem.emulator.Emulator` object to sample from, some observations and associated uncertainties. The only method it has to provide is `esem.sampler.Sampler.sample()` which should provide sample θ from the posterior.

Wherever possible these samplers should take advantage of the fact that the `esem.emulator.Emulator._predict()` method returns TensorFlow tensors and always prefer to use them directly rather than using `esem.emulator.Emulator.predict()` or calling `.numpy()` on them. This allows the sampling to happen on GPUs where available and can substantially speed-up sampling.

The `esem.abc_sampler.ABCSampler` extends this interface to include both `esem.abc_sampler.ABCSampler.get_implausibility()` and `esem.abc_sampler.ABCSampler.batch_constrain()` methods. The first allows inspection of the effect of different observations on the constraint and the second allows a streamlined approach for rejecting samples in batch, taking advantage of the large amounts of memory available on modern GPUs.

GLOSSARY

array_like Any object that can be treated as a numpy array, i.e. can be indexed to retrieve numerical values. Typically not a tensorflow Tensor.

INDICES AND TABLES

- `genindex`
- `modindex`
- `search`

Symbols

`__init__()` (*esem.abc_sampler.ABCSampler* method), 85
`__init__()` (*esem.data_processors.DataProcessor* method), 92
`__init__()` (*esem.data_processors.Flatten* method), 94
`__init__()` (*esem.data_processors.Log* method), 93
`__init__()` (*esem.data_processors.Normalise* method), 94
`__init__()` (*esem.data_processors.Recast* method), 95
`__init__()` (*esem.data_processors.Reshape* method), 94
`__init__()` (*esem.data_processors.Whiten* method), 93
`__init__()` (*esem.emulator.Emulator* method), 80
`__init__()` (*esem.model_adaptor.GPFlowModel* method), 92
`__init__()` (*esem.model_adaptor.KerasModel* method), 91
`__init__()` (*esem.model_adaptor.ModelAdaptor* method), 91
`__init__()` (*esem.model_adaptor.SKLearnModel* method), 91
`__init__()` (*esem.sampler.MCMCSampler* method), 83
`__init__()` (*esem.sampler.Sampler* method), 82
`__init__()` (*esem.wrappers.CubeWrapper* method), 89
`__init__()` (*esem.wrappers.DataArrayWrapper* method), 90
`__init__()` (*esem.wrappers.DataWrapper* method), 88
`__init__()` (*esem.wrappers.ProcessWrapper* method), 88
`_predict()` (*esem.emulator.Emulator* method), 81

A

`ABCSampler` (class in *esem.abc_sampler*), 85
`array_like`, 101

B

`batch_constrain()` (*esem.abc_sampler.ABCSampler* method), 87
`batch_stats()` (*esem.emulator.Emulator* method), 81

C

`cnn_model()` (in module *esem*), 78
`CubeWrapper` (class in *esem.wrappers*), 89

D

`data` (*esem.wrappers.DataWrapper* property), 89
`DataArrayWrapper` (class in *esem.wrappers*), 90
`DataProcessor` (class in *esem.data_processors*), 92
`DataWrapper` (class in *esem.wrappers*), 88
`dtype` (*esem.wrappers.DataWrapper* property), 89

E

`Emulator` (class in *esem.emulator*), 80
`ensemble_collocate()` (in module *esem.utils*), 96

F

`Flatten` (class in *esem.data_processors*), 94

G

`get_implausibility()`
 (*esem.abc_sampler.ABCSampler* method), 86
`get_param_mask()` (in module *esem.utils*), 97
`get_random_params()` (in module *esem.utils*), 96
`get_uniform_params()` (in module *esem.utils*), 96
`gp_model()` (in module *esem*), 77
`GPFlowModel` (class in *esem.model_adaptor*), 92

K

`KerasModel` (class in *esem.model_adaptor*), 91

L

`leave_one_out()` (in module *esem.utils*), 97
`Log` (class in *esem.data_processors*), 93

M

`MCMCSampler` (class in *esem.sampler*), 83
`model` (*esem.emulator.Emulator* attribute), 80
`ModelAdaptor` (class in *esem.model_adaptor*), 91

N

`name` (*esem.emulator.Emulator attribute*), 80
`name()` (*esem.wrappers.DataWrapper method*), 89
`Normalise` (*class in esim.data_processors*), 94

P

`plot_parameter_space()` (*in module esim.utils*), 96
`plot_results()` (*in module esim.utils*), 96
`predict()` (*esem.emulator.Emulator method*), 81
`ProcessWrapper` (*class in esim.wrappers*), 88

R

`Recast` (*class in esim.data_processors*), 95
`Reshape` (*class in esim.data_processors*), 94
`rf_model()` (*in module esim*), 79

S

`sample()` (*esem.abc_sampler.ABCSampler method*), 86
`sample()` (*esem.sampler.MCMCSampler method*), 84
`sample()` (*esem.sampler.Sampler method*), 83
`Sampler` (*class in esim.sampler*), 82
`SKLearnModel` (*class in esim.model_adaptor*), 91

T

`train()` (*esem.emulator.Emulator method*), 81
`training_data` (*esem.emulator.Emulator attribute*), 80

V

`validation_plot()` (*in module esim.utils*), 96

W

`Whiten` (*class in esim.data_processors*), 93
`wrap()` (*esem.wrappers.DataWrapper method*), 89