

---

# CatLearn Documentation

**SUNCAT**

**May 04, 2019**



---

## User Guide:

---

<b>1</b>	<b>Installation</b>	<b>3</b>
<b>2</b>	<b>Changelog</b>	<b>5</b>
<b>3</b>	<b>Version 0.6.1 (April 2019)</b>	<b>7</b>
<b>4</b>	<b>Version 0.6.0 (January 2019)</b>	<b>9</b>
<b>5</b>	<b>Version 0.5.0 (October 2018)</b>	<b>11</b>
<b>6</b>	<b>Version 0.4.4 (August 2018)</b>	<b>13</b>
<b>7</b>	<b>Version 0.4.3 (May 2018)</b>	<b>15</b>
<b>8</b>	<b>Version 0.4.2 (May 2018)</b>	<b>17</b>
<b>9</b>	<b>Version 0.4.1 (April 2018)</b>	<b>19</b>
<b>10</b>	<b>Version 0.4.0 (April 2018)</b>	<b>21</b>
<b>11</b>	<b>Version 0.3.1 (February 2018)</b>	<b>23</b>
<b>12</b>	<b>Version 0.3.0 (February 2018)</b>	<b>25</b>
<b>13</b>	<b>Version 0.2.1 (February 2018)</b>	<b>27</b>
<b>14</b>	<b>Version 0.2.0 (January 2018)</b>	<b>29</b>
<b>15</b>	<b>Version 0.1.0 (December 2017)</b>	<b>31</b>
<b>16</b>	<b>Contributing</b>	<b>33</b>
<b>17</b>	<b>catlearn.api</b>	<b>37</b>
<b>18</b>	<b>catlearn.cross_validation</b>	<b>41</b>
<b>19</b>	<b>catlearn.fingerprint package</b>	<b>45</b>
<b>20</b>	<b>catlearn.ga</b>	<b>57</b>

---

<b>21</b>	<b>catlearn.learning_curve</b>	<b>63</b>
<b>22</b>	<b>catlearn.preprocess</b>	<b>69</b>
<b>23</b>	<b>catlearn.regression</b>	<b>79</b>
<b>24</b>	<b>catlearn.active_learning package</b>	<b>93</b>
<b>25</b>	<b>catlearn.estimator package</b>	<b>99</b>
<b>26</b>	<b>catlearn.utilities</b>	<b>101</b>
<b>27</b>	<b>Indices and tables</b>	<b>109</b>
	<b>Python Module Index</b>	<b>111</b>

CatLearn provides utilities for building and testing atomistic machine learning models for surface science and catalysis.

---

**Note:** This is part of the SUNCAT centers code base for understanding materials for catalytic applications. Other code is hosted on the center's [Github](#) repository.

---

CatLearn provides an environment to facilitate utilization of machine learning within the field of materials science and catalysis. Workflows are typically expected to utilize the Atomic Simulation Environment ([ASE](#)), or [NetworkX](#) graphs. Through close coupling with these codes, CatLearn can generate numerous embeddings for atomic systems. As well as generating a useful feature space for numerous problems, CatLearn has functions for model optimization. Further, Gaussian processes (GP) regression machine learning routines are implemented with additional functionality over standard implementations such as that in scikit-learn. A more detailed explanation of how to utilize the code can be found in the [Tutorials](#) folder.

To featurize ASE atoms objects, the following lines of code can be used:

```
import ase
from ase.cluster.cubic import FaceCenteredCubic

from catlearn.featurize.setup import FeatureGenerator

# First generate an atoms object.
surfaces = [(1, 0, 0), (1, 1, 0), (1, 1, 1)]
layers = [6, 9, 5]
lc = 3.61000
atoms = FaceCenteredCubic('Cu', surfaces, layers, latticeconstant=lc)

# Then generate some features.
generator = FeatureGenerator(nprocs=1)
features = generator.return_vec([atoms], [generator.eigenspectrum_vec,
                                         generator.composition_vec])
```

In the most basic form, it is possible to set up a GP model and make some predictions using the following lines of code:

```
import numpy as np
from catlearn.regression import GaussianProcess

# Define some input data.
train_features = np.arange(200).reshape(50, 4)
target = np.random.random_sample((50,))
test_features = np.arange(100).reshape(25, 4)

# Setup the kernel.
kernel = [{'type': 'gaussian', 'width': 0.5}]

# Train the GP model.
gp = GaussianProcess(kernel_list=kernel, regularization=1e-3,
                     train_fp=train_features, train_target=target,
```

(continues on next page)

(continued from previous page)

```
optimize_hyperparameters=True)

# Get the predictions.
prediction = gp.predict(test_fp=test_features)
```

There is much functionality in CatLearn to assist in handling atom data and building optimal models. This includes:

- API to other codes:
  - Atomic simulation environment API
  - Magpie API
  - NetworkX API
- Fingerprint generators:
  - Bulk systems
  - Support/slab systems
  - Discrete systems
- Preprocessing routines:
  - Data cleaning
  - Feature elimination
  - Feature engineering
  - Feature extraction
  - Feature scaling
- Regression methods:
  - Regularized ridge regression
  - Gaussian processes regression
- Cross-validation:
  - K-fold cv
  - Ensemble k-fold cv
- Optimize:
  - Machine Learning Accelerated Nudged Elastic Band ML-NEB
- General utilities:
  - K-means clustering
  - Neighborlist generators
  - Penalty functions
  - SQLite db storage

# CHAPTER 1

---

## Installation

---

A number of different methods can be used to run the CatLearn code.

### 1.1 Requirements

- ase
- h5py
- networkx
- numpy
- pandas
- scikit-learn
- scipy
- tqdm

### 1.2 Installation using pip

The easiest way to install CatLearn is with:

```
$ pip install catlearn
```

This will automatically install the code as well as the dependencies.

### 1.3 Installation from source

To get the most up-to-date development version of the code, you can clone the git repository to a local directory with:

```
$ git clone https://github.com/SUNCAT-Center/CatLearn.git
```

And then put the <install\_dir>/ into your \$PYTHONPATH environment variable. If you are using Windows, there is some advice on how to do that [here](#).

Be sure to install dependencies in with:

```
$ pip install -r requirements.txt
```

## CHAPTER 2

---

Changelog

---



# CHAPTER 3

---

Version 0.6.1 (April 2019)

---

- Fixed compatibility issue with MLNEB and [GPAW](#)
- Various bugfixes



# CHAPTER 4

---

Version 0.6.0 (January 2019)

---

- Added ML-MIN algorithm for energy minimization.
- Added ML-NEB algorithm for transition state search.
- Changed input format for kernels in the GP.



# CHAPTER 5

---

Version 0.5.0 (October 2018)

---

- Restructure of fingerprint module
- Pandas DataFrame getter in FeatureGenerator
- CatMAP API using ASE database.
- New active learning module.
- Small fixes in adsorbate fingerprinter.



# CHAPTER 6

---

Version 0.4.4 (August 2018)

---

- Major modifications to adsorbates fingerprinter
- Bag of site neighbor coordinations numbers implemented.
- Bag of connections implemented for adsorbate systems.
- General bag of connections implemented.
- Data cleaning function now return a dictionary with ‘index’ of clean features.
- New clean function to discard features with excessive skewness.
- New adsorbate-chalcogenide fingerprint generator.
- Enhancements to automatic identification of adsorbate, site.
- Generalized coordination number for site.
- Formal charges utility.
- New sum electronegativity over bonds fingerprinter.



# CHAPTER 7

---

Version 0.4.3 (May 2018)

---

- `ConvolvedFingerprintGenerator` added for bulk and molecules.
- Dropped support for Python3.4 as it appears to start causing problems.



# CHAPTER 8

---

Version 0.4.2 (May 2018)

---

- Genetic algorithm feature selection can parallelize over population within each generation.
- Default fingerprinter function sets accessible using `catlearn.fingerprint.setup.default_fingerprinters`
- New surrogate model utility
- New utility for evaluating cutoff radii for connectivity based fingerprinting.
- `default_catlearn_radius` improved.



# CHAPTER 9

---

Version 0.4.1 (April 2018)

---

- AtoML renamed to CatLearn and moved to Github.
- Adsorbate fingerprinting again parallelizable.
- Adsorbate fingerprinting use atoms.tags to get layers if present.
- Adsorbate fingerprinting relies on connectivity matrix before neighborlist.
- New bond-electronegativity centered fingerprints for adsorbates.
- Fixed a bug that caused the negative log marginal likelihood to be attached to the gp class.
- Small speed improvement for initialize and updates to GaussianProcess.



# CHAPTER 10

---

Version 0.4.0 (April 2018)

---

- Added `autogen_info` function for list of atoms objects representing adsorbates.
  - This can auto-generate all atomic group information and attach it to `atoms.info`.
  - Parallelized fingerprinting is not yet supported for output from `autogen_info`.
- Added `database_to_list` for import of atoms objects from `ase.db` with formatted metadata.
- Added function to translate a connection matrix to a formatted neighborlist dict.
- `periodic_table_data.list_mendeleev_params` now returns a numpy array.
- Magpie api added, allows for Voronoi and prototype feature generation.
- A genetic algorithm added for feature optimization.
- Parallelism updated to be compatible with Python2.
- Added in better neighborlist generation.
  - Updated wrapper for `ase neighborlist`.
  - Updated CatLearn neighborlist generator.
  - Defaults cutoffs changed to `atomic_radius` plus a relative tolerance.
- Added basic NetworkX api.
- Added some general functions to clean data and build a GP.
- Added a test for dependencies. Will raise a warning in the CI if things get out of date.
- Added a custom docker image for the tests. This is compiled in the `setup/` directory in root.
- Modified uncertainty output. The user can ask for the uncertainty with and without adding noise parameter (regularization).
- Clean up some bits of code, fix some bugs.



# CHAPTER 11

---

Version 0.3.1 (February 2018)

---

- Added a parallel version of the greedy feature selection. **Python3 only!**
- Updated the k-fold cross-validation function to handle features and targets explicitly.
- Added some basic read/write functionality to the k-fold CV.
- A number of minor bugs have been fixed.



# CHAPTER 12

---

Version 0.3.0 (February 2018)

---

- Update the fingerprint generator functions so there is now a `FeatureGenerator` class that wraps round all type specific generators.
- Feature generation can now be performed in parallel, setting `nprocs` variable in the `FeatureGenerator` class. **Python3 only!**
- Add better handling when passing variable length/composition data objects to the feature generators.
- More acquisition functions added.
- Penalty functions added.
- Started adding a general api for ASE.
- Added some more test and changed the way test are called/handled.
- A number of minor bugs have been fixed.



# CHAPTER 13

---

Version 0.2.1 (February 2018)

---

- Update functions to compile features allowing for variable length of atoms objects.
- Added some tutorials for hierarchy cross-validation and prediction on organic molecules.



# CHAPTER 14

---

Version 0.2.0 (January 2018)

---

- Gradients added to hyperparameter optimization.
- More features added to the adsorbate fingerprint generator.
- Acquisition function structure updated. Added new functions.
- Add some standardized input/output functions to save and load models.
- The kernel setup has been made more modular.
- Better test coverage, the tests have also been optimized for speed.
- Better CI configuration. The new method is much faster and more flexible.
- Added Dockerfile and appropriate documentation in the README and CONTRIBUTING guidelines.
- A number of minor bugs have been fixed.



# CHAPTER 15

---

## Version 0.1.0 (December 2017)

---

- The first stable version of the code base!
- For those that used the previous development version, there are many big changes in the way the code is structured. Most scripts will need to be rewritten.
- A number of minor bugs have been fixed.



# CHAPTER 16

---

## Contributing

---

### 16.1 General

There are some general coding conventions that the CatLearn repository adheres to. These include the following:

- Code should support Python 2.7, 3.4 and higher.
- Code should adhere to the [pep8](#) and [pyflakes](#) style guides.
- Tests are run using [TravisCI](#) and coverage tracked using [Coveralls](#).
- When new functions are added, tests should be written and added to the CI script.
- Documentation is hosted on Read the Docs at <http://catlearn.readthedocs.io>.
- Should use NumPy style [docstrings](#).

### 16.2 Git Setup

We adhere to the git workflow described [here](#), if you are considering contributing, please familiarize yourself with this. It is a bad idea to develop directly on the main CatLearn repository. Instead, fork a version into your own namespace on Github with the following:

- Fork the repository and then clone it to your local machine.

```
$ git clone https://github.com/SUNCAT-Center/CatLearn.git
```

- Add and track upstream to the local copy.

```
$ git remote add upstream https://github.com/SUNCAT-Center/CatLearn.git
```

All development can then be performed on the fork and a merge request opened into the upstream when appropriate. It is normally best to open merge requests as soon as possible, as it will allow everyone to see what is being worked on and comment on any potential issues.

## 16.3 Development

The following workflow is recommended when adding some new functionality:

- Before starting any new work, always sync with the upstream version.

```
$ git fetch upstream
$ git checkout master
$ git merge upstream/master --ff-only
```

- It is a good idea to keep the remote repository up to date.

```
$ git push origin master
```

- Start a new branch to do work on.

```
$ git checkout -b branch-name
```

- Once a file has been changed/created, add it to the staging area.

```
$ git add file-name
```

- Now commit it to the local repository and push it to the remote.

```
$ git commit -m 'some descriptive message'
$ git push --set-upstream origin branch-name
```

- When the desired changes have been made on your fork of the repository, open up a merge request on Github.

## 16.4 Environment

It is highly recommended to use pipenv for handling dependencies and the virtual environment, more information can be found [here](#). Once installed, go to the root directory of CatLearn and use:

```
$ pipenv shell
```

From here it is possible to install and upgrade all the dependencies:

```
$ pipenv install --dev
$ pipenv update
```

There are a number of packages that may be important for the development cycle, these are installed with the `--dev` flag. There are then two ways to install additional dependencies required for new functionality, etc:

```
$ pipenv install package
$ pipenv install --dev package
```

The first command will install the package as a dependency for everyone using the code, e.g. people who install CatLearn with pip would be expected to also install this dependency. The second line will only install a package for developers. This workflow can even be used to keep the `requirements.txt` file up-to-date:

```
$ pipenv lock -r > requirements.txt
```

When complete, use `exit` to quit the virtualenv.

## 16.5 Docker

A `docker` image is included in the repository. It is sometimes easier to develop within a controlled environment such as this. In particular, it is possible for other developers to attain the same environment. To run CatLearn in the docker container, use the following commands:

```
$ docker build -t catlearn .
$ docker run -it catlearn bash
```

This will load up the CatLearn directory. To check that everything is working correctly simply run the following:

```
$ python2 test/test_suite.py
$ python3 test/test_suite.py
```

This will run the `test_suite.py` script with python version 2 and 3, respectively. If one version of python is preferred over the other, it is possible to create an alias as normal with:

```
$ alias python=python3
```

### Use **ctrl+d** to exit.

To make changes to this, it is possible to simply edit the `Dockerfile`. To list the images available on the local system, use the following:

```
$ docker images
$ docker inspect REPOSITORY
```

It is a good idea to remove old images. This can be performed using the following lines:

```
$ docker rm $(docker ps -q -f status=exited)
$ docker rmi $(docker images -q -f "dangling=true")
```

## 16.6 Testing

When writing new code, please add some tests to ensure functionality doesn't break over time. We look at test coverage when merge requests are opened and will expect that coverage does not decrease due to large portions of new code not being tested. In CatLearn we just use the built-in unittest framework.

When commits are made, the CI will also automatically test if dependencies are up to date. This test is allowed to fail and will simply return a warning if a module in `requirements.txt` is out of date. This shouldn't be of concern and is mostly in place for us to keep track of changes in other code bases that could cause problems.

If changes are being made that change some core functionality, please run the `tutorials/test_notebooks.py` script. In general, the tutorials involve more demanding computations and thus are not run with the CI. The `test_notebooks.py` script will run through the various tutorials and make sure that they do not fail.

## 16.7 Tutorials

Where appropriate please consider adding some tutorials for new functionality. It would be great if they were written in jupyter notebook form, allowing for some detailed discussion of what is going on in the code.



### 17.1 catlearn.api.ase\_atoms\_api

Functions that interface ase with CatLearn.

`catlearn.api.ase_atoms_api.database_to_list(fname, selection=None)`  
Return a list of atoms objects imported from an ase database.

#### Parameters

- `fname (str)` – path/filename of ase database.
- `selection (list)` – search filters to limit the import.

`catlearn.api.ase_atoms_api.extend_atoms_class(atoms)`  
A wrapper to add extra functionality to ase atoms objects.

#### Parameters `atoms (class)` – An ase atoms object.

`catlearn.api.ase_atoms_api.get_features(self)`  
Function to read feature vector from ase atoms object.

This function provides a uniform way in which to return a feature vector from an atoms object.

#### Parameters `self (class)` – An ase atoms object to attach feature vector to.

#### Returns `fp` – The feature vector attached to the atoms object.

#### Return type array

`catlearn.api.ase_atoms_api.get_graph(self)`  
Function to read networkx graph from ase atoms object.

This function provides a uniform way in which to return a graph object from an atoms object.

#### Parameters `self (class)` – An ase atoms object to attach feature vector to.

#### Returns `graph` – The networkx graph object attached to the atoms object.

#### Return type object

`catlearn.api.ase_atoms_api.get_neighborlist(self)`

Function to read neighborlist from ase atoms object.

This function provides a uniform way in which to return a neighborlist from an atoms object.

**Parameters** `self` (*class*) – An ase atoms object to attach feature vector to.

**Returns** `neighborlist` – The neighbor list attached to the atoms object.

**Return type** dict

`catlearn.api.ase_atoms_api.images_connectivity(images, check_cn_max=False)`

Return a list of atoms objects imported from an ase database.

**Parameters**

- `fname` (*str*) – path/filename of ase database.
- `selection` (*list*) – search filters to limit the import.

`catlearn.api.ase_atoms_api.images_pair_distances(images, mic=True)`

Return a list of atoms objects imported from an ase database.

**Parameters**

- `fname` (*str*) – path/filename of ase database.
- `selection` (*list*) – search filters to limit the import.

`catlearn.api.ase_atoms_api.set_features(self, fp)`

Function to write feature vector to ase atoms object.

This function provides a uniform way in which to attach a feature vector to an atoms object. Can be used in conjunction with the `get_features` function.

**Parameters**

- `self` (*class*) – An ase atoms object to attach feature vector to.
- `fp` (*array*) – The feature vector to attach.

`catlearn.api.ase_atoms_api.set_graph(self, graph)`

Function to write networkx graph to ase atoms object.

This function provides a uniform way in which to attach a graph object to an atoms object. Can be used in conjunction with the `ase_to_networkx` function.

**Parameters**

- `self` (*class*) – An ase atoms object to attach feature vector to.
- `graph` (*object*) – The networkx graph object to attach.

`catlearn.api.ase_atoms_api.set_neighborlist(self, neighborlist)`

Function to write neighborlist to ase atoms object.

This function provides a uniform way in which to attach a neighbor list to an atoms object. Can be used in conjunction with the `get_neighborlist` function.

**Parameters**

- `self` (*class*) – An ase atoms object to attach feature vector to.
- `neighborlist` (*dict*) – The neighbor list dict to attach.

## 17.2 catlearn.api.ase\_data\_setup

Data generation functions to interact with ASE atoms objects.

`catlearn.api.ase_data_setup.get_train(atoms, key, size=None, taken=None)`  
Return a training dataset.

### Parameters

- **atoms** (*list*) – A list of ASE atoms objects.
- **size** (*int*) – Size of training dataset.
- **taken** (*list*) – List of candidates that have been used in unique dataset.
- **key** (*string*) – Property on which to base the predictions stored in the atoms object as atoms.info[‘key\_value\_pairs’][key].

`catlearn.api.ase_data_setup.get_unique(atoms, size, key)`  
Return a unique test dataset.

### Parameters

- **atoms** (*list*) – A list of ASE atoms objects.
- **size** (*int*) – Size of unique dataset to be returned.
- **key** (*string*) – Property on which to base the predictions stored in the atoms object as atoms.info[‘key\_value\_pairs’][key].

## 17.3 catlearn.api.networkx\_graph\_api

API to convert from ASE and NetworkX.

`catlearn.api.networkx_graph_api.ase_to_networkx(atoms, cutoffs=None)`  
Make the NetworkX graph form ASE atoms object.

The graph is dependent on the generation of the neighborlist. Currently this is handled by the version implemented in ASE.

### Parameters

- **atoms** (*object*) – An ASE atoms object.
- **cutoffs** (*list*) – A list of distance parameters for each atom.

**Returns** `atoms_graph` – A networkx graph object.

**Return type** object

`catlearn.api.networkx_graph_api.matrix_to_nl(matrix)`

Returns a neighborlist as a dictionary. :param matrix: symmetric connection matrix. :type matrix: numpy array

**Returns** `nl` – neighborlist.

**Return type** dict

`catlearn.api.networkx_graph_api.networkx_to_adjacency(graph)`

Simple wrapper for graph to adjacency matrix.

**Parameters** `graph` (*object*) – The networkx graph object.

**Returns** `matrix` – The numpy adjacency matrix.

**Return type** array

# CHAPTER 18

---

## catlearn.cross\_validation

---

### 18.1 catlearn.cross\_validation.hierarchy\_cv

Cross validation routines to work with feature database.

```
class catlearn.cross_validation.hierarchy_cv.Hierarchy(file_name, db_name,
                                                       table='FingerVector',
                                                       file_format='pickle')
```

Bases: object

Class to form hierarchy crossvalidation setup.

This class is used to cross-validate with respect to data size. The initial dataset is split in two and subsequent datasets split further until a minimum size is reached. Predictions are made on all subsets of data giving averaged error and certainty at each data size.

```
get_subset_data(index_split, indices, split=None)
```

Make array with training data according to index.

#### Parameters

- **index\_split** (array) – Array with the index data.
- **indices** (array) – Index used to generate data.

```
globalscaledata(index_split)
```

Make an array with all data.

Parameters **index\_split** (array) – Array with the index data.

```
load_split()
```

Function to load the split from file.

```
split_index(min_split, max_split=None, all_index=None)
```

Function to split up the db index to form subsets of data.

#### Parameters

- **min\_split** (int) – Minimum size of a data subset.

- **max\_split** (*int*) – Maximum size of a data subset.
- **all\_index** (*list*) – List of indices in the feature database.

**split\_predict** (*index\_split, predict, \*\*kwargs*)

Function to make predictions looping over all subsets of data.

#### Parameters

- **index\_split** (*dict*) – All data for the split.
- **predict** (*function*) – The prediction function. Must return dict with ‘result’ in it.

#### Returns

- **result** (*list*) – A list of averaged errors for each subset of data.
- **size** (*list*) – A list of data sizes corresponding to the errors list.

**todb** (*features, targets*)

Function to convert numpy arrays to basic db.

**transform\_output** (*data*)

Function to compile results in a format for plotting average error.

**Parameters** **data** (*dict*) – The dictionary output from the split\_predict function.

#### Returns

- **size** (*list*) – A list of the data sizes used in the CV.
- **error** (*list*) – A list of the mean errors at each data size.

## 18.2 catlearn.cross\_validation.k\_fold\_cv

Setup k-fold array split for cross validation.

```
catlearn.cross_validation.k_fold_cv.k_fold(features, targets=None, nsplit=3, fix_size=None)
```

Routine to split feature matrix and return sublists.

#### Parameters

- **features** (*array*) – An n, d feature array.
- **targets** (*list*) – A list to target values.
- **nsplit** (*int*) – The number of bins that data should be devided into.
- **fix\_size** (*int*) – Define a fixed sample size, e.g. nsplit=5 fix\_size=100, generates 5 x 100 data split. Default is None, all available data is divided nsplit times.

#### Returns

- **features** (*list*) – A list of feature arrays of length nsplit.
- **targets** (*list*) – A list of targets lists of length nsplit.

```
catlearn.cross_validation.k_fold_cv.read_split(fname, fformat='pickle')
```

Function to read the k-fold split from file.

#### Parameters

- **fname** (*str*) – The name of the read file.
- **fformat** (*str*) – File format to read from. Can be json or pickle, default is pickle.

**Returns**

- **features** (*list*) – A list of feature arrays of length nsplit.
- **targets** (*list*) – A list of targets lists of length nsplit.

```
catlearn.cross_validation.k_fold_cv.write_split(features, targets, fname, fformat='pickle')
```

Function to write the k-fold split to file.

**Parameters**

- **features** (*array*) – An n, d feature array.
- **targets** (*list*) – A list to target values.
- **fname** (*str*) – The name of the write file.
- **fformat** (*str*) – File format to write to. Can be json or pickle, default is pickle.

Cross validation functions.



# CHAPTER 19

---

## catlearn.fingerprint package

---

### 19.1 Submodules

#### 19.2 catlearn.fingerprint.adsorbate module

Slab adsorbate fingerprint functions for machine learning.

**class** catlearn.fingerprint.adsorbate.**AdsorbateFingerprintGenerator**(*\*\*kwargs*)  
Bases: catlearn.featurize.base.BaseGenerator

**ads\_av**(*atoms=None*)

Function that takes an atoms objects and returns a fingerprint vector with averages of the atomic properties of the adsorbate.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**ads\_sum**(*atoms=None*)

Function that takes an atoms objects and returns a fingerprint vector with averages of the atomic properties of the adsorbate.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**bag\_atoms\_ads**(*atoms=None*)

Function that takes an atoms object and returns a fingerprint vector containing the count of each element in the adsorbate.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**bag\_cn** (*atoms*)

Count the number of neighbors of the site, which has a n number of neighbors. This is equivalent to a bag of coordination numbers over the site neighbors. These can be used in the “alpha parameters” linear model.

Please cite: Roling LT, Abild-Pedersen F. Structure-Sensitive Scaling Relations: Adsorption Energies from Surface Site Stability. *ChemCatChem*. 2018 Apr 9;10(7):1643-50.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**bag\_cn\_general** (*atoms*)

Count the number of neighbors of the site, which has a n number of neighbors. This is equivalent to a bag of coordination numbers over the site neighbors. These can be used in the “alpha parameters” linear model for alloys.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**bag\_edges\_ads** (*atoms*)

Returns bag of connections, counting only the bonds within the adsorbate.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**bag\_edges\_all** (*atoms*)

Returns bag of connections, counting all bonds within the adsorbate and between adsorbate atoms and surface. If we assign an energy to each type of bond, considering first neighbors only, this fingerprint would work independently in a linear model. The length of the vector is atom\_types \* ads\_atom\_types.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**bag\_edges\_chemi** (*atoms*)

Returns bag of connections, counting only the bonds within the adsorbate and the connections between adsorbate and surface.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**bulk** (*atoms=None*)

Return a fingerprint vector with properties averaged over the bulk atoms.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**count\_chemisorbed\_fragment (atoms=None)**

Function that takes an atoms objects and returns a fingerprint vector containing the count over atom types, that are neighbors to the chemisorbing atom.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**ctime (atoms=None)**

Return the contents of atoms.info[‘ctime’] as a feature.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**db\_size (atoms=None)**

Return a fingerprint containing the number of layers in the slab, the number of surface atoms in the unit cell and the adsorbate coverage.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**dbid (atoms=None)**

Return the contents of atoms.info[‘id’] as a feature.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**delta\_energy (atoms=None)**

Return the contents of atoms.info[‘key\_value\_pairs’][‘delta\_energy’] as a feature.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**en\_difference\_active (atoms=None)**

Returns a list of electronegativity metrics, squared and summed over adsorbate bonds including those with the surface.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**en\_difference\_ads (atoms=None)**

Returns a list of electronegativity metrics, squared and summed over bonds within the adsorbate atoms.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**en\_difference\_chemi (atoms=None)**

Returns a list of electronegativity metrics, squared and summed over adsorbate-site bonds.

**Parameters** `atoms` (*object*) – ASE Atoms object.

**Returns** `features` – If None was passed, the elements are strings, naming the feature.

**Return type** list

`generalized_cn` (*atoms*)

Returns the averaged generalized coordination number over the site. Calle-Vallejo et al. Angew. Chem. Int. Ed. 2014, 53, 8316-8319.

**Parameters** `atoms` (*object*) – ASE Atoms object.

**Returns** `features` – If None was passed, the elements are strings, naming the feature.

**Return type** list

`max_site` (*atoms=None*)

Function that takes an atoms objects and returns a fingerprint vector with properties averaged over the surface metal atoms closest to an add atom.

**Parameters** `atoms` (*object*) – ASE Atoms object.

**Returns** `features` – If None was passed, the elements are strings, naming the feature.

**Return type** list

`mean_chemisorbed_atoms` (*atoms=None*)

Function that takes an atoms objects and returns a fingerprint vector containing properties of the closest add atom to a surface metal atom.

**Parameters** `atoms` (*object*) – ASE Atoms object.

**Returns** `features` – If None was passed, the elements are strings, naming the feature.

**Return type** list

`mean_site` (*atoms=None*)

Function that takes an atoms objects and returns a fingerprint vector with properties averaged over the surface metal atoms closest to an add atom.

**Parameters** `atoms` (*object*) – ASE Atoms object.

**Returns** `features` – If None was passed, the elements are strings, naming the feature.

**Return type** list

`mean_surf_ligands` (*atoms=None*)

Function that takes an atoms objects and returns a fingerprint vector containing the count of nearest neighbors and properties of the nearest neighbors.

**Parameters** `atoms` (*object*) – ASE Atoms object.

**Returns** `features` – If None was passed, the elements are strings, naming the feature.

**Return type** list

`median_site` (*atoms=None*)

Function that takes an atoms objects and returns a fingerprint vector with properties averaged over the surface metal atoms closest to an add atom.

**Parameters** `atoms` (*object*) – ASE Atoms object.

**Returns** `features` – If None was passed, the elements are strings, naming the feature.

**Return type** list

**min\_site**(atoms=None)

Function that takes an atoms objects and returns a fingerprint vector with properties averaged over the surface metal atoms closest to an add atom.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**strain**(atoms=None)

Return a fingerprint with the expected strain of the site atoms and the termination atoms.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**sum\_site**(atoms=None)

Function that takes an atoms objects and returns a fingerprint vector with properties summed over the surface metal atoms closest to an add atom.

**Parameters** **atoms** (*object*) – ASE Atoms object.

**Returns** **features** – If None was passed, the elements are strings, naming the feature.

**Return type** list

**term**(atoms=None)

Return a fingerprint vector with properties averaged over the termination atoms.

**Parameters** **atoms** (*object*) –

## 19.3 catlearn.fingerprint.bulk module

Slab adsorbate fingerprint functions for machine learning.

**class** catlearn.fingerprint.bulk.BulkFingerprintGenerator(\*\*kwargs)

Bases: catlearn.featurize.base.BaseGenerator

**bulk\_average**(atoms=None)

Return a fingerprint vector with properties of the element name saved in the atoms.info['key\_value\_pairs'][‘bulk’]

**bulk\_std**(atoms=None)

Return a fingerprint vector with properties of the element name saved in the atoms.info['key\_value\_pairs'][‘bulk’]

**bulk\_summation**(atoms=None)

Return a fingerprint vector with properties of the element name saved in the atoms.info['key\_value\_pairs'][‘bulk’]

**xyz\_id**(atoms=None)

## 19.4 catlearn.fingerprint.chalcogenide module

Slab adsorbate fingerprint functions for machine learning.

```
class catlearn.fingerprint.chalcogenide.ChalcogenideFingerprintGenerator(**kwargs)
Bases: catlearn.featurize.base.BaseGenerator

formal_charges(atoms)
    Return a fingerprint based on formal charges.

    Parameters atoms(object) –

max_cation(atoms=None)
    Function that takes an atoms objects and returns a fingerprint vector with properties averaged over the surface metal atoms closest to an add atom.

    Parameters atoms(object) –

mean_cation(atoms=None)
    Function that takes an atoms objects and returns a fingerprint vector with properties averaged over the surface metal atoms closest to an add atom.

    Parameters atoms(object) –

median_cation(atoms=None)
    Function that takes an atoms objects and returns a fingerprint vector with properties averaged over the surface metal atoms closest to an add atom.

    Parameters atoms(object) –

min_cation(atoms=None)
    Function that takes an atoms objects and returns a fingerprint vector with properties averaged over the surface metal atoms closest to an add atom.

    Parameters atoms(object) –

sum_cation(atoms=None)
    Function that takes an atoms objects and returns a fingerprint vector with properties summed over the surface metal atoms closest to an add atom.

    Parameters atoms(object) –
```

## 19.5 catlearn.fingerprint.convoluted module

Slab adsorbate convoluted fingerprint functions for machine learning.

```
class catlearn.fingerprint.convoluted.ConvolvedFingerprintGenerator(**kwargs)
Bases: catlearn.featurize.base.BaseGenerator

conv_bulk(atoms=None)
    Return a fingerprint vector with propeties convoluted over the bulk atoms.

    Parameters atoms(object) – A single atoms object.

conv_term(atoms=None)
    Return a fingerprint vector with propeties convoluted over the terminal atoms.

    Parameters atoms(object) – A single atoms object.

catlearn.fingerprint.convoluted.check_length(labels, result, atoms)
Check that two lists have the same length.

If not, print an informative error message containing a databse id if present.

Parameters
    • labels(list) – A list of feature names.
```

- **result** (*list*) – A fingerprint.
- **atoms** (*object*) – A single atoms object.

## 19.6 catlearn.fingerprint.graph module

Functions to build a neighbor matrix feature representation.

**class** catlearn.fingerprint.graph.GraphFingerprintGenerator (\*\*kwargs)

Bases: catlearn.featurize.base.BaseGenerator

Function to build a fingerprint vector based on an atoms object.

**neighbor\_mean\_vec** (*data*)

Transform neighborlist into a neighbor averaged feature vector.

**Parameters** **data** (*object*) – Target data object from which to generate features.

**Returns** **features** – A 1d numpy array of the feature vector.

**Return type** array

**neighbor\_sum\_vec** (*data*)

Transform neighborlist into a neighbor sum feature vector.

**Parameters** **data** (*object*) – Target data object from which to generate features.

**Returns** **features** – A 1d numpy array of the feature vector.

**Return type** array

## 19.7 catlearn.fingerprint.molecule module

Functions to build a gas phase molecule fingerprint.

**class** catlearn.fingerprint.molecule.AutoCorrelationFingerprintGenerator (\*\*kwargs)

Bases: catlearn.featurize.base.BaseGenerator

Class for constructing an autocorrelation fingerprint.

**get\_autocorrelation** (*atoms*)

Return the autocorrelation fingerprint for a molecule.

## 19.8 catlearn.fingerprint.particle module

Nanoparticle fingerprint functions.

These functions will typically perform well at describing chemical ordering within alloyed nanoparticles. However, they may be applicable to other applications where bond counting or coordination numbers are important descriptors.

This class inherits from the catlearn.fingerprint.BaseGenerator function.

**class** catlearn.fingerprint.particle.ParticleFingerprintGenerator (\*\*kwargs)

Bases: catlearn.featurize.base.BaseGenerator

Function to build a fingerprint vector based on an atoms object.

**bond\_count\_vec** (*data*)

Bond counting with a distribution measure for coordination.

**Parameters** **data** (*object*) – Data object with atomic distances.

**Returns** **track\_nmat** – List with summed number of atoms with given coordination numbers.

**Return type** list

**connections\_vec** (*data*)

Sum atoms with a certain number of connections.

**distribution\_vec** (*data*)

Return atomic distribution measure.

**nearestneighbour\_vec** (*data*)

Nearest neighbour average, Topics in Catalysis, 2014, 57, 33.

This is a slightly modified version of the code found in the *ase.ga* module.

**Parameters** **data** (*object*) – Data object with atomic numbers available.

**Returns** **nnlist** – Feature vector that will be  $n^{**}2$  where n is the number of atomic species passed to the class.

**Return type** list

**rdf\_vec** (*data*)

Return list of partial rdbs for use as fingerprint vector.

## 19.9 catlearn.fingerprint.prototype module

Prototype fingerprint based on Magpie.

```
class catlearn.fingerprint.prototype.PrototypeFingerprintGenerator(atoms,
    sites, system_name='',
    target='id',
    delete_temp=True,
    properties=[])
```

Bases: object

Function to build prototype fingerprint in pandas.DataFrame.

Based on a list of *ase.atoms* object.

**generate()**

Generate Prototype fingerprint and return all the fingerprint.

**Returns FP**

**Return type** pandas.Frame

**generate\_all()**

**run\_proto()**

Call Magpie to generate Prototype FP and write to proto\_FP.csv.

**update\_str()**

---

```
write_proto_input()
    Write Prototype input for Magpie.

class catlearn.fingerprint.prototype.PrototypeSites (site_dict=None)
    Bases: object

    Prototype site objective for generating prototype input.
```

## 19.10 catlearn.fingerprint.standard module

Standard fingerprint functions.

These feature sets should perform relatively well on a variety of different systems. They are general descriptors based predominantly on the elemental properties and in some cases structure.

This class inherits from the catlearn.fingerprint.BaseGenerator function.

```
class catlearn.fingerprint.standard.StandardFingerprintGenerator (**kwargs)
    Bases: catlearn.featurize.base.BaseGenerator

    Function to build a fingerprint vector based on an atoms object.

bag_edges (atoms)
    Returns the bag of connections, defined as counting connections between types of elements pairs. We define the bag as a vector, e.g. return [Number of C-H connections, # C-C, # C-O, ..., # M-X]

        Parameters atoms (object) –
        Returns features
        Return type list

bag_edges_cn (atoms)
    Returns the bag of connections folded with coordination numbers of the node atoms.

        Parameters atoms (object) –
        Returns features
        Return type list

bag_element_cn (atoms)
    Bag elements folded with coordination numbers, e.g. number of C with CN = 4, number of C with CN = 3, ect.

        Parameters atoms (object) – ASE Atoms object.
        Returns features – If None was passed, the elements are strings, naming the feature.
        Return type list

bag_elements (atoms)
    Returns the bag of elements, defined as counting occurrence of elements in a given structure. This is mostly useful for subtracting atomization energies.

        Parameters atoms (object) –
        Returns features
        Return type list

composition_vec (data)
    Function to return a feature vector based on the composition.
```

**Parameters** `data` (*object*) – Data object with atomic numbers available.

**Returns** `features` – Vector containing a count of the different atomic types, e.g. for CH<sub>3</sub>OH the vector [1, 4, 1] would be returned.

**Return type** array

`distance_vec` (*data*)

Averaged distance between e.g. A-A atomic pairs.

**Parameters** `data` (*object*) – Data object with Cartesian coordinates and atomic numbers available.

**Returns** `features` – Vector of averaged distances between homoatomic atoms.

**Return type** ndarray

`eigenspectrum_vec` (*data*)

Sorted eigenspectrum of the Coulomb matrix.

**Parameters** `data` (*object*) – Data object with Cartesian coordinates and atomic numbers available.

**Returns** `features` – Sorted Eigen values of the coulomb matrix, n atoms is size.

**Return type** ndarray

`element_mass_vec` (*data*)

Function to return a vector based on mass parameter.

**Parameters** `data` (*object*) – Data object with atomic masses available.

**Returns** `features` – Vector of the summed mass.

**Return type** ndarray

`element_parameter_vec` (*data*)

Function to return a vector based on a defined parameter.

The vector is compiled based on the summed parameters for each elemental type as well as the sum for all atoms.

**Parameters** `data` (*object*) – Data object with atomic numbers available.

**Returns** `features` – An n + 1 array where n in the length of self.atom\_types.

**Return type** array

## 19.11 catlearn.fingerprint.voro module

Voronoi fingerprint based on Magpie.

`class` `catlearn.fingerprint.voro.VoronoiFingerprintGenerator` (*atoms*,  
*delete\_temp=True*)

Bases: object

Function to build voronoi fingerprint in pandas.DataFrame.

Based on a list of ase.atoms object.

`generate()`

Generate Voronoi fingerprint and return all the fingerprint.

**Returns** FP

**Return type** pandas.DataFrame

**run\_voro()**

Call Magpie to generate Voronoi FP and write to voro\_FP.csv.

**write\_voro\_input()**

Write Voronoi input for Magpie.

## 19.12 Module contents



# CHAPTER 20

---

catlearn.ga

---

## 20.1 catlearn.ga.algorithm

The GeneticAlgorithm class methods.

```
class catlearn.ga.algorithm.GeneticAlgorithm(fit_func, features, targets, population_size='auto', population=None, operators=None, fitness_parameters=1, nsplit=2, accuracy=None, nprocs=1, dmax=None)
```

Bases: object

Genetic algorithm for parameter optimization.

```
search(steps, natural_selection=True, convergence_operator=None, repeat=5, verbose=False, writefile=None)
```

Do the actual search.

### Parameters

- **steps** (*int*) – Maximum number of steps to be taken.
- **natural\_selection** (*bool*) – A flag that when set True will perform natural selection.
- **convergence\_operator** (*object*) – The function to perform the convergence check. If None is passed then the *no\_progress* function is used.
- **repeat** (*int*) – Number of repeat generations with no progress.
- **verbose** (*bool*) – If True, will print out the progress of the search. Default is False.
- **writefile** (*str*) – Name of a json file to save data too.

### population

The current population.

Type list

**fitness**

The fitness for the current population.

**Type** list

## 20.2 catlearn.ga.convergence

Functions to check for convergence in the GA.

**class** catlearn.ga.convergence.**Convergence**

Bases: object

Class to check convergence.

**no\_progress** (*fitness, repeat*)

Convergence based on a lack of any progress in the search.

**Parameters**

- **fitness** (*array*) – A List of fitnesses from the search.
- **repeat** (*int*) – Number of repeat generations with no progress.

**Returns** **converged** – True if convergence has been reached, False otherwise.

**Return type** bool

**stagnation** (*fitness, repeat*)

Convergence based on a stagnation of the population.

**Parameters**

- **fitness** (*array*) – A List of fitnesses from the search.
- **repeat** (*int*) – Number of repeat generations with no progress.

**Returns** **converged** – True if convergence has been reached, False otherwise.

**Return type** bool

## 20.3 catlearn.ga.initialize

Function to initialize a population.

catlearn.ga.initialize.**initialize\_population** (*pop\_size, dimension, dmax=None*)

Generate a random starting population.

**Parameters**

- **pop\_size** (*int*) – Population size.
- **d\_param** (*int*) – Dimension of parameters in model.

## 20.4 catlearn.ga.io

Functions to read and write GA data.

catlearn.ga.io.**read\_data** (*writefile*)

Funtion to read population and fitness.

**Parameters** `writefile` (*str*) – Name of the JSON file to read.

**Returns**

- **population** (*array*) – The population saved from a previous search.
- **fitness** (*array*) – The fitness associated with the saved population.

## 20.5 catlearn.ga.mating

Cut and splice mating function.

`catlearn.ga.mating.cut_and_splice` (*parent\_one*, *parent\_two*, *index='random'*)

Perform cut\_and\_splice between two parents.

**Parameters**

- **parent\_one** (*list*) – List of params for first parent.
- **parent\_two** (*list*) – List of params for second parent.
- **index** (*str*) – Define how to choose size of each cut index.

**Returns** `offspring` – A new child candidate from the two parents.

**Return type** array

## 20.6 catlearn.ga.mutate

Define some mutation functions.

`catlearn.ga.mutate.probability_include` (*parent\_one*)

A mutation that will include features with a certain probability.

**Parameters** `parent_one` (*list*) – List of params for first parent.

**Returns** `p1` – Mutated parameter list based on the parent parameters provided.

**Return type** list

`catlearn.ga.mutate.probability_remove` (*parent\_one*)

A mutation that will remove features with a certain probability.

**Parameters** `parent_one` (*list*) – List of params for first parent.

**Returns** `p1` – Mutated parameter list based on the parent parameters provided.

**Return type** list

`catlearn.ga.mutate.random_permutation` (*parent\_one*)

Perform a random permutation on a parameter index.

**Parameters** `parent_one` (*list*) – List of params for first parent.

**Returns** `p1` – Mutated parameter list based on the parent parameters provided.

**Return type** list

## 20.7 catlearn.ga.natural\_selection

Functions to perform some natural selection.

`catlearn.ga.natural_selection.population_reduction(pop, fit, population_size)`

Method to reduce population size to constant.

### Parameters

- **pop** (*list*) – Extended population.
- **fit** (*list*) – Extended fitness assignment.
- **population\_size** (*int*) – The population size.
- **pareto** (*bool*) – Flag to specify whether search is for Pareto optimal set.

### Returns

- **population** (*list*) – The population after natural selection.
- **fitness** (*list*) – The fitness for the current population.

`catlearn.ga.natural_selection.remove_duplicates(population, fitness, accuracy)`

Function to delete duplicate candidates based on fitness.

### Parameters

- **population** (*array*) – The current population.
- **fitness** (*array*) – The fitness for the current population.
- **accuracy** (*int*) – Number of decimal places to include when finding unique.

### Returns

- **population** (*list*) – The population after duplicates deleted.
- **fitness** (*list*) – The fitness for the population after duplicates deleted.

## 20.8 catlearn.ga.predictors

Some generic prediction functions.

`catlearn.ga.predictors.minimize_error(train_features, train_targets, test_features, test_targets)`

A generic fitness function.

This fitness function will minimize the cost function.

### Parameters

- **train\_features** (*array*) – The training features.
- **train\_targets** (*array*) – The training targets.
- **test\_features** (*array*) – The test features.
- **test\_targets** (*array*) – The test targets.

`catlearn.ga.predictors.minimize_error_descriptors(train_features, train_targets, test_features, test_targets)`

A generic fitness function.

This fitness function will minimize the cost function as well as the number of descriptors. This will provide a Pareto optimial set of solutions upon convergence.

#### Parameters

- **train\_features** (*array*) – The training features.
- **train\_targets** (*array*) – The training targets.
- **test\_features** (*array*) – The test feaatures.
- **test\_targets** (*array*) – The test targets.

```
catlearn.ga.predictors.minimize_error_time(train_features, train_targets, test_features,  
                                         test_targets)
```

A generic fitness function.

This fitness function will minimize the cost function as well as the time to train the model. This will provide a Pareto optimial set of solutions upon convergence.

#### Parameters

- **train\_features** (*array*) – The training features.
- **train\_targets** (*array*) – The training targets.
- **test\_features** (*array*) – The test features.
- **test\_targets** (*array*) – The test targets.



# CHAPTER 21

---

## catlearn.learning\_curve

---

### 21.1 catlearn.learning\_curve.data\_process

Processing of data for HierarchyValidation.

```
class catlearn.learning_curve.data_process(features,      min_split,
                                           max_split,      scale=True,
                                           normalization=True,
                                           ridge=True,    loocv=True,
                                           batchfarm=False)
```

Bases: object

Class to glue different function used for HierarchyValidation.

This class pick up data from HierarchyValidation. The data is then modified if requested with “feature\_preprocess”, and “predict”. The data is then fitted with regression model for example with “ridge\_regression”. The error of the fit is then measured.

**average\_nested**(*Y, X*)

Calculate statistics for prediciton.

#### Parameters

- **data\_size** (*list*) – Data\_size for where the prediction were made.
- **p\_error** (*list*) – Error for where the prediction were made.

**get\_statistic**(*data\_size, p\_error*)

Generate statistics for prediciton.

#### Parameters

- **data\_size** (*list*) – Data\_size for where the prediction were made.
- **p\_error** (*list*) – Error for where the prediction were made.

**globalscaling**(*globalscaledata, train\_features*)

All sub-groups of traindata are scaled same.

**Parameters** `globalscaledata` (*string*) – The data will be scaled globally if requested.

**`prediction_error`** (*test\_features*, *test\_targets*, *coef*, *s\_tar*, *m\_tar*)

Calculate the error of the prediction with the model.

#### Parameters

- **`test_features`** (*array*) – Independet data for testing the model.
- **`test_targets`** (*array*) – Dependent data to test the model.
- **`coef`** (*array*) – The coeffieiceints which makes up the model.
- **`s_tar`** (*string*) – Standard devation or (max-min), for the dependent train\_targets.
- **`m_tar`** (*array*) – Mean for the dependent train\_targets.

**`scaling_data`** (*train\_features*, *train\_targets*, *test\_features*, *s\_tar*, *m\_tar*, *s\_feat*, *m\_feat*)

Scaling the data if requested.

#### Parameters

- **`train_feature`** (*array*) – Independent data used to train model.
- **`train_targets`** (*array*) – Dependent data used to train model.
- **`test_features`** (*array*) – Independent data used to test the model.
- **`s_tar`** (*array*) – Standard devation or (max-min), for the dependent train\_targets.
- **`m_tar`** (*array*) – Mean for the dependent train\_targets.
- **`s_feat`** (*array*) – Standard devation or (max-min), for the independent train\_features.
- **`m_feat`** (*array*) – Mean for the independent train\_features.

## 21.2 catlearn.learning\_curve.feature\_selection

Feature selection with lasso.

**`class catlearn.learning_curve.feature_selection`** (*train\_features*, *train\_targets*)

Bases: object

Class made to make it possible to select features.

Used with hierarchy cross-validation.

**`alpha_finder`** (*feat\_vec*, *alpha\_vec*, *feat*)

Find the alpha corresponding to the number of features.

#### Parameters

- **`feat_vec`** (*list*) – Features within the interval.
- **`alpha_vec`** (*list*) – Alphas within the interval.
- **`feat`** (*int*) – The group of feature searched.

**`alpha_refinement`** (*alpha*, *feat*, *splits=10*, *refsteps=1*, *upper=1.5*)

Find a more stringent alpha for the number of feature searched for.

#### Parameters

- **`alpha`** (*int*) – Initial alpha found for the nuumber of feature searched for. Will be used as a lower limit.

- **feat** (*int*) – The number of feature searched for.
- **splits** (*int*) – Increase of Number of alphas under inspection within interval.
- **refsteps** (*int*) – Number of refinements.
- **upper** – How many times alpha the upper limit should be.

**feature\_inspection** (*lower=0, upper=1, interval=100, alpha\_list=None*)

Generate interval used to search for the alpha.

#### Parameters

- **lower** (*int*) – Lower bound for the interval search.
- **upper** (*int*) – Upper bound for the interval search.
- **interval** (*int*) – Number of alphas in interval inspected.

**interval\_modifier** (*feat\_vec, alpha\_vec, feat, splits, int\_expand*)

Modify the interval under inspection by reduction or expansion.

#### Parameters

- **feat\_vec** (*list*) – Features within the interval.
- **alpha\_vec** (*list*) – Alphas within the interval.
- **feat** (*int*) – The group of feature searched.
- **splits** (*int*) – Increase of Number of alphas under inspection within interval.
- **int\_expand** (*int*) – Number of times the number of alphas in interval is increased.

**selection** (*select\_limit*)

Select the the fature/s that works best wtig L1.

## 21.3 catlearn.learning\_curve.learning\_curve

Generate the learning curve.

**class** catlearn.learning\_curve.learning\_curve.**LearningCurve** (*nprocs=1*)  
Bases: object

Learning curve class. Test a model while varying the density of the training data.

**run** (*model, train, target, test, test\_target, step=1, min\_data=2*)  
Evaluate a model versus training data size.

#### Parameters

- **model** (*object*) – A function that will train or load a regression model or classifier and make predictions for testing. model should accept the parameters:
  - train\_features : array test\_features : array train\_targets : list test\_targets : list
 model should return either a float or a list of floats. The float or the first value of the list will be used as the fitness score.
- **train** (*array*) – An n, d array of training examples.
- **targets** (*test*) – A list of the target values.
- **test** (*array*) – An n, d array of test data.
- **targets** – A list of the test target values.

- **step** (*int*) – Incrementent the data set size by this many examples.
- **min\_data** (*int*) – Smallest number of training examples to test.

**Returns** **output** – Each row is the output from the model object.

**Return type** array

```
catlearn.learning_curve.learning_curve.feature_frequency(cv, features,  
min_split, max_split,  
smallest=False,  
new_data=True,  
ridge=True, scale=True,  
globalscale=True, nor-  
malization=True, featse-  
lect_featvar=False, feat-  
select_featconst=True,  
select_limit=None,  
feat_sub=15)
```

Function to extract raw data from the database.

#### Parameters

- **features** (*int*) – Number of features used for regression.
- **min\_split** (*int*) – Number of datasplit in the smallest sub-set.
- **max\_split** (*int*) – Number of datasplit in the largest sub-set.
- **new\_data** (*string*) – Use new data or the previous data.
- **ridge** (*string*) – Ridge regulazer is deafult. If False, lasso is used.
- **scale** (*string*) – If the data are supposed to be scaled or not.
- **globalscale** (*string*) – Using global scaleing or not.
- **normalization** (*string*) – If scaled, normalized or standardized. Normalized is default.
- **feature\_selection** (*string*) – Using feature selection with ridge, or plain vanilla ridge.
- **select\_limit** (*int*) – Up to have many number of features used for feature selection.

```
catlearn.learning_curve.learning_curve.hierarchy(cv, features, min_split, max_split,  
new_data=True, ridge=True,  
scale=True, globalscale=True,  
normalization=True, feat-  
select_featvar=False, feat-  
select_featconst=True, se-  
lect_limit=None, feat_sub=15)
```

Start the hierarchy.

#### Parameters

- **features** (*int*) – Number of features used for regression.
- **min\_split** (*int*) – Number of datasplit in the smallest sub-set.
- **max\_split** (*int*) – Number of datasplit in the largest sub-set.
- **new\_data** (*string*) – Use new data or the previous data.
- **ridge** (*string*) – Ridge regulazer is deafult. If False, lasso is used.

- **scale** (*string*) – If the data are supposed to be scaled or not.
- **globalscale** (*string*) – Using global scaleing or not.
- **normalization** (*string*) – If scaled, normalized or standardized. Normalized is default.
- **feature\_selection** (*string*) – Using feature selection with ridge, or plain vanilla ridge.
- **select\_limit** (*int*) – Up to have many number of features used for feature selection.

## 21.4 catlearn.learning\_curve.placeholder

Placeholder for now.

```
class catlearn.learning_curve.placeholder.placeholder (PC, index_split, hv, indices,
hier_level, featselect_featvar,
featselect_featconst, s_feat,
m_feat, feat_sub=15,
s_tar=None, m_tar=None,
select_limit=None, selected_features=None,
glob_feat1=None,
glob_tar1=None,
new_training=True)
```

Bases: object

Used to make the hierarchey more easy to follow.

Placeholder for now.

**get\_data\_scale** (*split, set\_size=None, p\_error=None, result=None*)

Get the data for each sub-set of data and scales it accordingly.

### Parameters

- **split** (*int*) – Which sub-set od data within hierarchy level.
- **result** (*list*) – Contain all the coefficien and omega2 for all training data.
- **set\_size** (*list*) – Size of sub-set of data/features which the model is based on.
- **p\_error** (*list*) – The prediction error for plain vanilla ridge.

**getstats()**

Used to get features for the frequencyplots.

**predict\_subsets** (*result=None, set\_size=None, p\_error=None*)

Run the prediction on each sub-set of data on the hierarchy level.

### Parameters

- **result** (*list*) – Contain all the coefficien and omega2 for all training data.
- **set\_size** (*list*) – Size of sub-set of data/features which the model is based on.
- **p\_error** (*list*) – The prediction error for plain vanilla ridge.

**reg\_data\_var** (*train\_features, train\_targets, test\_features, test\_targets, ridge, set\_size, p\_error, result*)

Ridge regression and calculation of prediction error.

### Parameters

- **train\_features** (*array*) – Independent data used to train the model.
- **train\_targets** (*array*) – Dependent data used to train model.
- **test\_features** (*array*) – Independent data used to test model.
- **test\_target** (*array*) – Dependent data used to test model.
- **ridge** (*object*) – Generates the model based on the training data.
- **set\_size** (*list*) – Size of sub-set of data/features which the model is based on.
- **p\_error** (*list*) – The prediction error for plain vanilla ridge.
- **result** (*list*) – Contain all the coefficien and omega2 for all training data.

**reg\_feat\_var** (*train\_features*, *train\_targets*, *test\_features*, *test\_targets*, *ridge*, *set\_size*, *p\_error*, *result*)

Regression within a dataset with varying feature.

### Parameters

- **train\_features** (*array*) – Independent data used to train the model.
- **train\_targets** (*array*) – Dependent data used to train model.
- **test\_features** (*array*) – Independent data used to test model.
- **test\_target** (*array*) – Dependent data used to test model.
- **ridge** (*object*) – Generates the model based on the training data.
- **p\_error** (*list*) – The prediction error for feature selection corresponding to different feature set.
- **set\_size** (*list*) – Different data/feature set used for feature selection.
- **result** (*list*) – Contain all the coefficien and omega2 for all training data.

## 21.5 catlearn.learning\_curve.pltfile

# CHAPTER 22

---

## catlearn.preprocess

---

### 22.1 catlearn.preprocess.clean\_data

Functions to clean data.

```
catlearn.preprocess.clean_data.clean_infinite(train,      test=None,      targets=None,
                                              labels=None,      mask=None,
                                              max_impute_fraction=0,      strat-
                                              egy='mean')
```

Remove features that have non finite values in the training data.

Optionally removes features in test data with non fininte values. Returns a dictionary with the clean ‘train’, ‘test’ and ‘index’ that were removed from the original data.

#### Parameters

- **train** (*array*) – Feature matrix for the traing data.
- **test** (*array*) – Optional feature matrix for the test data. Default is None passed.
- **targets** (*array*) – An array of training targets.
- **labels** (*array*) – Optional list of feature labels. Default is None passed.
- **mask** (*list*) – Indices of features that are not subject to cleaning.
- **max\_impute\_fraction** (*float*) – Maximum fraction of values in a column that can be imputed. Columns with higher fractions of nans values will be discarded.
- **strategy** (*str*) – Imputation strategy.

#### Returns

**data** –

key value pairs

- ‘train’ [array] Clean training data matrix.
- ‘test’ [array] Clean test data matrix

- **'targets'** [list] Boolean list on whether targets are finite.
- **'labels'** [list] Feature labels of clean data set.

**Return type** dict

```
catlearn.preprocess.clean_data.clean_skewness(train, test=None, labels=None,  
mask=None, skewness=3.0)
```

Discards features that are excessively skewed.

**Parameters**

- **train** (array) – Feature matrix for the training data.
- **test** (array) – Optional feature matrix for the test data. Default is None passed.
- **labels** (array) – Optional list of feature labels. Default is None passed.
- **mask** (list) – Indices of features that are not subject to cleaning.
- **skewness** (float) – Maximum allowed skewness threshold.

```
catlearn.preprocess.clean_data.clean_variance(train, test=None, labels=None,  
mask=None)
```

Remove features that contribute nothing to the model.

Removes a feature if there is zero variance in the training data. If this is the case, then the model won't learn anything new from adding this feature as it will just act as a scalar.

**Parameters**

- **train** (array) – Feature matrix for the training data.
- **test** (array) – Optional feature matrix for the test data. Default is None passed.
- **labels** (array) – Optional list of feature labels. Default is None passed.
- **mask** (list) – Indices of features that are not subject to cleaning.

```
catlearn.preprocess.clean_data.remove_outliers(features, targets, con=1.4826, dev=3.0,  
constraint=None)
```

Preprocessing routine to remove outliers by median absolute deviation.

This will take the training feature and target arrays, calculate any outliers, then return the reduced arrays. It is possible to set a constraint key ('high', 'low', None) in order to allow for outliers that are e.g. very low in energy, as this may be the desired outcome of the study.

**Parameters**

- **features** (array) – Feature matrix for training data.
- **targets** (list) – List of target values for the training data.
- **con** (float) – Constant scale factor dependent on the distribution. Default is 1.4826 expecting the data is normally distributed.
- **dev** (float) – The number of deviations from the median to account for.
- **constraint** (str) – Can be set to 'low' to remove candidates with targets that are too small/negative or 'high' for outliers that are too large/positive. Default is to remove all.

## 22.2 catlearn.preprocess.feature\_elimination

Functions to select features for the fingerprint vectors.

---

```
class catlearn.preprocess.feature_elimination.FeatureScreening(correlation='pearson',
                                                               iterative=True,
                                                               regression='ridge',
                                                               random_state=None,
                                                               domain_check=False)
```

Bases: object

Class for feature elimination based on correlation screening.

**eliminate\_features** (target, train\_features, test\_features, size=None, step=None, order=None)  
Function to eliminate features from training/test data.

#### Parameters

- **target** (*list*) – The target values for the training data.
- **train\_features** (*array*) – Array of training data to eliminate features from.
- **test\_features** (*array*) – Array of test data to eliminate features from.
- **size** (*int*) – Number of features after elimination.
- **step** (*int*) – Number of features to eliminate at each step.
- **order** (*list*) – Precomputed ordered indices for features.

#### Returns

- **reduced\_train** (*array*) – Reduced training feature matrix, now n x size shape.
- **reduced\_test** (*array*) – Reduced test feature matrix, now m x size shape.

**iterative\_screen** (target, feature\_matrix, size=None, step=None)  
Function iteratively screen feautes.

#### Parameters

- **target** (*list*) – The target values for the training data.
- **feature\_matrix** (*array*) – The feature matrix for the training data.
- **size** (*int*) – Number of features to be returned. Default is number of data.
- **step** (*int*) – Step size by which to reduce the number of features. Default is n / log(n).

#### Returns

- **index** (*list*) – The ordered list of feature indices, top index[:size] will be indices for best features.
- **size** (*int*) – Number of accepted features.

**screen** (target, feature\_matrix)

Feature selection based on SIS.

Further discussion on this topic can be found in Fan, J., Lv, J., J. R. Stat. Soc.: Series B, 2008, 70, 849.

#### Parameters

- **target** (*list*) – The target values for the training data.
- **feature\_matrix** (*array*) – The feature matrix for the training data.

#### Returns

- **index** (*list*) – The ordered list of feature indices.
- **correlation** (*list*) – The ordered list of correlations between features and targets.

- **size** (*int*) – Number of accepted features following screening.

## 22.3 catlearn.preprocess.feature\_engineering

Functions for feature engineering.

```
catlearn.preprocess.feature_engineering.generate_features(p, max_num=2,  
max_den=1, log=False,  
sqrt=False, exclude=False, s=False)
```

Generate composite features from a combination of input features.

developer note: This is currently scales *quite slowly* with max\_den. There's surely a better way to do this, but it's apparently currently functional.

### Parameters

- **p** (*list*) – User-provided list of physical features to be combined.
- **max\_num** (*integer*) – The maximum order of the polynomial in the numerator of the composite features. Must be non-negative.
- **max\_den** (*integer*) – The maximum order of the polynomial in the denominator of the composite features. Must be non-negative.
- **log** (*boolean (not currently supported)*) – Set to True to include terms involving the logarithm of the input features. Default is False.
- **sqrt** (*boolean (not currently supported)*) – Set to True to include terms involving the square root of the input features. Default is False.
- **exclude** (*bool*) – Set exclude=True to avoid returning 1 to represent the zeroth power. Default is False.
- **s** (*bool*) – Set True to return a list of strings and False to evaluate each element in the list. Default is False.

**Returns** **features** – A list of combinations of the input features to meet the required specifications.

### Return type

```
catlearn.preprocess.feature_engineering.generate_positive_features(p, N, exclude=False,  
s=False)
```

Generate list of polynomial combinations in list p up to order N.

Example: p = (a,b,c) ; N = 3

returns (order not preserved) [a\*a\*a, a\*a\*b, a\*a\*c, a\*b\*b, a\*b\*c, a\*c\*c, b\*b\*b, b\*b\*c, b\*c\*c, c\*c\*c, a\*a, a\*b, a\*c, b\*b, b\*c, c\*c, a, b, c]

### Parameters

- **p** (*list*) – Features to be combined.
- **N** (*integer*) – The maximum polynomial coefficient for combinations. Must be non-negative.
- **exclude** (*bool*) – Set True to avoid returning 1 to represent the zeroth power. Default is False.
- **s** (*bool*) – Set True to return a list of strings and False to evaluate each element in the list. Default is False.

**Returns all\_powers** – A list of combinations of the input features to meet the required specifications.

**Return type** list

`catlearn.preprocess.feature_engineering.get_ablog(A, a, b)`

Get all combinations  $x_{ij} = a * \log(x_i) + b * \log(x_j)$ .

The sorting order in dimension 0 is preserved.

**Parameters**

- **A** (*array*) – An  $n \times m$  matrix, where  $n$  is the number of training examples and  $m$  is the number of features.
- **a** (*float*) –
- **b** (*float*) –

**Returns new\_features** – The  $n \times \text{triangular}(m)$  matrix of new features.

**Return type** array

`catlearn.preprocess.feature_engineering.get_div_order_2(A)`

Get all combinations  $x_{ij} = x_i / x_j$ , where  $x_{i,j}$  are features.

The sorting order in dimension 0 is preserved. If a denominator is 0, Inf is returned.

**Parameters** **A** (*array*) –  $n \times m$  matrix, where  $n$  is the number of training examples and  $m$  is the number of features.

**Returns new\_features** – The  $n \times m^{**2}$  matrix of new features.

**Return type** array

`catlearn.preprocess.feature_engineering.get_labels_ablog(l, a, b)`

Get all combinations  $ij$ , where  $i,j$  are feature labels.

**Parameters**

- **a** (*float*) –
- **b** (*float*) –

**Returns new\_features** – List of new feature names.

**Return type** list

`catlearn.preprocess.feature_engineering.get_labels_order_2(l, div=False)`

Get all combinations  $ij$ , where  $i,j$  are feature labels.

**Parameters** **x** (*list*) – Length  $m$  vector, where  $m$  is the number of features.

**Returns new\_features** – List of new feature names.

**Return type** list

`catlearn.preprocess.feature_engineering.get_labels_order_2ab(l, a, b)`

Get all combinations  $ij$ , where  $i,j$  are feature labels.

**Parameters** **x** (*list*) – Length  $m$  vector, where  $m$  is the number of features.

**Returns new\_features** – List of new feature names.

**Return type** list

`catlearn.preprocess.feature_engineering.get_order_2(A)`

Get all combinations  $x_{ij} = x_i * x_j$ , where  $x_{i,j}$  are features.

The sorting order in dimension 0 is preserved.

**Parameters** `A` (*array*) – n x m matrix, where n is the number of training examples and m is the number of features.

**Returns** `new_features` – The n x triangular(m) matrix of new features.

**Return type** array

`catlearn.preprocess.feature_engineering.get_order_2ab(A, a, b)`

Get all combinations  $x_{ij} = x_i^{**a} * x_j^{**b}$ , where  $x_{i,j}$  are features.

The sorting order in dimension 0 is preserved.

**Parameters**

- `A` (*array*) – n x m matrix, where n is the number of training examples and m is the number of features.
- `a` (*float*) –
- `b` (*float*) –

**Returns** `new_features` – The n x triangular(m) matrix of new features.

**Return type** array

`catlearn.preprocess.feature_engineering.single_transform(A)`

Perform single variable transform  $x^2$ ,  $x^{0.5}$  and  $\log(x)$ .

**Parameters** `A` (*array*) – n x m matrix, where n is the number of training examples and m is the number of features.

**Returns** `new_features` – The n x m\*3 matrix of new features.

**Return type** array

## 22.4 catlearn.preprocess.feature\_extraction

Some feature extraction routines.

`catlearn.preprocess.feature_extraction.catlearn_pca(components, train_features, test_features=None, cleanup=False, scale=False)`

Principal component analysis varient that doesn't require scikit-learn.

**Parameters**

- `components` (*int*) – Number of principal components to transform the feature set by.
- `test_fpv` (*array*) – The feature matrix for the testing data.

`catlearn.preprocess.feature_extraction.pca(components, train_matrix, test_matrix)`

Principal component analysis routine.

**Parameters**

- `components` (*int*) – The number of components to be returned.
- `train_matrix` (*array*) – The training features.
- `test_matrix` (*array*) – The test features.

**Returns**

- **new\_train** (*array*) – Extracted training features.
- **new\_test** (*array*) – Extracted test features.

`catlearn.preprocess.feature_extraction.pls (components, train_matrix, target, test_matrix)`  
Projection of latent structure routine.

**Parameters**

- **components** (*int*) – The number of components to be returned.
- **train\_matrix** (*array*) – The training features.
- **test\_matrix** (*array*) – The test features.

**Returns**

- **new\_train** (*array*) – Extracted training features.
- **new\_test** (*array*) – Extracted test features.

`catlearn.preprocess.feature_extraction.sPCA (components, train_matrix, test_matrix)`  
Sparse principal component analysis routine.

**Parameters**

- **components** (*int*) – The number of components to be returned.
- **train\_matrix** (*array*) – The training features.
- **test\_matrix** (*array*) – The test features.

**Returns**

- **new\_train** (*array*) – Extracted training features.
- **new\_test** (*array*) – Extracted test features.

## 22.5 catlearn.preprocess.greedy\_elimination

Greedy feature selection routines.

`class catlearn.preprocess.greedy_elimination.GreedyElimination (nprocs=1,  
verbose=True,  
save_file=None)`

Bases: object

The greedy feature elimination class.

`greedy_elimination (predict, features, targets, nsplit=2, step=1)`  
Greedy feature elimination.

Function to iterate through feature set, eliminating worst feature in each pass. This is the backwards greedy algorithm.

**Parameters**

- **predict** (*object*) – A function that will make the predictions. predict should accept the parameters:

train\_features : array test\_features : array train\_targets : list test\_targets : list

`predict` should return either a float or a list of floats. The float or the first value of the list will be used as the fitness score.

- **features** (*array*) – An n, d array of features.
- **targets** (*list*) – A list of the target values.
- **nsplit** (*int*) – Number of folds in k-fold cross-validation.

#### Returns

**output** – First column is the index of features in the order they were eliminated.

Second column are corresponding cost function values, averaged over the k fold split.

Following columns are any additional values returned by `predict`, averaged over the k fold split.

**Return type** array

## 22.6 catlearn.preprocess.importance\_testing

Functions to check feature significance.

```
class catlearn.preprocess.importance_testing.ImportanceElimination(transform,
                                                                     nprocs=1,
                                                                     verbose=True)
```

Bases: object

The feature importance elimination class.

```
importance_elimination(train_predict, test_predict, features, targets, nsplit=2, step=1)
```

Importance feature elimination.

Function to iterate through feature set, eliminating least important feature in each pass. This is the backwards elimination algorithm.

#### Parameters

- **train\_predict** (*object*) – A function that will train a model. The function should accept the parameters:  
  
    train\_features : array train\_targets : list  
  
    predict should return a function that can be passed to `test_predict`.
- **test\_predict** (*object*) – A function that will accept a trained model object and return a float or a list of test metrics. The first returned metric will be used to eliminate features.
- **features** (*array*) – An n, d array of features.
- **targets** (*list*) – A list of the target values.
- **nsplit** (*int*) – Number of folds in k-fold cross-validation.
- **step** (*int*) – Optional number of features to eliminate in each round.

#### Returns

**output** – First column is the index of features in the order they were eliminated.

Second column are corresponding cost function values, averaged over the k fold split.

Following columns are any additional values returned by test\_predict, averaged over the k fold split.

**Return type** array

`catlearn.preprocess.importance_testing.feature_invariance(args)`

Make a feature invariant.

**Parameters** `args` (*list*) – A list of arguments:

`index` [int] The index of the feature to be shuffled.

`train_features` [array] The original training data matrix.

`test_features` [array] The original test data matrix.

**Returns**

- `train (array)` – Feature matrix with a shuffled feature column in matrix.

- `test (array)` – Feature matrix with a shuffled feature column in matrix.

`catlearn.preprocess.importance_testing.feature_randomize(args)`

Make a feature random noise.

**Parameters** `args` (*list*) – A list of arguments:

`index` [int] The index of the feature to be shuffled.

`train_features` [array] The original training data matrix.

`test_features` [array] The original test data matrix.

**Returns**

- `train (array)` – Feature matrix with a shuffled feature column in matrix.

- `test (array)` – Feature matrix with a shuffled feature column in matrix.

`catlearn.preprocess.importance_testing.feature_shuffle(args)`

Shuffle a feature.

The method has a number of advantages for measuring feature importance. Notably the original values and scale of the feature are maintained.

**Parameters** `args` (*list*) – A list of arguments:

`index` [int] The index of the feature to be shuffled.

`train_features` [array] The original training data matrix.

`test_features` [array] The original test data matrix.

**Returns**

- `train (array)` – Feature matrix with a shuffled feature column in matrix.

- `test (array)` – Feature matrix with a shuffled feature column in matrix.

## 22.7 catlearn.preprocess.scaling

Functions to process the raw feature matrix.

`catlearn.preprocess.scaling.min_max(train_matrix, test_matrix=None, local=True)`

Normalize each feature relative to the min and max.

### Parameters

- **train\_matrix** (*list*) – Feature matrix for the training dataset.
- **test\_matrix** (*list*) – Feature matrix for the test dataset.
- **local** (*boolean*) – Define whether to scale locally or globally.

```
catlearn.preprocess.scaling.normalize(train_matrix, test_matrix=None, mean=None, dif=None, local=True)
```

Normalize each feature relative to mean and min/max variance.

### Parameters

- **train\_matrix** (*list*) – Feature matrix for the training dataset.
- **test\_matrix** (*list*) – Feature matrix for the test dataset.
- **local** (*boolean*) – Define whether to scale locally or globally.
- **mean** (*list*) – List of mean values for each feature.
- **dif** (*list*) – List of max-min values for each feature.

```
catlearn.preprocess.scaling.standardize(train_matrix, test_matrix=None, mean=None, std=None, local=True)
```

Standardize each feature relative to the mean and standard deviation.

### Parameters

- **train\_matrix** (*array*) – Feature matrix for the training dataset.
- **test\_matrix** (*array*) – Feature matrix for the test dataset.
- **mean** (*list*) – List of mean values for each feature.
- **std** (*list*) – List of standard deviation values for each feature.
- **local** (*boolean*) – Define whether to scale locally or globally.

```
catlearn.preprocess.scaling.target_center(target)
```

Return a list of normalized target values.

**Parameters** **target** (*list*) – A list of the target values.

```
catlearn.preprocess.scaling.target_normalize(target)
```

Return a list of normalized target values.

**Parameters** **target** (*list*) – A list of the target values.

```
catlearn.preprocess.scaling.target_standardize(target)
```

Return a list of standardized target values.

**Parameters** **target** (*list*) – A list of the target values.

```
catlearn.preprocess.scaling.unit_length(train_matrix, test_matrix=None, local=True)
```

Normalize each feature vector relative to the Euclidean length.

### Parameters

- **train\_matrix** (*list*) – Feature matrix for the training dataset.
- **test\_matrix** (*list*) – Feature matrix for the test dataset.
- **local** (*boolean*) – Define whether to scale locally or globally.

# CHAPTER 23

---

catlearn.regression

---

## 23.1 catlearn.regression.gpfunctions

### 23.1.1 catlearn.regression.gpfunctions.covariance

Generation of covariance matrix.

```
catlearn.regression.gpfunctions.covariance.get_covariance(kernel_list, log_scale,
                                                               matrix1,          ma-
                                                               trix2=None,        reg-
                                                               ularization=None,
                                                               eval_gradients=False)
```

Return the covariance matrix of training dataset.

#### Parameters

- **kernel\_list** (*dict of dicts*) – A dict containing all dictionaries for the kernels.
- **log\_scale** – Flag to define if the hyperparameters are log scale.
- **train\_matrix** (*list*) – A list of the training fingerprint vectors.
- **test\_matrix** (*list*) – A list of the test fingerprint vectors.
- **regularization** (*None or float*) – Smoothing parameter for the Gramm matrix.

### 23.1.2 catlearn.regression.gpfunctions.default\_scale

Scale everything within regression functions.

```
class catlearn.regression.gpfunctions.default_scale.ScaleData(train_features,
                                                               train_targets)
```

Bases: object

Class to perform default scaling in the regression functions.

Will standardize both the features and the targets. These can then be rescaled before being returned. The parameters can be accessed from the class with:

```
ScaleData.feature_data['mean']
```

This can be accessed from the gp with:

```
gp = GaussianProcess(...) gp.scaling.feature_data['mean']
```

**rescale\_targets** (*predictions*)

Rescale predictions.

**Parameters** ***predictions*** (*list*) – The predicted values from the GP.

**Returns** **p** – The rescaled predictions.

**Return type** array

**test** (*test\_features*)

Scale the test features.

**Parameters** ***test\_features*** (*array*) – Feature matrix for the test data.

**Returns** **scaled\_features** – The scaled features for the test data.

**Return type** array

**train()**

Scale the training features and targets.

**Returns**

- **feature\_data** (*array*) – The scaled features for the training data.

- **target\_data** (*array*) – The scaled targets for the training data.

### 23.1.3 catlearn.regression.gpfunctions.hyperparameter\_scaling

Utility to scale hyperparameters.

```
catlearn.regression.gpfunctions.hyperparameter_scaling.hyperparameters(scaling,  
ker-  
nel_list)
```

Scale the hyperparameters.

```
catlearn.regression.gpfunctions.hyperparameter_scaling.rescale_hyperparameters(scaling,  
ker-  
nel_list)
```

Rescale hyperparameters.

### 23.1.4 catlearn.regression.gpfunctions.io

Functions to read and write models to file.

```
catlearn.regression.gpfunctions.io.read(filename, ext='pkl')
```

Function to read a pickle of model object.

**Parameters**

- **filename** (*str*) – The name of the save file.

- **ext** (*str*) – Format to save GP, can be pkl or hdf5. Default is pkl.

**Returns** **model** – Python GaussianProcess object.

**Return type** obj

```
catlearn.regression.gpfunctions.io.read_train_data(filename)
Function to read raw training data.
```

**Parameters** `filename` (*str*) – The name of the save file.

**Returns**

- `train_features` (*arr*) – Arry of the training features.
- `train_targets` (*list*) – A list of the training targets.
- `regularization` (*float*) – The regularization parameter.
- `kernel_list` (*list*) – The dictionary containing parameters for the kernels.

```
catlearn.regression.gpfunctions.io.write(filename, model, ext='pkl')
Function to write a pickle of model object.
```

**Parameters**

- `filename` (*str*) – The name of the save file.
- `model` (*obj*) – Python GaussianProcess object.
- `ext` (*str*) – Format to save GP, can be pkl or hdf5. Default is pkl.

```
catlearn.regression.gpfunctions.io.write_train_data(filename, train_features,
                                                    train_targets, regularization,
                                                    kernel_list)
```

Function to write raw training data.

**Parameters**

- `filename` (*str*) – The name of the save file.
- `train_features` (*arr*) – Arry of the training features.
- `train_targets` (*list*) – A list of the training targets.
- `regularization` (*float*) – The regularization parameter.
- `kernel_list` (*dict*) – The list containing dictionaries for the kernels.

### 23.1.5 catlearn.regression.gpfunctions.kernel\_scaling

Function to scale kernel hyperparameters.

```
catlearn.regression.gpfunctions.kernel_scaling.kernel_scaling(scale_data, kernel_list, rescale)
Base hyperparameter scaling function.
```

**Parameters**

- `scale_data` (*object*) – Output from the default scaling function.
- `kernel_list` (*list*) – Dictionary containing all dictionaries for the kernels.
- `rescale` (*boolean*) – Flag for whether to scale or rescale the data.

### 23.1.6 catlearn.regression.gpfunctions.kernel\_setup

Functions to prepare and return kernel data.

`catlearn.regression.gpfunctions.kernel_setup.kdict2list (kdict, N_D=None)`

Return ordered list of hyperparameters.

Assumes function is given a dictionary containing properties of a single kernel. The dictionary must contain either the key ‘hyperparameters’ or ‘theta’ containing a list of hyperparameters or the keys ‘type’ containing the type name in a string and ‘width’ in the case of a ‘gaussian’ or ‘laplacian’ type or the keys ‘degree’ and ‘slope’ in the case of a ‘quadratic’ type.

#### Parameters

- **kdict** (*dict*) – A kernel dictionary containing the keys ‘type’ and optional keys containing the hyperparameters of the kernel.
- **N\_D** (*none or int*) – The number of descriptors if not specified in the kernel dict, by the length of the lists of hyperparameters.

`catlearn.regression.gpfunctions.kernel_setup.kdicts2list (kernel_list, N_D=None)`

Return ordered list of hyperparameters given the kernel dictionary.

The kernel dictionary must contain one or more dictionaries, each specifying the type and hyperparameters.

#### Parameters

- **kernel\_list** (*dict*) – A dictionary containing kernel dictionaries.
- **N\_D** (*int*) – The number of descriptors if not specified in the kernel dict, by the length of the lists of hyperparameters.

`catlearn.regression.gpfunctions.kernel_setup.list2kdict (hyperparameters, kernel_list)`

Return updated kernel dictionary with updated hyperparameters from list.

Assumed an ordered list of hyperparameters and the previous kernel dictionary. The kernel dictionary must contain a dictionary for each kernel type in the same order as their respective hyperparameters in the list hyperparameters.

#### Parameters

- **hyperparameters** (*list*) – All hyperparameters listed in the order they are specified in the kernel dictionary.
- **kernel\_list** (*dict*) – A dictionary containing kernel dictionaries.

`catlearn.regression.gpfunctions.kernel_setup.prepare_kernels (kernel_list, regularization_bounds, eval_gradients, N_D)`

Format kernel listionary and stores bounds for optimization.

#### Parameters

- **kernel\_list** (*dict*) – List containing all dictionaries for the kernels.
- **regularization\_bounds** (*tuple*) – Optional to change the bounds for the regularization.
- **eval\_gradients** (*boolean*) – Flag to change kernel setup based on gradients being defined.
- **N\_D** (*int*) – Number of dimensions of the original data.

### 23.1.7 catlearn.regression.gpfunctions.kernels

Contains kernel functions and gradients of kernels.

```
catlearn.regression.gpfunctions.kernels.AA_kernel(theta, log_scale, m1, m2=None,
                                                eval_gradients=False)
```

Generate the covariance between data with a Aichinson & Aitken kernel.

#### Parameters

- **theta** (*list*) – [l, n, c]
- **log\_scale** (*boolean*) – Scaling hyperparameters in the kernel can be useful for optimization.
- **m1** (*list*) – A list of the training fingerprint vectors.
- **m2** (*list*) – A list of the training fingerprint vectors.

**Returns** **k** – The covariance matrix.

**Return type** array

```
catlearn.regression.gpfunctions.kernels.constant_kernel(theta, log_scale,
                                                       m1, m2=None,
                                                       eval_gradients=False)
```

Return constant to add to the kernel.

#### Parameters

- **theta** (*list*) – A list of widths for each feature.
- **log\_scale** (*boolean*) – Scaling hyperparameters in the kernel can be useful for optimization.
- **eval\_gradients** (*boolean*) – Analytical gradients of the training features can be included.
- **m1** (*list*) – A list of the training fingerprint vectors.
- **m2** (*list*) – A list of the training fingerprint vectors.

**Returns** **k** – The covariance matrix.

**Return type** array

```
catlearn.regression.gpfunctions.constant_multi_kernel(theta, log_scale,
                                                       m1, m2=None,
                                                       eval_gradients=True)
```

Return constant to add to the kernel.

#### Parameters

- **theta** (*list*) – A list containing the constants.
- **log\_scale** (*boolean*) – Scaling hyperparameters in the kernel can be useful for optimization.
- **eval\_gradients** (*boolean*) – Analytical gradients of the training features can be included.
- **m1** (*list*) – A list of the training fingerprint vectors.
- **m2** (*list*) – A list of the training fingerprint vectors.

**Returns** **k** – The covariance matrix.

**Return type** array

```
catlearn.regression.gpfunctions.kernels.gaussian_dk_dwidth(k, m1, kwidth,  
log_scale=False)
```

Return gradient of the gaussian kernel with respect to the j'th width.

**Parameters**

- **k** (*array*) – n by n array. The (not scaled) gaussian kernel.
- **m1** (*list*) – A list of the training fingerprint vectors.
- **kwidth** (*float*) – The full list of widths
- **log\_scale** (*boolean*) – Scaling hyperparameters in kernel can be useful for optimization.

```
catlearn.regression.gpfunctions.kernels.gaussian_kernel(theta, log_scale,  
m1, m2=None,  
eval_gradients=False)
```

Generate the covariance between data with a Gaussian kernel.

**Parameters**

- **theta** (*list*) – A list of widths for each feature.
- **log\_scale** (*boolean*) – Scaling hyperparameters in the kernel can be useful for optimization.
- **eval\_gradients** (*boolean*) – Analytical gradients of the training features can be included.
- **m1** (*list*) – A list of the training fingerprint vectors.
- **m2** (*list*) – A list of the training fingerprint vectors.

**Returns** **k** – The covariance matrix.

**Return type** array

```
catlearn.regression.gpfunctions.kernels.gaussian_xx_gradients(m1, kwidth, k)  
Gradient for k(x, x).
```

**Parameters**

- **m1** (*array*) – Feature matrix.
- **kwidth** (*list*) – List of lengthscales for the gaussian kernel.
- **k** (*array*) – Upper left portion of the overall covariance matrix.

```
catlearn.regression.gpfunctions.kernels.gaussian_xxp_gradients(m1, m2, kwidth,  
k)  
Gradient for k(x, x').
```

**Parameters**

- **m1** (*array*) – Feature matrix.
- **m2** (*array*) – Feature matrix typically associated with the test data.
- **kwidth** (*list*) – List of lengthscales for the gaussian kernel.
- **k** (*array*) – Upper left portion of the overall covariance matrix.

```
catlearn.regression.gpfunctions.kernels.laplacian_dk_dwidth(k, m1, kwidth,  
log_scale=False)
```

```
catlearn.regression.gpfunctions.kernels.laplacian_kernel(theta,           log_scale,
                                                       m1,                 m2=None,
                                                       eval_gradients=False)
```

Generate the covariance between data with a laplacian kernel.

#### Parameters

- **theta** (*list*) – A list of widths for each feature.
- **log\_scale** (*boolean*) – Scaling hyperparameters in the kernel can be useful for optimization.
- **m1** (*list*) – A list of the training fingerprint vectors.
- **m2** (*list or None*) – A list of the training fingerprint vectors.

**Returns** **k** – The covariance matrix.

**Return type** array

```
catlearn.regression.gpfunctions.kernels.linear_kernel(theta,           log_scale,
                                                       m1,                 m2=None,
                                                       eval_gradients=False)
```

Generate the covariance between data with a linear kernel.

#### Parameters

- **theta** (*list*) – A list containing constant offset.
- **log\_scale** (*boolean*) – Scaling hyperparameters in the kernel can be useful for optimization.
- **eval\_gradients** (*boolean*) – Analytical gradients of the training features can be included.
- **m1** (*list*) – A list of the training fingerprint vectors.
- **m2** (*list or None*) – A list of the training fingerprint vectors.

**Returns** **k** – The covariance matrix.

**Return type** array

```
catlearn.regression.gpfunctions.kernels.noise_multi_kernel(theta,           log_scale,
                                                       m1,                 m2=None,
                                                       eval_gradients=False)
```

Return constant to add to the kernel.

#### Parameters

- **theta** (*list*) – A list containing the constants to be added in the diagonal of the covariance matrix .
- **eval\_gradients** (*boolean*) – Analytical gradients of the training features can be included.
- **m1** (*list*) – A list of the training fingerprint vectors.
- **m2** (*list*) – A list of the training fingerprint vectors.

**Returns** **k** – The covariance matrix.

**Return type** array

```
catlearn.regression.gpfunctions.kernels.quadratic_dk_ddegree(k,    m1,    degree,
                                                               log_scale=False)
```

```
catlearn.regression.gpfunctions.kernels.quadratic_dk_dslope(k, m1, slope,  
log_scale=False)  
catlearn.regression.gpfunctions.kernels.quadratic_kernel(theta, log_scale,  
m1, m2=None,  
eval_gradients=False)
```

Generate the covariance between data with a quadratic kernel.

### Parameters

- **theta** (*list*) – A list containing slope and degree for quadratic.
- **log\_scale** (*boolean*) – Scaling hyperparameters in the kernel can be useful for optimization.
- **m1** (*list*) – A list of the training fingerprint vectors.
- **m2** (*list or None*) – A list of the training fingerprint vectors.

**Returns** **k** – The covariance matrix.

**Return type** array

```
catlearn.regression.gpfunctions.kernels.scaled_sqe_kernel(theta, log_scale,  
m1, m2=None,  
eval_gradients=False)
```

Generate the covariance between data with a Gaussian kernel.

### Parameters

- **theta** (*list*) – A list of hyperparameters.
- **log\_scale** (*boolean*) – Scaling hyperparameters in the kernel can be useful for optimization.
- **m1** (*list*) – A list of the training fingerprint vectors.
- **m2** (*list*) – A list of the training fingerprint vectors.

**Returns** **k** – The covariance matrix.

**Return type** array

```
catlearn.regression.gpfunctions.kernels.sqe_kernel(theta, log_scale, m1, m2=None,  
eval_gradients=False)
```

Generate the covariance between data with a Gaussian kernel.

### Parameters

- **theta** (*list*) – A list of widths for each feature.
- **log\_scale** (*boolean*) – Scaling hyperparameters in the kernel can be useful for optimization.
- **m1** (*list*) – A list of the training fingerprint vectors.
- **m2** (*list*) – A list of the training fingerprint vectors.

**Returns** **k** – The covariance matrix.

**Return type** array

## 23.1.8 `catlearn.regression.gpfunctions.log_marginal_likelihood`

Log marginal likelihood calculator function.

```
catlearn.regression.gpfunctions.log_marginal_likelihood.dk_dtheta_j(theta,
train_matrix,
ker-
nel_list,
Q)
```

Return the jacobian of the log marginal likelihood.

This is calculated with respect to the hyperparameters, as in: Equation 5.9 in C. E. Rasmussen and C. K. I. Williams, 2006

#### Parameters

- **theta** (*list*) – A list containing the hyperparameters.
- **train\_matrix** (*list*) – A list of the test fingerprint vectors.
- **kernel\_list** (*list*) – A list of kernel dictionaries.
- **Q** (*array*.*–*

```
catlearn.regression.gpfunctions.log_marginal_likelihood.log_marginal_likelihood(theta,
train_matrix,
tar-
gets,
ker-
nel_list,
scale_optimizer,
eval_gradients,
cinv=None,
eval_jac=False)
```

Return the negative of the log marginal likelihood.

Equation 5.8 in C. E. Rasmussen and C. K. I. Williams, 2006

#### Parameters

- **theta** (*list*) – A list containing the hyperparameters.
- **train\_matrix** (*list*) – A list of the test fingerprint vectors.
- **targets** (*list*) – A list of target values.
- **kernel\_list** (*dict*) – A list of kernel dictionaries.
- **scale\_optimizer** (*boolean*) – Flag to define if the hyperparameters are log scale for optimization.
- **eval\_gradients** (*boolean*) – Flag to specify whether to compute gradients in covariance.
- **cinv** (*array*) – Pre-computed inverted covariance matrix.
- **eval\_jac** (*boolean*) – Flag to specify whether to calculate gradients for hyperparameter optimization.

### 23.1.9 catlearn.regression.gpfunctions.sensitivity

Function performing GP sensitivity analysis.

```
class catlearn.regression.gpfunctions.sensitivity.SensitivityAnalysis(train_matrix,
    train_targets,
    test_matrix,
    ker-
    nel_list,
    init_reg=0.001,
    init_width=10.0)
```

Bases: object

Perform sensitivity analysis to estimate important features.

```
backward_selection(predict=False, test_targets=None, selection=None)
```

Feature selection with backward elimination.

#### Parameters

- **predict** (boolean) – Specify whether to make predictions on test data.
- **test\_targets** (list) – A list of test targets to calculate errors, if known.
- **selection** (int, list) – Specify the number or range of features to consider.

### 23.1.10 catlearn.regression.gpfunctions.uncertainty

Function performing uncertainty analysis.

```
catlearn.regression.gpfunctions.uncertainty.get_uncertainty(kernel_list, test_fp,
    ktb, cinv, log_scale)
```

Function to calculate uncertainty.

#### Parameters

- **kernel\_list** (list) – List containing all dictionaries for the kernels.
- **test\_fp** (array) – Test feature set.
- **ktb** (array) – Covariance matrix for test and training data.
- **cinv** (array) – Covariance matrix for training dataset.
- **log\_scale** (boolean) – Flag to define if the hyperparameters are log scale.

**Returns** **uncertainty** – The uncertainty on each prediction in the test data. By default, this includes a measure of the noise on the data.

**Return type** list

## 23.2 catlearn.regression.cost\_function

Functions to calculate the cost statistics.

```
catlearn.regression.cost_function.get_error(prediction, target, metrics=None, epsilon=None, return_percentiles=True)
```

Return error for predicted data.

Discussed in: Rosasco et al, Neural Computation, (2004), 16, 1063-1076.

#### Parameters

- **prediction** (list) – A list of predicted values.
- **target** (list) – A list of target values.

- **metrics** (*list*) – Define a list of additional cost functions to be returned. Can currently be ‘log’ and ‘insensitive’.
- **epsilon** (*float*) – insensitivity value.
- **return\_percentiles** (*boolean*) – Return some percentile statistics with the predictions.

## 23.3 catlearn.regression.gaussian\_process

Functions to make predictions with Gaussian Processes machine learning.

```
class catlearn.regression.gaussian_process.GaussianProcess (train_fp, train_target,  

    kernel_list, gradi-  

    ents=None, regu-  

    larization=None,  

    regulariza-  

    tion_bounds=None,  

    opti-  

    mize_hyperparameters=False,  

    scale_optimizer=False,  

    scale_data=False)
```

Bases: object

Gaussian processes functions for the machine learning.

```
optimize_hyperparameters (global_opt=False, algomin='L-BFGS-B', eval_jac=False,  

    loss_function='lml')
```

Optimize hyperparameters of the Gaussian Process.

This function assumes that the descriptors in the feature set remain the same. Optimization is performed with respect to the log marginal likelihood. Optimized hyperparameters are saved in the kernel dictionary. Finally, the covariance matrix is updated.

### Parameters

- **global\_opt** (*boolean*) – Flag whether to do basin hopping optimization of hyperparameters. Default is False.
- **algomin** (*str*) – Define scipy minimizer method to call. Default is L-BFGS-B.

```
predict (test_fp, test_target=None, uncertainty=False, basis=None, get_validation_error=False,  

    get_training_error=False, epsilon=None)
```

Function to perform the prediction on some training and test data.

### Parameters

- **test\_fp** (*list*) – A list of testing fingerprint vectors.
- **test\_target** (*list*) – A list of the the test targets used to generate the prediction errors.
- **uncertainty** (*boolean*) – Return data on the predicted uncertainty if True. Default is False.
- **basis** (*function*) – Basis functions to assess the reliability of the uncertainty predictions. Must be a callable function that takes a list of descriptors and returns another list.
- **get\_validation\_error** (*boolean*) – Return the error associated with the prediction on the test set of data if True. Default is False.

- **get\_training\_error** (*boolean*) – Return the error associated with the prediction on the training set of data if True. Default is False.
- **epsilon** (*float*) – Threshold for insensitive error calculation.

### Returns

**data** – Gaussian process predictions and meta data:

**prediction** [vector] Predicted mean.

**uncertainty** [vector] Predicted standard deviation of the Gaussian posterior.

**training\_error** [dictionary] Error metrics on training targets.

**validation\_error** [dictionary] Error metrics on test targets.

### Return type

dictionary

**predict\_uncertainty** (*test\_fp*)

Return uncertainty only.

**Parameters** **test\_fp** (*list*) – A list of testing fingerprint vectors.

**update\_data** (*train\_fp*, *train\_target=None*, *gradients=None*, *scale\_optimizer=False*)

Update the training matrix, targets and covariance matrix.

This function assumes that the descriptors in the feature set remain the same. That it is just the number of data points that is changing. For this reason the hyperparameters are not updated, so this update process should be fast.

### Parameters

- **train\_fp** (*list*) – A list of training fingerprint vectors.
- **train\_target** (*list*) – A list of training targets used to generate the predictions.
- **scale\_optimizer** (*boolean*) – Flag to define if the hyperparameters are log scale for optimization.

**update\_gp** (*train\_fp=None*, *train\_target=None*, *kernel\_list=None*, *scale\_optimizer=False*, *gradients=None*, *regularization\_bounds=(1e-06, None)*, *optimize\_hyperparameters=False*)

Potentially optimize the full Gaussian Process again.

This allows for the definition of a new kernel as a result of changing descriptors in the feature space. Other parts of the model can also be changed. The hyperparameters will always be reoptimized.

### Parameters

- **train\_fp** (*list*) – A list of training fingerprint vectors.
- **train\_target** (*list*) – A list of training targets used to generate the predictions.
- **kernel\_list** (*dict*) – This dict can contain many other dictionaries, each one containing parameters for separate kernels. Each kernel dict contains information on a kernel such as: - The ‘type’ key containing the name of kernel function. - The hyperparameters, e.g. ‘scaling’, ‘lengthscale’, etc.
- **scale\_optimizer** (*boolean*) – Flag to define if the hyperparameters are log scale for optimization.
- **regularization\_bounds** (*tuple*) – Optional to change the bounds for the regularization.

## 23.4 catlearn.regression.ridge\_regression

Modified ridge regression function from Keld Lundgaard.

```
class catlearn.regression.ridge_regression.RidgeRegression(W2=None, Vh=None, cv='loocv', Ns=100, wsteps=15, rsteps=3)
```

Bases: object

Ridge regression class to find an optimal model.

Regularization fitting can be performed with either the loocv or bootstrap.632 method. The loocv method is faster, but it is better to use bootstrap when there is highly correlated training data.

```
RR(X, Y, omega2, p=0.0, featselect_featvar=False)
```

Ridge Regression (RR) solver.

Cost is  $(Xa-y)^*2 + \text{omega2}*(a-p)^*2$ , SVD of  $X.T X$ , where T is the transpose V,  $W2, Vh = X.T*X$

### Parameters

- **x** (*array*) – Feature matrix for the training data.
- **y** (*list*) – Target data for the training sample.
- **p** (*float*) – Define the prior function.
- **omega2** (*float*) – Regularization strength.

### Returns

- **coefs** (*list*) – Optimal coefficients.
- **neff** (*float*) – Number of effective parameters.

```
bootstrap_calc(X, Y, p, omega2, samples, W2_samples, Vh_samples)
```

Calculate optimal omega2 from bootstrap.

### Parameters

- **x** (*array*) – Feature matrix for the training data.
- **y** (*list*) – Target data for the training sample.
- **p** (*float*) – Define the prior function.
- **omega2** (*float*) – Regularization strength.
- **samples** (*list*) – Sample index for bootstrap.
- **W2\_samples** (*array*) – Singular values for samples.
- **Vh\_samples** (*array*) – Right hand side of singular matrix for samples.

```
find_optimal_regularization(X, Y, p=0.0)
```

Find regularization value to minimize Expected Prediction Error.

### Parameters

- **x** (*array*) – Feature matrix for the training data.
- **y** (*list*) – Target data for the training sample.
- **p** (*float*) – Define the prior function. Default is zero.

**Returns** **omega2\_min** – Regularization corresponding to the minimum EPE.

**Return type** float

**get\_coefficients** (*train\_targets*, *train\_features*, *reg=None*, *p=0.0*)  
Generate the omgea2 and coef value's.

#### Parameters

- **train\_targets** (*array*) – Dependent data used for training.
- **train\_features** (*array*) – Independent data used for training.
- **reg** (*float*) – Precomputed optimal regularization.
- **p** (*float*) – Define the prior function. Default is zero.

**predict** (*train\_matrix*, *train\_targets*, *test\_matrix*, *test\_targets=None*, *coefficients=None*, *reg=None*, *p=0.0*)  
Function to do ridge regression predictions.

**regularization** (*train\_targets*, *train\_features*, *coef=None*, *featselect\_featvar=False*)  
Generate the omgea2 and coef value's.

Parameters **train\_targets** (*array*) – Dependent data used for training.

**train\_features** [*array*] Independent data used for training.

**coef** [*int*] List of indices in the feature database.

## 23.5 catlearn.regression.scikit\_wrapper

Regression models to assess features using scikit-learn framework.

**class** `catlearn.regression.scikit_wrapper.RegressionFit` (*train\_matrix*, *train\_target*,  
*test\_matrix=None*,  
*test\_target=None*,  
*method='ridge'*, *predict=False*)

Bases: `object`

Class to perform a fit to specified regression model.

**feature\_select** (*size=None*, *iterations=100000.0*, *steps=None*, *line\_search=False*, *min\_alpha=1e-08*, *max\_alpha=0.1*, *eps=0.001*)

Find index of important featurts.

#### Parameters

- **size** (*int*) – Number best features to return.
- **iterations** (*float*) – Maximum number of iterations taken minimizing the regression function. Implemented in elastic net and lasso.
- **steps** (*int*) – Number of steps to be taken in the penalty function of LASSO.
- **min\_alpha** (*float*) – Starting penalty when searching over range. Default is 1.e-8.
- **max\_alpha** (*float*) – Final penalty when searching over range. Default is 1.e-1.

# CHAPTER 24

---

## catlearn.active\_learning package

---

### 24.1 Submodules

#### 24.2 catlearn.active\_learning.acquisition\_functions module

GP acquisition functions.

```
catlearn.active_learning.acquisition_functions.EI(y_best, predictions, uncertainty, objective='max')
```

Return expected improvement acq. function.

##### Parameters

- **y\_best** (*float*) – Condition
- **predictions** (*list*) – Predicted means.
- **uncertainty** (*list*) – Uncertainties associated with the predictions.

```
catlearn.active_learning.acquisition_functions.PI(y_best, predictions, uncertainty, objective)
```

Probability of improvement acq. function.

##### Parameters

- **y\_best** (*float*) – Condition
- **predictions** (*list*) – Predicted means.
- **uncertainty** (*list*) – Uncertainties associated with the predictions.

```
catlearn.active_learning.acquisition_functions.UCB(predictions, uncertainty, objective='max', kappa=1.5)
```

Upper-confidence bound acq. function.

##### Parameters

- **predictions** (*list*) – Predicted means.

- **uncertainty** (*list*) – Uncertainties associated with the predictions.
- **kappa** (*float*) – Constant that controls the exploitation/exploration ratio in UCB.

```
catlearn.active_learning.acquisition_functions.classify(classifier, train_atoms,  
                  test_atoms, targets, predictions, uncertainty,  
                  train_features=None,  
                  test_features=None,  
                  objective='max',  
                  k_means=3,   kappa=1.5,  
                  metrics=['optimistic',  
                  'UCB', 'EI', 'PI'])
```

Classify ranked predictions based on acquisition function.

#### Parameters

- **classifier** (*func*) – User defined function to classify an atoms object.
- **train\_atoms** (*list*) – List of atoms objects from training data upon which to base classification.
- **test\_atoms** (*list*) – List of atoms objects from test data upon which to base classification.
- **targets** (*list*) – List of known target values.
- **predictions** (*list*) – List of predictions from the GP.
- **uncertainty** (*list*) – List of variance on the GP predictions.
- **train\_features** (*array*) – Feature matrix for the training data.
- **test\_features** (*array*) – Feature matrix for the test data.
- **k\_means** (*int*) – Number of cluster to generate with clustering.
- **kappa** (*float*) – Constant that controls the exploitation/exploration ratio in UCB.
- **metrics** (*list*) – list of strings. Accepted values are ‘cdf’, ‘UCB’, ‘EI’, ‘PI’, ‘optimistic’ and ‘pdf’.

**Returns** *res* – A dictionary of lists containing the fitness of each test point for the different acquisition functions.

#### Return type dict

```
catlearn.active_learning.acquisition_functions.cluster(train_features,       targets,  
                  test_features,      predictions,  
                  k_means=3)
```

Penalize test points that are too clustered.

#### Parameters

- **train\_features** (*array*) – Feature matrix for the training data.
- **targets** (*list*) – Training targets.
- **test\_features** (*array*) – Feature matrix for the test data.
- **predictions** (*list*) – Predicted means.
- **k\_means** (*int*) – Number of clusters.

```
catlearn.active_learning.acquisition_functions.optimistic(y_best, predictions, uncertainty)
```

Find predictions that will optimistically lead to progress.

**Parameters**

- **y\_best** (*float*) – Condition
- **predictions** (*list*) – Predicted means.
- **uncertainty** (*list*) – Uncertainties associated with the predictions.

```
catlearn.active_learning.acquisition_functions.optimistic_proximity(y_best,
                                                               predictions,
                                                               uncertainty)
```

Return uncertainties minus distances to y\_best.

**Parameters**

- **y\_best** (*float*) – Condition
- **predictions** (*list*) – Predicted means.
- **uncertainty** (*list*) – Uncertainties associated with the predictions.

```
catlearn.active_learning.acquisition_functions.probability_density(y_best,
                                                               predictions,
                                                               uncertainty)
```

Return probability densities at y\_best.

**Parameters**

- **y\_best** (*float*) – Condition
- **predictions** (*list*) – Predicted means.
- **uncertainty** (*list*) – Uncertainties associated with the predictions.

```
catlearn.active_learning.acquisition_functions.proximity(y_best, predictions, uncertainty=None)
```

Return negative distances to y\_best.

**Parameters**

- **y\_best** (*float*) – Condition
- **predictions** (*list*) – Predicted means.
- **uncertainty** (*list*) – Uncertainties associated with the predictions.

```
catlearn.active_learning.acquisition_functions.random_acquisition(y_best,
                                                               predictions,
                                                               uncertainty=None)
```

Return random numbers for control experiments.

**Parameters**

- **y\_best** (*float*) – Condition
- **predictions** (*list*) – Predicted means.
- **uncertainty** (*list*) – Uncertainties associated with the predictions.

```
catlearn.active_learning.acquisition_functions.rank(targets, predictions, uncertainty, train_features=None, test_features=None, objective='max', k_means=3, kappa=1.5, metrics=['optimistic', 'UCB', 'EI', 'PI'])
```

Rank predictions based on acquisition function.

#### Parameters

- **targets** (*list*) – List of known target values.
- **predictions** (*list*) – List of predictions from the GP.
- **uncertainty** (*list*) – List of variance on the GP predictions.
- **train\_features** (*array*) – Feature matrix for the training data.
- **test\_features** (*array*) – Feature matrix for the test data.
- **k\_means** (*int*) – Number of cluster to generate with clustering.
- **kappa** (*float*) – Constant that controls the exploitation/exploration ratio in UCB.
- **metrics** (*list*) – list of strings. Accepted values are ‘cdf’, ‘UCB’, ‘EI’, ‘PI’, ‘optimistic’ and ‘pdf’.

**Returns** **res** – A dictionary of lists containing the fitness of each test point for the different acquisition functions.

**Return type** dict

## 24.3 catlearn.active\_learning.algorithm module

Class to automate building a surrogate model.

```
class catlearn.active_learning.algorithm.ActiveLearning(surrogate_model, train_data, target)
```

Bases: object

Active learning class, intended for screening or optimizing in a predefined and finite search space.

```
acquire(unlabeled_data, batch_size=1)
```

Return indices of datapoints to acquire, from a predefined, finite search space.

#### Parameters

- **unlabeled\_data** (*array*) – Data matrix representing an unlabeled search space.
- **initial\_subset** (*list*) – Row indices of data to train on in the first iteration.
- **batch\_size** (*int*) – Number of training points to acquire (move from test to training) in every iteration.

#### Returns

- **to\_acquire** (*list*) – Row indices of unlabeled data to acquire.
- **score** – User defined output from predict.

```
ensemble_test(size, initial_subset=None, batch_size=1, n_max=None, seed_list=None, nprocs=None)
```

Return a 3d array of test results for a surrogate model. The third dimension expands the ensemble of tests.

**Parameters**

- **size** (*int*) – How many tests to run.
- **initial\_subset** (*list*) – Row indices of data to train on in the first iteration.
- **batch\_size** (*int*) – Number of training points to acquire (move from test to training) in every iteration.
- **n\_max** (*int*) – Max number of training points to test.
- **seed\_list** (*list*) – List of integer seeds for shuffling training data.
- **nprocs** (*int*) – Number of processors for parallelization

**Returns** `ensemble` – size by iterations by number of metrics array of test results.

**Return type** array

**test\_acquisition** (*initial\_subset=None*, *batch\_size=1*, *n\_max=None*, *seed=None*)

Return an array of test results for a surrogate model.

**Parameters**

- **initial\_subset** (*list*) – Row indices of data to train on in the first iteration.
- **batch\_size** (*int*) – Number of training points to acquire (move from test to training) in every iteration.
- **n\_max** (*int*) – Max number of training points to test.

## 24.4 Module contents



# CHAPTER 25

---

## catlearn.estimator package

---

### 25.1 Submodules

### 25.2 catlearn.estimator.general\_gp module

Function to setup a general GP.

```
class catlearn.estimator.general_gp.GeneralGaussianProcess(clean_type='eliminate',
                                                               dimension='single',
                                                               kernel='general')
```

Bases: object

Define a general setup for the Gaussian process.

This should not be used to try and obtain highly accurate solutions. Though it should give a reasonable model.

```
gaussian_process_predict(test_features)
```

Function to make GP predictions on tests data.

**Parameters** `test_features` (`array`) – The array of test features.

**Returns** `prediction` – The prediction data generated by the Gaussian process.

**Return type** dict

```
train_gaussian_process(train_features, train_targets)
```

Generate a general gaussian process model.

**Parameters**

- `train_features` (`array`) – The array of training features.
- `train_targets` (`array`) – A list of training target values.

**Returns** `gp` – The trained Gaussian process.

**Return type** object

## 25.3 catlearn.estimator.general\_kernel module

Setup a generic kernel.

```
catlearn.estimator.general_kernel.default_lengthscale(features, dimension='single')  
    Generate defaults for the kernel lengthscale.
```

### Parameters

- **features** (*array*) – The feature matrix for the training data.
- **dimension** (*str*) – The number of parameters to return. Can be ‘single’, or ‘features’.

**Returns** `std` – The standard deviation of the features.

**Return type** array

```
catlearn.estimator.general_kernel.general_kernel(features, dimension='single')  
    Generate a default kernel.
```

```
catlearn.estimator.general_kernel.smooth_kernel(features, dimension='single')  
    Generate a default kernel.
```

## 25.4 catlearn.estimator.general\_preprocess module

A default setup for data preprocessing.

```
class catlearn.estimator.general_preprocess.GeneralPreprocess(clean_type='eliminate')  
    Bases: object
```

A general purpose data preprocessing class.

```
process(train_features, train_targets, test_features=None)  
    Processing function.
```

### Parameters

- **train\_features** (*array*) – The array of training features.
- **train\_targets** (*array*) – A list of training target values.
- **test\_features** (*array*) – The array of test features.

```
transform(features)
```

Function to transform a new set of features.

**Parameters** `features` (*array*) – A new array of features to clean. This will most likely be the new test features.

**Returns** `processed` – A cleaned and scaled feature set.

**Return type** array

## 25.5 Module contents

# CHAPTER 26

---

## catlearn.utilities

---

### 26.1 catlearn.utilities.clustering

Simple k-means clustering.

```
catlearn.utilities.clustering.cluster_features(train_matrix, train_target, k=2,
                                               test_matrix=None, test_target=None)
```

Function to perform k-means clustering in the feature space.

#### Parameters

- **train\_matrix** (*list*) – Feature matrix for the training dataset.
- **train\_target** (*list*) – List of target values for training data.
- **k** (*int*) – Number of clusters to divide data into.
- **test\_matrix** (*list*) – Feature matrix for the test dataset.
- **test\_target** (*list*) – List of target values for test data.

### 26.2 catlearn.utilities.database\_functions

Functions to create databases storing feature matrix.

```
class catlearn.utilities.database_functions.DescriptorDatabase(db_name='descriptor_store.sqlite',
                                                               ta-
                                                               ble='Descriptors')
```

Bases: object

Store sets of descriptors for a given atoms object assigned a unique ID.

The descriptors for a given system can be stored in the ase.atoms object, though we typically find this method to be slower.

### `create_column(new_column)`

Function to create a new column in the table.

The new column will be initialized with None values.

**Parameters** `new_column(str)` – Name of new feature or target.

### `create_db(names)`

Function to setup a database storing descriptors.

**Parameters** `names(list)` – List of heading names for features and targets.

### `fill_db(descriptor_names, data)`

Function to fill the descriptor database.

#### Parameters

- `descriptor_names(list)` – List of descriptor names for features and targets.
- `data(array)` – First row should contain string of UUIDs, thereafter array should contain floats corresponding to the descriptor names provided.

### `get_column_names()`

Function to get the of a supplied table column names.

### `query_db(unique_id=None, names=None)`

Return single row based on uuid or all rows.

#### Parameters

- `unique_id(str)` – If specified, the data corresponding to the given UUID will be returned. If None, all rows will be returned.
- `names(list)` – If specified, only the data corresponding to provided column names will be returned. If None, all columns will be returned.

### `update_descriptor(descriptor, new_data, unique_id)`

Function to update a descriptor based on a given uuid.

#### Parameters

- `descriptor(str)` – Name of descriptor to be updated.
- `new_data(float)` – New value to be entered into table.
- `unique_id(str)` – The UUID of the entry to be updated.

`class catlearn.utilities.database_functions.FingerprintDB(db_name='fingerprints.db', verbose=False)`

A class for accessing a temporary SQLite database.

This function works as a context manager and should be used as follows:

`with FingerprintDB() as fpdb:` (Perform operation here)

This syntax will automatically construct the temporary database, or access an existing one. Upon exiting the indentation, the changes to the database will be automatically committed.

### `create_table()`

Create the database table framework used in SQLite.

This includes 3 tables: images, parameters, and fingerprints.

The images table currently stores ase\_id information and a unqiue string. This can be adapted in the future to support atoms objects.

The parameters table stores a symbol (10 character maximum) for convenient reference and a description of the parameter.

The fingerprints table holds a unique image and parameter ID along with a float value for each. The ID pair must be unique.

### **fingerprint\_entry (ase\_id, param\_id, value)**

Enter fingerprint value to database for given ase and parameter ID.

#### **Parameters**

- **ase\_id** (*int*) – The ase unique ID associated with an atoms object in the database.
- **param\_id** (*int or str*) – The parameter ID or symbol associated with and entry in the parameters table.
- **value** (*float*) – The value of the parameter for the atoms object.

### **get\_fingerprints (ase\_ids, params=[])**

Return values of provided parameters for each ase\_id provided.

#### **Parameters**

- **ase\_id** (*list*) – The ase ID(s) associated with an atoms object in the database.
- **params** (*list*) – List of symbols or int in parameters table to be selected.

**Returns** **fingerprint** – An array of values associated with the given parameters (a fingerprint) for each ase\_id.

#### **Return type** array

### **get\_parameters (selection=None, display=False)**

Return integer values corresponding to parameter IDs.

The array returned will be for a set of provided symbols. If no selection is provided, return all symbols.

#### **Parameters**

- **selection** (*list*) – List of symbols in parameters table to be selected.
- **display** (*bool*) – If True, print parameter descriptions.

**Returns** **res** – Return the integer values of selected parameters.

#### **Return type** array

### **image\_entry (asedb\_entry=None, identity=None)**

Enter a single ase-db image into the fingerprint database.

This table can be expanded to contain atoms objects in the future.

#### **Parameters**

- **d** (*object*) – An ase-db object which can be parsed.
- **identity** (*str*) – An identifier of the users choice.

**Returns** **d.id** – The ase ID collected for the ase-db object.

#### **Return type** int

### **parameter\_entry (symbol=None, description=None)**

Function for entering unique parameters into the database.

#### **Parameters**

- **symbol** (*str*) – A unique symbol the entry can be referenced by. If None, the symbol will be the ID of the parameter as a string.
- **description** (*str*) – A description of the parameter.

## 26.3 catlearn.utilities.distribution

Pair distribution function.

```
catlearn.utilities.distribution.pair_deviation(images, cutoffs, bins=33, bounds=None,  
                                              mic=True, element=None)
```

Return distribution of deviations from atom-pair nominal bond length.

### Parameters

- **images** (*list*) – List of atoms objects.
- **cutoffs** (*dictionary*) – Subtract elemental cutoff radii from distances. This is a useful for testing cutoff radii.
- **bins** (*int*) – Number of bins
- **bounds** (*tuple*) – Optional upper and lower bound of distances.
- **mic** (*boolean*) – Use minimum image convention. Set to False for non-periodic structures.
- **subset** (*list*) – Optionally select a subset of atomic indices to include.

```
catlearn.utilities.distribution.pair_distribution(images, bins=101, bounds=None,  
                                                mic=True, element=None)
```

Return the pair distribution function from a list of atoms objects.

### Parameters

- **images** (*list*) – List of atoms objects.
- **bins** (*int*) – Number of bins
- **bounds** (*tuple*) – Optional upper and lower bound of distances.
- **mic** (*boolean*) – Use minimum image convention. Set to False for non-periodic structures.
- **subset** (*list*) – Optionally select a subset of atomic indices to include.

## 26.4 catlearn.utilities.neighborlist

Functions to generate the neighborlist.

```
catlearn.utilities.neighborlist.ase_connectivity(atoms,                      cutoffs=None,  
                                               count_bonds=True)
```

Return a connectivity matrix calculated of an atoms object.

If no neighborlist or connectivity matrix is attached to the atoms object, a new one will be generated. Multiple connections are counted.

### Parameters

- **atoms** (*object*) – An ase atoms object.
- **cutoffs** (*list*) – A list of cutoff radii for the atoms, ordered by atom index.

**Returns** `conn` – An n by n, where n is len(atoms).

**Return type** array

```
catlearn.utilities.neighborlist.ase_neighborlist(atoms, cutoffs=None)
```

Make dict of neighboring atoms using ase function.

This provides a wrapper for the ASE neighborlist generator. Currently default values are used.

#### Parameters

- `atoms (object)` – Target ase atoms object on which to get neighbor list.
- `cutoffs (list)` – A list of radii for each atom in atoms.
- `rtol (float)` – The tolerance factor to allow for small variation in the cutoff radii.

**Returns** `neighborlist` – A dictionary containing the atom index and each neighbor index.

**Return type** dict

```
catlearn.utilities.neighborlist.catlearn_neighborlist(atoms, dx=None, max_neighbor=1, mic=True)
```

Make dict of neighboring atoms for discrete system.

Possible to return neighbors from defined neighbor shell e.g. 1st, 2nd, 3rd by changing the neighbor number.

#### Parameters

- `atoms (object)` – Target ase atoms object on which to get neighbor list.
- `dx (dict)` – Buffer to calculate nearest neighbor pairs in dict format: `dx = {atomic_number: buffer}`.
- `max_neighbor (int or str)` – Maximum neighbor shell. If int is passed this will define how many shells to consider. If ‘full’ is passed then all neighbor combinations will be included. This might get expensive for particularly large systems.

**Returns** `connection_matrix` – An array of the neighbor shell each atom index is located in.

**Return type** array

## 26.5 catlearn.utilities.penalty\_functions

Class with penalty functions.

```
class catlearn.utilities.penalty_functions.PenaltyFunctions(targets=None, predictions=None, uncertainty=None, train_features=None, test_features=None)
```

Bases: object

Base class for penalty functions.

```
penalty_close(c_min_crit=100000.0, d_min_crit=1e-05)
```

Penalize data that is too close.

Pass an array of test features and train features and returns an array of penalties due to ‘too short distance’ ensuring no duplicates are added.

#### Parameters

- `d_min_crit (float)` – Critical distance.

- **c\_min\_crit** (*float*) – Constant for penalty minimum distance.
- **penalty\_min** (*array*) – Array containing the penalty to add.

**penalty\_far** (*c\_max\_crit=100.0, d\_max\_crit=10.0*)

Penalize data that is too far.

Pass an array of test features and train features and returns an array of penalties due to ‘too far distance’. This prevents to explore configurations that are unrealistic.

#### Parameters

- **d\_max\_crit** (*float*) – Critical distance.
- **c\_max\_crit** (*float*) – Constant for penalty minimum distance.
- **penalty\_max** (*array*) – Array containing the penalty to add.

## 26.6 catlearn.utilities.sammon

Function to compute Sammon’s error between original and reduced features.

`catlearn.utilities.sammon.sammons_error(original, reduced)`

Sammon error.

#### Parameters

- **original** (*array*) – The original feature set.
- **reduced** (*array*) – The reduced feature set.

**Returns** `error` – Sammon’s error value.

**Return type** float

## 26.7 catlearn.utilities.utilities

Some useful utilities.

`catlearn.utilities.utilities.formal_charges(atoms, ion_number=8, ion_charge=-2)`

Return a list of formal charges on atoms.

#### Parameters

- **atoms** (*object*) – ase.Atoms object representing a chalcogenide. The default parameters are relevant for an oxide.
- **anion\_number** (*int*) – atomic number of anion.
- **anion\_charge** (*int*) – formal charge of anion.

**Returns** `all_charges` – Formal charges ordered by atomic index.

**Return type** list

`catlearn.utilities.utilities.geometry_hash(atoms)`

A hash based strictly on the geometry features of an atoms object.

Uses positions, cell, and symbols.

This is intended for planewave basis set calculations, so pbc is not considered.

Each element is sorted in the algorithm to help prevent new hashes for identical geometries.

```
catlearn.utilities.utilities.holdout_set(data, fraction, target=None, seed=None)
```

Return a dataset split in a hold out set and a training set.

#### Parameters

- **matrix** (*array*) – n by d array
- **fraction** (*float*) – fraction of data to hold out for testing.
- **target** (*list*) – optional list of targets or separate feature.
- **seed** (*float*) – optional float for reproducible splits.

```
catlearn.utilities.utilities.target_correlation(train, target, correlation=['pearson', 'spearman', 'kendall'])
```

Return the correlation of all columns of train with a target feature.

#### Parameters

- **train** (*array*) – n by d training data matrix.
- **target** (*list*) – target for correlation.

**Returns** **metric** – len(metric) by d matrix of correlation coefficients.

**Return type** array



# CHAPTER 27

---

## Indices and tables

---

- genindex
- modindex



---

## Python Module Index

---

### C

catlearn, 107  
catlearn.active\_learning, 97  
catlearn.active\_learning.acquisition\_functions, 93  
catlearn.active\_learning.algorithm, 96  
catlearn.api, 40  
catlearn.api.ase\_atoms\_api, 37  
catlearn.api.ase\_data\_setup, 39  
catlearn.api.networkx\_graph\_api, 39  
catlearn.cross\_validation, 43  
catlearn.cross\_validation.hierarchy\_cv, 41  
catlearn.cross\_validation.k\_fold\_cv, 42  
catlearn.estimator, 100  
catlearn.estimator.general\_gp, 99  
catlearn.estimator.general\_kernel, 100  
catlearn.estimator.general\_preprocess, 100  
catlearn.fingerprint, 55  
catlearn.fingerprint.adsorbate, 45  
catlearn.fingerprint.bulk, 49  
catlearn.fingerprint.chalcogenide, 49  
catlearn.fingerprint.convoluted, 50  
catlearn.fingerprint.graph, 51  
catlearn.fingerprint.molecule, 51  
catlearn.fingerprint.particle, 51  
catlearn.fingerprint.prototype, 52  
catlearn.fingerprint.standard, 53  
catlearn.fingerprint.voro, 54  
catlearn.ga, 61  
catlearn.ga.algorithm, 57  
catlearn.ga.convergence, 58  
catlearn.ga.initialize, 58  
catlearn.ga.io, 58  
catlearn.ga.mating, 59  
catlearn.ga.mutate, 59  
catlearn.ga.natural\_selection, 60  
catlearn.ga.predictors, 60

catlearn.learning\_curve, 68  
catlearn.learning\_curve.data\_process, 63  
catlearn.learning\_curve.feature\_selection, 64  
catlearn.learning\_curve.learning\_curve, 65  
catlearn.learning\_curve.placeholder, 67  
catlearn.preprocess, 78  
catlearn.preprocess.clean\_data, 69  
catlearn.preprocess.feature\_elimination, 70  
catlearn.preprocess.feature\_engineering, 72  
catlearn.preprocess.feature\_extraction, 74  
catlearn.preprocess.greedy\_elimination, 75  
catlearn.preprocess.importance\_testing, 76  
catlearn.preprocess.scaling, 77  
catlearn.regression, 92  
catlearn.regression.cost\_function, 88  
catlearn.regression.gaussian\_process, 89  
catlearn.regression.gpfunctions, 88  
catlearn.regression.gpfunctions.covariance, 79  
catlearn.regression.gpfunctions.default\_scale, 79  
catlearn.regression.gpfunctions.hyperparameter\_scale, 80  
catlearn.regression.gpfunctions.io, 80  
catlearn.regression.gpfunctions.kernel\_scaling, 81  
catlearn.regression.gpfunctions.kernel\_setup, 82  
catlearn.regression.gpfunctions.kernels, 83  
catlearn.regression.gpfunctions.log\_marginal\_likelihood, 83

```
    86
catlearn.regression.gpfunctions.sensitivity,
    87
catlearn.regression.gpfunctions.uncertainty,
    88
catlearn.regression.ridge_regression,
    91
catlearn.regression.scikit_wrapper, 92
catlearn.utilities, 107
catlearn.utilities.clustering, 101
catlearn.utilities.database_functions,
    101
catlearn.utilities.distribution, 104
catlearn.utilities.neighborlist, 104
catlearn.utilities.penalty_functions,
    105
catlearn.utilities.sammon, 106
catlearn.utilities.utilities, 106
```

---

## Index

---

### A

AA\_kernel() (in module `catlearn.regression.gpfunctions.kernels`), 83  
acquire() (`catlearn.active_learning.algorith.ActiveLearning` method), 96  
ActiveLearning (class in `catlearn.active_learning.algorith`), 96  
ads\_av() (`catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator` method), 45  
ads\_sum() (`catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator` method), 45  
AdsorbateFingerprintGenerator (class in `catlearn.fingerprint.adsorbate`), 45  
alpha\_finder() (`catlearn.learning_curve.feature_selection`.`feature_selection`.`method`), 64  
alpha\_refinement() (`catlearn.learning_curve.feature_selection`.`feature_selection`.`method`), 64  
ase\_connectivity() (in module `catlearn.utilities.neighborlist`), 104  
ase\_neighborlist() (in module `catlearn.utilities.neighborlist`), 105  
ase\_to\_networkx() (in module `catlearn.api.networkx_graph_api`), 39  
AutoCorrelationFingerprintGenerator (class in `catlearn.fingerprint.molecule`), 51  
average\_nested() (`catlearn.learning_curve.data_process`.`data_process`.`method`), 63  
bag\_edges() (in module `catlearn.fingerprint.standard.StandardFingerprintGenerator` method), 46  
bag\_edges\_ads() (in module `catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator` method), 46  
bag\_edges\_all() (in module `catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator` method), 46  
bag\_edges\_chemi() (in module `catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator` method), 46  
bag\_element\_cn() (in module `catlearn.fingerprint.standard.StandardFingerprintGenerator` method), 53  
bag\_elements() (in module `catlearn.fingerprint.standard.StandardFingerprintGenerator` method), 53  
bond\_count\_vec() (in module `catlearn.fingerprint.particle.ParticleFingerprintGenerator` method), 51  
bootstrap\_calc() (in module `catlearn.regression.ridge_regression.RidgeRegression` method), 91  
bulk() (in module `catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator` method), 46  
bulk\_average() (in module `catlearn.fingerprint.bulk.BulkFingerprintGenerator` method), 49  
bulk\_std() (in module `catlearn.fingerprint.bulk.BulkFingerprintGenerator` method), 49  
bulk\_summation() (in module `catlearn.fingerprint.bulk.BulkFingerprintGenerator` method), 49  
BulkFingerprintGenerator (class in `catlearn.fingerprint.bulk`), 49

### B

backward\_selection() (`catlearn.regression.gpfunctions.sensitivity.SensitivityAnalysis`.`module`), 107  
bag\_atoms\_ads() (`catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator`.`module`), 93  
bag\_cn() (`catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator`.`module`), 96  
bag\_cn\_general() (`catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator`.`module`), 40  
catlearn\_active\_learning.acquisition\_functions (in module `catlearn.active_learning`), 97  
catlearn\_active\_learning.algorith (in module `catlearn.active_learning.algorith`), 96  
catlearn\_api (in module `catlearn.api`), 40  
catlearn\_ase\_atoms\_api (in module `catlearn.api.ase_atoms_api`), 37

### C

catlearn.api.ase\_data\_setup (*module*), 39  
catlearn.api.networkx\_graph\_api (*module*), 39  
catlearn.cross\_validation (*module*), 43  
catlearn.cross\_validation.hierarchy\_cv (*module*), 41  
catlearn.cross\_validation.k\_fold\_cv (*module*), 42  
catlearn.estimator (*module*), 100  
catlearn.estimator.general\_gp (*module*), 99  
catlearn.estimator.general\_kernel (*module*), 100  
catlearn.estimator.general\_preprocess (*module*), 100  
catlearn.fingerprint (*module*), 55  
catlearn.fingerprint.adsorbate (*module*), 45  
catlearn.fingerprint.bulk (*module*), 49  
catlearn.fingerprint.chalcogenide (*module*), 49  
catlearn.fingerprint.convoluted (*module*), 50  
catlearn.fingerprint.graph (*module*), 51  
catlearn.fingerprint.molecule (*module*), 51  
catlearn.fingerprint.particle (*module*), 51  
catlearn.fingerprint.prototype (*module*), 52  
catlearn.fingerprint.standard (*module*), 53  
catlearn.fingerprint.voro (*module*), 54  
catlearn.ga (*module*), 61  
catlearn.ga.algorithm (*module*), 57  
catlearn.ga.convergence (*module*), 58  
catlearn.ga.initialize (*module*), 58  
catlearn.ga.io (*module*), 58  
catlearn.ga.mating (*module*), 59  
catlearn.ga.mutate (*module*), 59  
catlearn.ga.natural\_selection (*module*), 60  
catlearn.ga.predictors (*module*), 60  
catlearn.learning\_curve (*module*), 68  
catlearn.learning\_curve.data\_process (*module*), 63  
catlearn.learning\_curve.feature\_selectionatlearn.utilities.database\_functions (*module*), 64  
catlearn.learning\_curve.learning\_curve (*module*), 65  
catlearn.learning\_curve.placeholder (*module*), 67  
catlearn.preprocess (*module*), 78  
catlearn.preprocess.clean\_data (*module*), 69  
catlearn.preprocess.feature\_elimination (*module*), 70  
catlearn.preprocess.feature\_engineering (*module*), 72  
catlearn.preprocess.feature\_extraction (*module*), 74  
catlearn.preprocess.greedy\_elimination (*module*), 75  
catlearn.preprocess.importance\_testing (*module*), 76  
catlearn.preprocess.scaling (*module*), 77  
catlearn.regression (*module*), 92  
catlearn.regression.cost\_function (*module*), 88  
catlearn.regression.gaussian\_process (*module*), 89  
catlearn.regression.gpfunctions (*module*), 88  
catlearn.regression.gpfunctions.covariance (*module*), 79  
catlearn.regression.gpfunctions.default\_scale (*module*), 79  
catlearn.regression.gpfunctions.hyperparameter\_scale (*module*), 80  
catlearn.regression.gpfunctions.io (*module*), 80  
catlearn.regression.gpfunctions.kernel\_scaling (*module*), 81  
catlearn.regression.gpfunctions.kernel\_setup (*module*), 82  
catlearn.regression.gpfunctions.kernels (*module*), 83  
catlearn.regression.gpfunctions.log\_marginal\_likelihood (*module*), 86  
catlearn.regression.gpfunctions.sensitivity (*module*), 87  
catlearn.regression.gpfunctions.uncertainty (*module*), 88  
catlearn.regression.ridge\_regression (*module*), 91  
catlearn.regression.scikit\_wrapper (*module*), 92  
catlearn.utilities (*module*), 107  
catlearn.utilities.clustering (*module*), 101  
catlearn.utilities.database\_functions (*module*), 101  
catlearn.utilities.distribution (*module*), 104  
catlearn.utilities.neighborlist (*module*), 104  
catlearn.utilities.penalty\_functions (*module*), 105  
catlearn.utilities.sammon (*module*), 106  
catlearn.utilities.utilities (*module*), 106  
catlearn\_neighborlist() (*in module* catlearn.utilities.neighborlist), 105  
catlearn\_pca() (*in module*

*catlearn.preprocess.feature\_extraction), 74*

*ChalcogenideFingerprintGenerator (class in catlearn.fingerprint.chalcogenide), 49*

*check\_length() (in module catlearn.fingerprint.convolved), 50*

*classify() (in module catlearn.active\_learning.acquisition\_functions), 94*

*clean\_infinite() (in module catlearn.preprocess.clean\_data), 69*

*clean\_skewness() (in module catlearn.preprocess.clean\_data), 70*

*clean\_variance() (in module catlearn.preprocess.clean\_data), 70*

*cluster() (in module catlearn.active\_learning.acquisition\_functions), 94*

*cluster\_features() (in module catlearn.utilities.clustering), 101*

*composition\_vec() (catlearn.fingerprint.standard.StandardFingerprintGenerator method), 53*

*connections\_vec() (catlearn.fingerprint.particle.ParticleFingerprintGenerator method), 52*

*constant\_kernel() (in module catlearn.regression.gpfunctions.kernels), 83*

*constant\_multi\_kernel() (in module catlearn.regression.gpfunctions.kernels), 83*

*conv\_bulk() (catlearn.fingerprint.convolved.ConvolvedFingerprintGenerator method), 50*

*conv\_term() (catlearn.fingerprint.convolved.ConvolvedFingerprintGenerator method), 50*

*Convergence (class in catlearn.ga.convergence), 58*

*ConvolvedFingerprintGenerator (class in catlearn.fingerprint.convolved), 50*

*count\_chemisorbed\_fragment() (catlearn.fingerprint.adsorbate.AdsoarbateFingerprintGenerator method), 46*

*create\_column() (catlearn.utilities.database\_functions.DescriptorDatabaseFingerprint.adsorbate.AdsoarbateFingerprintGenerator method), 101*

*create\_db() (catlearn.utilities.database\_functions.DescriptorDatabaseFingerprint.adsorbate.AdsoarbateFingerprintGenerator method), 102*

*create\_table() (catlearn.utilities.database\_functions.FingerprintDBBiod), 47*

*ctime() (catlearn.fingerprint.adsorbate.AdsoarbateFingerprintGenerator method), 47*

*cut\_and\_splice() (in module catlearn.ga.mating), 59*

*data\_process (class in catlearn.learning\_curve.data\_process), 63*

*database\_to\_list() (in module catlearn.api.ase\_atoms\_api), 37*

*db\_size() (catlearn.fingerprint.adsorbate.AdsoarbateFingerprintGenerator method), 47*

*dbid() (catlearn.fingerprint.adsorbate.AdsoarbateFingerprintGenerator method), 47*

*default\_lengthscale() (in module catlearn.estimator.general\_kernel), 100*

*delta\_energy() (catlearn.fingerprint.adsorbate.AdsoarbateFingerprintGenerator method), 47*

*DescriptorDatabase (class in catlearn.utilities.database\_functions), 101*

*distance\_vec() (catlearn.fingerprint.standard.StandardFingerprintGenerator method), 54*

*distribution\_vec() (catlearn.fingerprint.particle.ParticleFingerprintGenerator method), 52*

*dK\_dtheta\_j() (in module catlearn.regression.gpfunctions.log\_marginal\_likelihood), 86*

**E**

*eigenspectrum\_vec() (catlearn.fingerprint.standard.StandardFingerprintGenerator method), 93*

*element\_mass\_vec() (catlearn.fingerprint.standard.StandardFingerprintGenerator method), 54*

*eliminate\_features() (catlearn.preprocess.feature\_elimination.FeatureScreening method), 71*

*en\_difference\_active() (catlearn.fingerprint.adsorbate.AdsoarbateFingerprintGenerator method), 47*

*en\_difference\_ads() (catlearn.fingerprint.adsorbate.AdsoarbateFingerprintGenerator method), 47*

*ensemble\_test() (catlearn.active\_learning.algorithm.ActiveLearning method), 96*

*extend\_atoms\_class() (in module catlearn.api.ase\_atoms\_api), 37*

**F**

*feature\_frequency() (in module catlearn.learning\_curve.learning\_curve),*

66  
 feature\_inspection() (in module `catlearn.learning_curve.feature_selection.feature_selection_kernel()` (in module `catlearn.estimator.general_kernel`), 65)  
 feature\_invariance() (in module `catlearn.preprocess.importance_testing`), 77  
 feature\_randomize() (in module `catlearn.preprocess.importance_testing`), 77  
 feature\_select() (in module `catlearn.regression.scikit_wrapper.RegressionFit` (in module `catlearn.fingerprint.prototype.PrototypeFingerprintGenerator`), method), 92  
 feature\_selection (class in `catlearn.learning_curve.feature_selection`), 64  
 feature\_shuffle() (in module `catlearn.preprocess.importance_testing`), 77  
 FeatureScreening (class in `catlearn.preprocess.feature_elimination`), 70  
 fill\_db() (in module `catlearn.utilities.database_functions.DescriptorDatabase` (class in `catlearn.ga.algorithm`), method), 102  
 find\_optimal\_regularization() (in module `catlearn.regression.ridge_regression.RidgeRegression` (class in `catlearn.utilities.utilities`), method), 91  
 fingerprint\_entry() (in module `catlearn.utilities.database_functions.FingerprintDB` (class in `catlearn.fingerprint.chalcogenide.ChalcogenideFingerprintGenerator`), method), 103  
 FingerprintDB (class in `catlearn.utilities.database_functions`), 102  
 fitness (in module `catlearn.ga.algorithm.GeneticAlgorithm` attribute), 57  
 formal\_charges() (in module `catlearn.fingerprint.chalcogenide.ChalcogenideFingerprintGenerator` (class in `catlearn.utilities.utilities`), method), 50  
 formal\_charges() (in module `catlearn.utilities.utilities`), 106

**G**

gaussian\_dk\_dwidth() (in module `catlearn.regression.gpfunctions.kernels`), 84  
 gaussian\_kernel() (in module `catlearn.regression.gpfunctions.kernels`), 84  
 gaussian\_process\_predict() (in module `catlearn.estimator.general_gp.GeneralGaussianProcess` (class in `catlearn.regression.cost_function`), method), 99  
 gaussian\_xx\_gradients() (in module `catlearn.regression.gpfunctions.kernels`), 84  
 gaussian\_xxp\_gradients() (in module `catlearn.regression.gpfunctions.kernels`), 84

GaussianProcess (class in `catlearn.regression.gaussian_process`), 89  
 GeneralGaussianProcess (class in `catlearn.estimator.general_gp`), 99  
 generalized\_cn() (in module `catlearn.fingerprint.adsorbate.AdssorbateFingerpr` method), 48  
 GeneralPreprocess (class in `catlearn.estimator.general_preprocess`), 100  
 generate() (in module `catlearn.fingerprint.voro.VoronoiFingerprintGenerator` method), 54  
 generate\_all() (in module `catlearn.fingerprint.prototype.PrototypeFingerprintG` method), 52  
 generate\_features() (in module `catlearn.preprocess.feature_engineering`), 72  
 generate\_positive\_features() (in module `catlearn.preprocess.feature_engineering`), 72  
 geometry\_hash() (in module `catlearn.utilities.utilities`), 106  
 get\_ablog() (in module `catlearn.preprocess.feature_engineering`), 73  
 get\_acorrelation() (in module `catlearn.fingerprint.molecule.AutoCorrelationFingerprintGenerat` method), 51  
 get\_coefficients() (in module `catlearn.regression.ridge_regression.RidgeRegression` (class in `catlearn.utilities.utilities`), method), 102  
 get\_column\_names() (in module `catlearn.utilities.database_functions.DescriptorDatabase` (class in `catlearn.utilities.utilities`), method), 102  
 get\_covariance() (in module `catlearn.regression.gpfunctions.covariance`), 79  
 get\_data\_scale() (in module `catlearn.learning_curve.placeholder.placeholder` (class in `catlearn.utilities.utilities`), method), 67  
 get\_div\_order\_2() (in module `catlearn.preprocess.feature_engineering`), 73  
 get\_error() (in module `catlearn.regression.cost_function`), 88  
 get\_features() (in module `catlearn.api.ase_atoms_api`), 37  
 get\_fingerprints() (in module `catlearn.utilities.database_functions.FingerprintDB` (class in `catlearn.fingerprint.chalcogenide.ChalcogenideFingerprintGenerator`), method), 103  
 get\_graph() (in module `catlearn.api.ase_atoms_api`), 37

get\_labels\_ablog() (in module `catlearn.preprocess.feature_engineering`), 73  
 get\_labels\_order\_2() (in module `catlearn.preprocess.feature_engineering`), 73  
 get\_labels\_order\_2ab() (in module `catlearn.preprocess.feature_engineering`), 73  
 get\_neighborlist() (in module `catlearn.api.ase_atoms_api`), 37  
 get\_order\_2() (in module `catlearn.preprocess.feature_engineering`), 73  
 get\_order\_2ab() (in module `catlearn.preprocess.feature_engineering`), 74  
 get\_parameters() (`catlearn.utilities.database_functions.FingerprintDBn.preprocess.importance_testing`), method), 103  
 get\_statistic() (`catlearn.learning_curve.data_process`), data\_population() (in module `catlearn.ga.initialize`), 58  
 get\_subset\_data() (in module `catlearn.cross_validation.hierarchy_cv.Hierarchy`), method), 41  
 get\_train() (in module `catlearn.api.ase_data_setup`), 39  
 get\_uncertainty() (in module `catlearn.regression.gpfunctions.uncertainty`), 88  
 get\_unique() (in module `catlearn.api.ase_data_setup`), 39  
 getstats() (in module `catlearn.learning_curve.placeholder`), placeholder, method), 67  
 globalscaledata() (in module `catlearn.cross_validation.hierarchy_cv.Hierarchy`), method), 41  
 globalscaling() (in module `catlearn.learning_curve.data_process`), data\_process, method), 63  
 GraphFingerprintGenerator (class in `catlearn.fingerprint.graph`), 51  
 greedy\_elimination() (in module `catlearn.preprocess.greedy_elimination`), GreedyElimination, method), 75  
 GreedyElimination (class in `catlearn.preprocess.greedy_elimination`), 75

**H**

Hierarchy (class in `catlearn.cross_validation.hierarchy_cv`), 41  
 hierarchy() (in module `catlearn.learning_curve.learning_curve`), 66

**I**

holdout\_set() (in module `catlearn.utilities.utilities`), 107  
 hyperparameters() (in module `catlearn.regression.gpfunctions.hyperparameter_scaling`), 80

**J**

image\_entry() (`catlearn.utilities.database_functions.FingerprintDB`), method), 103  
 images\_connectivity() (in module `catlearn.api.ase_atoms_api`), 38  
 images\_pair\_distances() (in module `catlearn.api.ase_atoms_api`), 38  
 importance\_elimination() (`catlearn.preprocess.importance_testing`.ImportanceElimination method), 76  
 ImportanceElimination (class in `catlearn.utilities.database_functions.FingerprintDBn.preprocess.importance_testing`), 76  
 interval\_modifier() (in module `catlearn.learning_curve.feature_selection`.feature\_selection method), 65  
 iterative\_screen() (`catlearn.preprocess.feature_elimination`.FeatureScreening method), 71

**K**

k\_fold() (in module `catlearn.cross_validation.k_fold_cv`), 42  
 kdicts2list() (in module `catlearn.regression.gpfunctions.kernel_setup`), 82  
 kernel\_scaling() (in module `catlearn.regression.gpfunctions.kernel_scaling`), 81  
 laplacian\_dk\_dwidth() (in module `catlearn.regression.gpfunctions.kernels`), 84  
 laplacian\_kernel() (in module `catlearn.regression.gpfunctions.kernels`), 84

**L**

LearningCurve (class in `catlearn.learning_curve.learning_curve`), 65  
 linear\_kernel() (in module `catlearn.regression.gpfunctions.kernels`), 85

```

list2kdict()           (in module no_progress() (catlearn.ga.convergence.Convergence
    catlearn.regression.gpfunctions.kernel_setup), method), 58
    82
    noise_multi_kernel()      (in module
load_split() (catlearn.cross_validation.hierarchy_cv.Hierarchy catlearn.regression.gpfunctions.kernels),
    method), 41
    85
log_marginal_likelihood() (in module normalize() (in module catlearn.preprocess.scaling),
    catlearn.regression.gpfunctions.log_marginal_likelihood), 78
    87

```

## O

```

M
matrix_to_nl()           (in module optimistic() (in module
    catlearn.api.networkx_graph_api), 39
    catlearn.active_learning.acquisition_functions),
    94
max_cation() (catlearn.fingerprint.chalcogenide.ChalcogenideFingerprintGenerator) (in module
    method), 50
    catlearn.active_learning.acquisition_functions,
max_site() (catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator
    method), 48
    optimize_hyperparameters()
mean_cation() (catlearn.fingerprint.chalcogenide.ChalcogenideFingerprintGenerator gaussian_process.GaussianProcess
    method), 50
    method), 89
mean_chemisorbed_atoms()
    (catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator
        method), 48
    pair_deviation() (in module
mean_site() (catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator catlearn.utilities.distribution), 104
    method), 48
    pair_distribution() (in module
mean_surf_ligands()
    (catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator catlearn.utilities.distribution), 104
        method), 48
    parameter_entry()
        (catlearn.utilities.database_functions.FingerprintDB
median_cation() (catlearn.fingerprint.chalcogenide.ChalcogenideFingerprintGenerator
    method), 50
    ParticleFingerprintGenerator (class in
median_site() (catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator catlearn.fingerprint.particle), 51
    method), 48
    pca() (in module catlearn.preprocessing.feature_extraction),
min_cation() (catlearn.fingerprint.chalcogenide.ChalcogenideFingerprintGenerator
    method), 50
    penalty_close() (catlearn.utilities.penalty_functions.PenaltyFunction
min_max() (in module catlearn.preprocessing.scaling), 77
    method), 105
min_site() (catlearn.fingerprint.adsorbate.AdsorbateFingerprintGenerator penality_generator (catlearn.utilities.penalty_functions.PenaltyFunctions
    method), 48
    method), 106
minimize_error() (in module PenaltyFunctions (class in
    catlearn.ga.predictors), 60
    catlearn.utilities.penalty_functions), 105
minimize_error_descriptors() (in module PI() (in module catlearn.active_learning.acquisition_functions),
    catlearn.ga.predictors), 60
    93
minimize_error_time() (in module placeholder (class in
    catlearn.ga.predictors), 61
    catlearn.learning_curve.placeholder), 67
    pls() (in module catlearn.preprocessing.feature_extraction),
    75

```

## N

```

nearestneighbour_vec() population (catlearn.ga.algorithm.GeneticAlgorithm
    (catlearn.fingerprint.particle.ParticleFingerprintGenerator attribute), 57
    method), 52
    population_reduction() (in module
neighbor_mean_vec() (catlearn.fingerprint.graph.GraphFingerprintGenerator predict() (catlearn.regression.gaussian_process.GaussianProcess
    method), 51
    method), 89
    predict())
neighbor_sum_vec() (catlearn.fingerprint.graph.GraphFingerprintGenerator predict() (catlearn.regression.ridge_regression.RidgeRegression
    method), 92
    method), 92
    predict_subsets())
networkx_to_adjacency() (in module predict_subsets() (catlearn.learning_curve.placeholder.placeholder
    catlearn.api.networkx_graph_api), 39
    method), 67
    method), 67

```

```

predict_uncertainty()                               read_train_data()           (in      module
    (catlearn.regression.gaussian_process.GaussianProcess   catlearn.regression.gpfunctions.io), 81
    method), 90                                         reg_data_var() (catlearn.learning_curve.placeholder.placeholder
                                                               method), 67
prediction_error()                                feat_var() (catlearn.learning_curve.placeholder.placeholder
    (catlearn.learning_curve.data_process.data_process     method), 68
    method), 64
prepare_kernels()        (in      module  RegressionFit       (class      in
    catlearn.regression.gpfunctions.kernel_setup),       catlearn.regression.scikit_wrapper), 92
    82                                             regularization() (catlearn.regression.ridge_regression.RidgeRegress
                                                               method), 92
probability_density()   (in      module  remove_duplicates() (in      module
    catlearn.active_learning.acquisition_functions),       catlearn.ga.natural_selection), 60
    95
probability_include()  (in      module  remove_outliers() (in      module
    catlearn.ga.mutate), 59                           catlearn.preprocess.clean_data), 70
probability_remove()  (in      module  rescale_hyperparameters() (in      module
    catlearn.ga.mutate), 59                          catlearn.regression.gpfunctions.hyperparameter_scaling),
process() (catlearn.estimator.general_preprocess.GeneralPreprocess
method), 100                                         rescale_targets()
PrototypeFingerprintGenerator (class  in      (catlearn.regression.gpfunctions.default_scale.ScaleData
    catlearn.fingerprint.prototype), 52               method), 80
PrototypeSites          (class  in      RidgeRegression       (class      in
    catlearn.fingerprint.prototype), 53               catlearn.regression.ridge_regression), 91
proximity()            (in      module  RR() (catlearn.regression.ridge_regression.RidgeRegression
    catlearn.active_learning.acquisition_functions),       method), 91
    95
Q
quadratic_dk_ddegree() (in      module  run() (catlearn.learning_curve.learning_curve.LearningCurve
    catlearn.regression.gpfunctions.kernels),       method), 65
    85
quadratic_dk_dslope() (in      module  run_proto() (catlearn.fingerprint.prototype.PrototypeFingerprintGener
    catlearn.regression.gpfunctions.kernels),       method), 52
    85
quadratic_kernel()     (in      module  run_voro() (catlearn.fingerprint.voro.VoronoiFingerprintGenerator
    catlearn.regression.gpfunctions.kernels),       method), 55
    86
query_db() (catlearn.utilities.database_functions.DescriptorDatabase
method), 102
R
random_acquisition() (in      module  S
    catlearn.active_learning.acquisition_functions),  scaling_data() (catlearn.learning_curve.data_process.data_process
    95                                         method), 64
random_permutation()   (in      module  screen() (catlearn.preprocess.feature_elimination.FeatureScreening
    catlearn.ga.mutate), 59                         method), 71
rank() (in module catlearn.active_learning.acquisition_functions, method), 57
selection() (catlearn.learning_curve.feature_selection.feature_selection
method), 65
rdf_vec() (catlearn.fingerprint.particle.ParticleFingerprintGenerator
method), 52
read() (in module catlearn.regression.gpfunctions.io), 80
read_data() (in module catlearn.ga.io), 58
read_split() (in      module  SensitivityAnalysis (class      in
    catlearn.cross_validation.k_fold_cv), 42           catlearn.regression.gpfunctions.sensitivity),
                                                               87
set_features() (in      module  set_features() (in      module
    catlearn.api.ase_atoms_api), 38

```

set\_graph() (in module `catlearn.api.ase_atoms_api`), 38  
 set\_neighborhood() (in module `catlearn.api.ase_atoms_api`), 38  
 single\_transform() (in module `catlearn.preprocess.feature_engineering`), 74  
 smooth\_kernel() (in module `catlearn.estimator.general_kernel`), 100  
 spca() (in module `catlearn.preprocess.feature_extraction`)  
 split\_index() (in module `catlearn.cross_validation.hierarchy_cv.Hierarchy`)  
 split\_predict() (in module `catlearn.cross_validation.hierarchy_cv.Hierarchy`)  
 sqe\_kernel() (in module `catlearn.regression.gpfunctions.kernels`), 86  
 stagnation() (in module `catlearn.ga.convergence`), 58  
 StandardFingerprintGenerator (class in `catlearn.fingerprint.standard`), 53  
 standardize() (in module `catlearn.preprocess.scaling`), 78  
 strain() (in module `catlearn.fingerprint.adsorbate`), 49  
 sum\_cation() (in module `catlearn.fingerprint.chalcogenide`), 50  
 sum\_site() (in module `catlearn.fingerprint.adsorbate`), 49

**T**

target\_center() (in module `catlearn.preprocess.scaling`), 78  
 target\_correlation() (in module `catlearn.utilities.utilities`), 107  
 target\_normalize() (in module `catlearn.preprocess.scaling`), 78  
 target\_standardize() (in module `catlearn.preprocess.scaling`), 78  
 term() (in module `catlearn.fingerprint.adsorbate`), 49  
 test() (in module `catlearn.regression.gpfunctions.default_scale`), 80  
 test\_acquisition()  
     (catlearn.active\_learning.algorithm.ActiveLearning method), 97  
 todb() (in module `catlearn.cross_validation.hierarchy_cv.Hierarchy`), 42  
 train() (in module `catlearn.regression.gpfunctions.default_scale`), 80  
 train\_gaussian\_process()  
     (catlearn.estimator.general\_gp.GeneralGaussianProcess method), 99

transform() (in module `catlearn.estimator.general_preprocess`), 100  
 transform\_output() (in module `catlearn.cross_validation.hierarchy_cv.Hierarchy`)  
 unit\_length() (in module `catlearn.preprocess.scaling`), 78  
 update\_data() (in module `catlearn.regression.gaussian_process.GaussianProcess`)  
 update\_descriptor() (in module `catlearn.utilities.database_functions`), 102  
 update\_gp() (in module `catlearn.regression.gaussian_process.GaussianProcess`), 90  
 update\_str() (in module `catlearn.fingerprint.prototype`), 52

**U**

UCB() (in module `catlearn.active_learning.acquisition_functions`), 93

**V**

VoronoiFingerprintGenerator (class in `catlearn.fingerprint.voro`), 54

**W**

write() (in module `catlearn.regression.gpfunctions.io`), 81  
 write\_proto\_input()  
     (catlearn.fingerprint.prototype.PrototypeFingerprintGenerator method), 52  
 write\_split() (in module `catlearn.cross_validation.k_fold_cv`), 43  
 write\_train\_data() (in module `catlearn.regression.gpfunctions.io`), 81  
 write\_voro\_input()  
     (catlearn.fingerprint.voro.VoronoiFingerprintGenerator method), 55

**X**

xyz\_id() (in module `catlearn.fingerprint.bulk`), 49