

---

# PicturedRocks Documentation

*Release 0.2.0+0.g87ae8a0.dirty*

**Anna Gilbert, Alexander Vargo, Umang Varma**

Sep 11, 2018



---

## Contents:

---

<b>1</b>	<b>Installing</b>	<b>3</b>
1.1	Reading data . . . . .	3
1.2	Preprocessing . . . . .	4
1.3	Plotting . . . . .	4
1.4	Selecting Markers . . . . .	5
1.5	Measuring Feature Selection Performance . . . . .	7
<b>2</b>	<b>Indices and tables</b>	<b>11</b>
	<b>Python Module Index</b>	<b>13</b>



PicturedRocks is a tool for the analysis of single cell RNA-seq data. Currently, we implement two marker selection approaches:

- a 1-bit compressed sensing based sparse SVM algorithm, and
- a mutual information-based greedy feature selection algorithm.



# CHAPTER 1

---

## Installing

---

Please ensure you have Python 3.6 or newer and have *numba* and *scikit-learn* installed. The best way to get Python and various dependencies is with Anaconda or Miniconda. Once you have a conda environment, run `conda install numba scikit-learn`. Then use pip to install PicturedRocks and all additional dependencies:

```
pip install picturedrocks
```

To install the latest code from github, clone our github repository. Once inside the project directory, instal by running `pip install -e ..`. The `-e` option will point a symbolic link to the current directory instead of installing a copy on your system.

## 1.1 Reading data

In addition to various functions for reading input data in `scanpy`, various methods in `picturedrocks` need cluster labels.

`picturedrocks.read.process_clusts(adata, copy=False)`

Annotate with information about clusters

Precomputes cluster indices, number of clusters, etc.

### Parameters

- `adata (anndata.Anndata)` –
- `copy (bool)` – determines whether a copy of `AnnData` object is returned

**Returns** object with annotation

**Return type** `anndata.Anndata`

### Notes

The information computed here is lost when saving as a `.loom` file. If a `.loom` file has cluster information, you should run this function immediately after `sc.read_loom`.

```
picturedrocks.read.read_clusts(adata, filename, sep=', ', copy=False)
```

Read cluster labels from a csv

#### Parameters

- **adata** (`anndata.Anndata`) – the *AnnData* object to read labels into
- **filename** (`str`) – filename of the csv file with labels
- **sep** (`str, optional`) – csv delimiter
- **copy** (`bool`) – determines whether a copy of *AnnData* object is returned

**Returns** object with cluster labels

**Return type** `anndata.Anndata`

#### Notes

- Cluster ids will automatically be changed so they are 0-indexed
- csv can either be two columns (in which case the first column is treated as observation label and merging handled by pandas) or one column (only cluster labels, ordered as in `adata`)

## 1.2 Preprocessing

The preprocessing module provides basic preprocessing tools. To avoid reinventing the wheel, we won't repeat methods already in `scipy` unless we need functionality not available there.

```
picturedrocks.preprocessing.pca(data, dim=3, center=True, copy=False)
```

Runs PCA

#### Parameters

- **data** (`anndata.Anndata`) – input data
- **dim** (`int, optional`) – number of PCs to compute
- **center** (`bool, optional`) – determines whether to center data before running PCA
- **copy** – determines whether a copy of *AnnData* object is returned

**Returns** object with `obsm["X_pca"]`, `varm["PCs"]`, and `uns["num_pcs"]` set

**Return type** `anndata.Anndata`

## 1.3 Plotting

```
picturedrocks.plot.genericplot(celldata, coords, **scatterkwargs)
```

Generate a figure for some embedding of data

This function supports both 2D and 3D plots. This may be used to plot data for any embedding (e.g., PCA or t-SNE). For example usage, see code for `pcafigure`.

#### Parameters

- **celldata** (`anndata.Anndata`) – data to plot
- **coords** (`numpy.ndarray`) – (N, 2) or (N, 3) shaped coordinates of the embedded data

- **\*\*scatterkwargs** – keyword arguments to pass to Scatter or Scatter3D in *plotly* (dictionaries are merged recursively)

`picturedrocks.plot.genericwrongplot(celldata, coords, yhat, labels=None, **scatterkwargs)`

Plot figure with incorrectly classified points highlighted

This can be used with any 2D or 3D embedding (e.g., PCA or t-SNE). For example code, see `pcawrongplot`.

#### Parameters

- **celldata** (`anndata.Anndata`) – data to plot
- **coords** (`numpy.ndarray`) – (N, 2) or (N, 3) shaped array with coordinates to plot
- **yhat** (`numpy.ndarray`) – (N, 1) shaped array of predicted y values
- **labels** (`list, optional`) – list of axis titles
- **\*\*scatterkwargs** – keyword arguments to pass to Scatter or Scatter3D in *plotly* (dictionaries are merged recursively)

`picturedrocks.plot.pcaffigure(celldata, **scatterkwargs)`

Make a 3D PCA figure for an AnnData object

#### Parameters

- **celldata** (`anndata.Anndata`) – data to plot
- **\*\*scatterkwargs** – keyword arguments to pass to Scatter or Scatter3D in *plotly* (dictionaries are merged recursively)

`picturedrocks.plot.pcawrongplot(celldata, yhat, **scatterkwargs)`

Generate a 3D PCA figure with incorrectly classified points highlighted

#### Parameters

- **celldata** (`anndata.Anndata`) – data to plot
- **yhat** (`numpy.ndarray`) – (N, 1) shaped array of predicted y values
- **\*\*scatterkwargs** – keyword arguments to pass to Scatter or Scatter3D in *plotly* (dictionaries are merged recursively)

`picturedrocks.plot.plotgeneheat(celldata, coords, genes, hide_clusts=False, **scatterkwargs)`

Generate gene heat plot for some embedding of AnnData

This generates a figure with multiple dropdown options. The first option is “Clust” for a plot similar to `genericplot`, and the remaining dropdown options correspond to genes specified in `genes`. When `celldata.genes` is defined, these drop downs are labeled with the gene names.

#### Parameters

- **celldata** (`anndata.Anndata`) – data to plot
- **coords** (`numpy.ndarray`) – (N, 2) or (N, 3) shaped coordinates of the embedded data (e.g., PCA or tSNE)
- **genes** (`list`) – list of gene indices or gene names
- **hide\_clusts** (`bool`) – Determines if cluster labels are ignored even if they are available

## 1.4 Selecting Markers

PicturedRocks current implements two categories of marker selection algorithms:

- mutual information-based algorithms
- 1-bit compressed sensing based algorithms

### 1.4.1 Mutual information

TODO: Explanation of how these work goes here.

Before running any mutual information based algorithms, we need a discretized version of the gene expression matrix, with a limited number of discrete values (because we do not make any assumptions about the distribution of gene expression). Such data is stored in `picturedrocks.markers.InformationSet`, but by default, we suggest using `picturedrocks.markers.makeinfoset()` to generate such an object after appropriate normalization

```
class picturedrocks.markers.mutualinformation.iterative.MIM(infoset)
class picturedrocks.markers.mutualinformation.iterative.CIFE(infoset)
class picturedrocks.markers.mutualinformation.iterative.JMI(infoset)
class picturedrocks.markers.mutualinformation.iterative.UniEntropy(infoset)
class picturedrocks.markers.mutualinformation.iterative.CIFEUnsup(infoset)
```

## Auxiliary Classes and Methods

`class picturedrocks.markers.InformationSet(X, has_y=False)`  
Stores discrete gene expression matrix

### Parameters

- `X (numpy.ndarray)` – a (num\_obs, num\_vars) shape array with dtype `int`
- `has_y (bool)` – whether the array `X` has a target label column (a `y` column) as its last column

`picturedrocks.markers.makeinfoset(adata, include_y)`  
Discretize data

### Parameters

- `adata (anndata.Anndata)` – The data to discretize. By default data is discretized as  $\text{round}(\log_2(X + 1))$ .
- `include_y (bool)` – Determines if the `y` (cluster label) column is included in the `InformationSet` object

**Returns** An object that can be used to perform information theoretic calculations.

**Return type** `picturedrocks.markers.mutualinformation.infoset.InformationSet`

## 1.4.2 Interactive Marker Selection

```
class picturedrocks.markers.interactive.InteractiveMarkerSelection(adata,
    fea-
    ture_selection,
    disp_genes=10,
    con-
    nected=True,
    show_cells=True,
    show_genes=True)
```

Run an interactive marker selection GUI inside a jupyter notebook

### Parameters

- **adata** (`anndata.Anndata`) – The data to run marker selection on. If you want to restrict to a small number of genes, slice your anndata object.
- **feature\_selection** (`picturedrocks.markers.mutualinformation.iterative.IterativeFeatureSelection`) – An instance of a interative feature selection algorithm class that corresponds to `adata` (i.e., the column indices in `feature_selection` should correspond to the column indices in `adata`)
- **disp\_genes** (`int`) – Number of genes to display as options (by default, number of genes plotted on the tSNE plot is  $3 * \text{disp\_genes}$ , but can be changed by setting the `plot_genes` property after initializing).
- **connected** (`bool`) – Parameter to pass to `plotly.offline.init_notebook_mode`. If your browser does not have internet access, you should set this to False.
- **show\_cells** (`bool`) – Determines whether to display a tSNE plot of the cells with a drop-down menu to look at gene expression levels for candidate genes.
- **show\_genes** (`bool`) – Determines whether to display a tSNE plot of genes to visualize gene similarity

**Warning:** This class requires modules not explicitly listed as dependencies of picturedrocks. Specifically, please ensure that you have `ipywidgets` installed and that you use this class only inside a jupyter notebook.

`picturedrocks.markers.interactive.cife_obj(H, i, S)`

The CIFE objective function for feature selection

### Parameters

- **H** (`function`) – an entropy function, typically the bound method `H` on an instance of `InformationSet`. For example, if `infoset` is of type `picturedrocks.markers.InformationSet`, then pass `infoset.H`
- **i** (`int`) – index of candidate gene
- **S** (`list`) – list of features already selected

**Returns** the candidate feature's score relative to the selected gene set  $S$

**Return type** `float`

## 1.5 Measuring Feature Selection Performance

This module can be used to evaluate feature selection methods via K-fold cross validation.

**class** `picturedrocks.performance.FoldTester(adata)`

Performs K-fold Cross Validation for Marker Selection

`FoldTester` can be used to evaluate various marker selection algorithms. It can split the data in  $K$  folds, run marker selection algorithms on these folds, and classify data based on testing and training data.

**Parameters** `adata (anndata.Anndata)` – data to slice into folds

**classify (classifier)**

Classify each cell using training data from other folds

For each fold, we project the data onto the markers selected for that fold, which we treat as test data. We also project the complement of the fold and treat that as training data.

**Parameters** `classifier` – a classifier that trains with a training data set and predicts labels of test data. See `NearestCentroidClassifier` for an example.

**loadfolds (file)**

Load folds from a file

The file can be one saved either by `FoldTester.savefolds()` or `FoldTester.savefoldsandmarkers()`. In the latter case, it will not load any markers.

**See also:**

`FoldTester.loadfoldsandmarkers()`

**loadfoldsandmarkers (file)**

Load folds and markers

Loads a folds and markers file saved by `FoldTester.savefoldsandmarkers()`

**Parameters** `file (str)` – filename to load from (typically with a .npz extension)

**See also:**

`FoldTester.loadfolds()`

**makefolds ( $k=5$ , random=False)**

Makes folds

**Parameters**

- `k (int)` – the value of  $K$
- `random (bool)` – If true, `makefolds` will make folds randomly. Otherwise, the folds are made in order (i.e., the first  $\text{ceil}(N / k)$  cells in the first fold, etc.)

**savefolds (file)**

Save folds to a file

**Parameters** `file (str)` – filename to save (typically with a .npz extension)

**savefoldsandmarkers (file)**

Save folds and markers for each fold

This saves folds, and for each fold, the markers previously found by `FoldTester.selectmarkers()`.

**Parameters** `file (str)` – filename to save to (typically with a .npz extension)

**selectmarkers (select\_function)**

Perform a marker selection algorithm on each fold

**Parameters** `select_function (function)` – a function that takes in an `AnnData` object and outputs a list of gene markers, given by their index

**validatefolds()**

Ensure that all observations are in exactly one fold

**Returns****Return type** `bool`**class** `picturedrocks.performance.NearestCentroidClassifier`

Nearest Centroid Classifier for Cross Validation

Computes the centroid of each cluster label in the training data, then predicts the label of each test data point by finding the nearest centroid.

**test** (`Xtest`)**train** (`adata`)**class** `picturedrocks.performance.PerformanceReport` (`y, yhat`)

Report actual vs predicted statistics

**Parameters**

- **y** (`numpy.ndarray`) – actual cluster labels, (N, 1)-shaped numpy array
- **yhat** (`numpy.ndarray`) – predicted cluster labels, (N, 1)-shaped numpy array

**confusionmatrixfigure()**

Compute and make a confusion matrix figure

**Returns** confusion matrix**Return type** `plotly.figure`**getconfusionmatrix()**

Get the confusion matrix for the latest run

**Returns** array of shape (K, K), with the [i, j] entry being the fraction of cells in cluster i that were predicted to be in cluster j

**Return type** `numpy.ndarray`**printscore()**

Print a message with the score

**show()**

Print a full report

This uses `iplot`, so we assume this will only be run in a Jupyter notebook and that `init_notebook_mode` has already been run.

**wrong()**

Returns the number of cells misclassified.

**picturedrocks.performance.kfoldindices** (`n, k, random=False`)

Generate indices for k-fold cross validation

**Parameters**

- **n** (`int`) – number of observations
- **k** (`int`) – number of folds
- **random** (`bool`) – determines whether to randomize the order

**Yields** `numpy.ndarray` – array of indices in each fold



# CHAPTER 2

---

## Indices and tables

---

- genindex
- modindex
- search



---

## Python Module Index

---

### p

`picturedrocks.performance`, 7  
`picturedrocks.plot`, 4  
`picturedrocks.preprocessing`, 4  
`picturedrocks.read`, 3



---

## Index

---

### C

CIFE (class in picture-  
drocks.markers.mutualinformation.iterative),  
6  
cife\_obj() (in module picturedrocks.markers.interactive),  
7  
CIFEUnsup (class in picture-  
drocks.markers.mutualinformation.iterative),  
6  
classify() (picturedrocks.performance.FoldTester  
method), 8  
confusionmatrixfigure() (picture-  
drocks.performance.PerformanceReport  
method), 9

### F

FoldTester (class in picturedrocks.performance), 7

### G

genericplot() (in module picturedrocks.plot), 4  
genericwrongplot() (in module picturedrocks.plot), 5  
getconfusionmatrix() (picture-  
drocks.performance.PerformanceReport  
method), 9

### I

InformationSet (class in picturedrocks.markers), 6  
InteractiveMarkerSelection (class in picture-  
drocks.markers.interactive), 7

### J

JMI (class in picture-  
drocks.markers.mutualinformation.iterative),  
6

### K

kfoldindices() (in module picturedrocks.performance), 9

### L

loadfolds() (picturedrocks.performance.FoldTester  
method), 8  
loadfoldsandmarkers() (picture-  
drocks.performance.FoldTester  
method), 8

### M

makefolds() (picturedrocks.performance.FoldTester  
method), 8  
makeinfoset() (in module picturedrocks.markers), 6  
MIM (class in picture-  
drocks.markers.mutualinformation.iterative),  
6

### N

NearestCentroidClassifier (class in picture-  
drocks.performance), 9

### P

pca() (in module picturedrocks.preprocessing), 4  
pcafigure() (in module picturedrocks.plot), 5  
pcawrongplot() (in module picturedrocks.plot), 5  
PerformanceReport (class in picturedrocks.performance),  
9  
picturedrocks.performance (module), 7  
picturedrocks.plot (module), 4  
picturedrocks.preprocessing (module), 4  
picturedrocks.read (module), 3  
plotgeneheat() (in module picturedrocks.plot), 5  
printscore() (picturedrocks.performance.PerformanceReport  
method), 9  
process\_clusts() (in module picturedrocks.read), 3

### R

read\_clusts() (in module picturedrocks.read), 3

### S

savefolds() (picturedrocks.performance.FoldTester  
method), 8

savefoldsandmarkers() (picture-  
drocks.performance.FoldTester  
method),  
8

selectmarkers() (picturedrocks.performance.FoldTester  
method), 8

show() (picturedrocks.performance.PerformanceReport  
method), 9

## T

test() (picturedrocks.performance.NearestCentroidClassifier  
method), 9

train() (picturedrocks.performance.NearestCentroidClassifier  
method), 9

## U

UniEntropy (class in picture-  
drocks.markers.mutualinformation.iterative),  
6

## V

validatefolds() (picturedrocks.performance.FoldTester  
method), 8

## W

wrong() (picturedrocks.performance.PerformanceReport  
method), 9