
PicturedRocks Documentation

Release 1.0.1+0.g9706071.dirty

Umang Varma, Anna Gilbert

Jan 27, 2020

Contents:

1	Installing	3
1.1	Feature Selection Tutorial	3
1.2	Reading data	7
1.3	Preprocessing	8
1.4	Plotting	8
1.5	Selecting Markers	9
1.6	Interactive Marker Selection	12
1.7	Measuring Feature Selection Performance	16
2	Indices and tables	19
	Python Module Index	21
	Index	23

PicturedRocks is a python package that implements information-theoretic feature selection algorithms for scRNA-seq analysis.

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. This allows you to update picturedrocks by pulling the latest commit from github.

1.1 Feature Selection Tutorial

In this Jupyter notebook, we'll walk through the information-theoretic feature selection algorithms in PicturedRocks and demonstrate the interactive marker selection user interface.

If you are viewing this notebook inside the PicturedRocks documentation, the interactive marker selection tool will not work (as it needs a python backend to perform the computations. You can download this notebook from [GitHub](#) and run in on your own computer to try the interactive tool.)

```
[1]: import numpy as np
import scanpy as sc
import picturedrocks as pr
```

```
[2]: adata = sc.datasets.paul15()
```

```
WARNING: In Scanpy 0.*, this returned logarithmized data. Now it returns non-
↳ logarithmized data.
```

```
... storing 'paul15_clusters' as categorical
Trying to set attribute `.uns` of view, making a copy.
```

```
[3]: adata
[3]: AnnData object with n_obs × n_vars = 2730 × 3451
      obs: 'paul15_clusters'
      uns: 'iroot'
```

The `process_clusts` method copies the cluster column and precomputes various indices, etc. If you have multiple columns that can be used as target labels (e.g., different treatments, clusters via different clustering algorithms or parameters, or demographics), this sets and processes the given columns as the one we're currently examining.

This is necessary for supervised analysis and visualization tools in PicturedRocks that use cluster labels.

```
[4]: pr.read.process_clusts(adata, "paul15_clusters")
[4]: AnnData object with n_obs × n_vars = 2730 × 3451
      obs: 'paul15_clusters', 'clust', 'y'
      uns: 'iroot', 'num_clusts', 'clusterindices'
```

The `make_infoset` method creates a `SparseInformationSet` object with a discretized version of the data matrix. It is useful to have only a small number of discrete states that each gene can take so that entropy is a reasonable measurement. By default, `make_infoset` performs an adaptive transform that we call a recursive quantile transform. This is implemented in `pr.markers.mutualinformation.infoset.quantile_discretize`. If you have a different discretization transformation, you can pass a transformed matrix directly to `SparseInformationSet`.

```
[5]: infoset = pr.markers.makeinfoset(adata, True)
```

Because this dataset only has 3451 features, it is computationally easy to do feature selection without restricting the number of features. If we wanted to, we could do either supervised or unsupervised univariate feature selection (i.e., without considering any interactions between features).

```
[6]: # supervised
      mim = pr.markers.mutualinformation.iterative.MIM(infoset)
      most_relevant_genes = mim.autoselect(1000)
```

```
[7]: # unsupervised
      ue = pr.markers.mutualinformation.iterative.UniEntropy(infoset)
      most_variable_genes = ue.autoselect(1000)
```

At this stage we can slice our `adata` object as `adata[:,most_relevant_genes]` or `adata[:,most_variable_genes]` and create a new `InformationSet` object for this sliced object. We don't need to do that here since there are not a lot of genes but will do so anyway for demonstration purposes.

1.1.1 Supervised Feature Selection

Let's jump straight into supervised feature selection. Here we will use the CIFE objective

```
[8]: adata_mr = adata[:,most_relevant_genes].copy()
      infoset_mr = pr.markers.makeinfoset(adata_mr, True)
[9]: cife = pr.markers.CIFE(infoset_mr)
[10]: cife.score[:20]
```



```
[10]: array([0.95022366, 0.93749845, 0.88470651, 0.86819372, 0.8634894 ,
          0.80903075, 0.75775072, 0.75361203, 0.71991963, 0.7106652 ,
          0.70321104, 0.6821289 , 0.67109598, 0.65202536, 0.65192364,
          0.6458561 , 0.64569101, 0.63526239, 0.62452935, 0.62346646])
```

```
[11]: top_genes = np.argsort(cife.score)[::-1]
      print(adata_mr.var_names[top_genes[:10]])

Index(['Mpo', 'Prtn3', 'Ctsg', 'Car2', 'Elane', 'Car1', 'Klf1', 'Blvrb',
       'Ermap', 'Mt2'],
      dtype='object')
```

Let's select 'Mpo'

```
[12]: ind = adata_mr.var_names.get_loc('Mpo')
```

```
[13]: cife.add(ind)
```

Now, the top genes are

```
[14]: top_genes = np.argsort(cife.score)[::-1]
      print(adata_mr.var_names[top_genes[:10]])

Index(['Car2', 'Car1', 'Gnb2l1', 'Fth1', 'Atpif1', 'AK158095', 'Ncl', 'Blvrb',
       'Rpl4', 'Atp5b'],
      dtype='object')
```

Observe that the order has changed based on redundancy (or lack thereof) with 'Mpo'. Let's add 'Car1'

```
[15]: ind = adata_mr.var_names.get_loc('Car1')
      cife.add(ind)
```

```
[16]: top_genes = np.argsort(cife.score)[::-1]
      print(adata_mr.var_names[top_genes[:10]])

Index(['Actb', 'Gpx1', 'Hsp90ab1', 'Ftl1', 'Ybx1', 'AK158095', 'Ncl', 'Rps3',
       'hnRNP A2/B1', 'Tubalb'],
      dtype='object')
```

If we want to select the top gene repeatedly, we can use `autoselect`

```
[17]: cife.autoselect(5)
```

To look at the markers we've selected, we can examine `cife.S`

```
[18]: cife.S
```

```
[18]: [0, 5, 187, 23, 49, 931, 306]
```

```
[19]: adata_mr.var_names[cife.S]
```

```
[19]: Index(['Mpo', 'Car1', 'Actb', 'H2afy', 'Hsp90ab1', 'Gpr56', 'Ly6e'], dtype='object')
```

User Interface

This process can also be done manually with a user-interface allowing you to incorporate domain knowledge in this process. Use the `View` dropdown to look at heatmaps for candidate genes and already selected genes.

Normalize per cell and log transform the data. We are doing this here only to generate familiar features. We do not recommend performing these transformations before `make_infoaset`.

```
[20]: sc.pp.normalize_per_cell(adata_mr)
      sc.pp.log1p(adata_mr)
```

```
[21]: im = pr.markers.InteractiveMarkerSelection(adata_mr, cife, ['tsne', 'violin'])
```

```
Running tsne on cells...
```

```
[22]: im.show()
```

```
Output()
```

Note, that because we passed the same `cife` object, any genes added/removed in the interface will affect the `cife` object.

```
[23]: adata_mr.var_names[cife.S]
```

```
[23]: Index(['Mpo', 'Car1', 'Actb', 'H2afy', 'Hsp90ab1', 'Gpr56', 'Ly6e'], dtype='object')
```

1.1.2 Unsupervised Feature Selection

This works very similarly. In the example below, we'll autoselect 5 genes and then run the interface. Note that although the previous section would not work without cluster labels, the following code will.

```
[24]: cife_unsup = pr.markers.CIFEUnsup(infoaset)
```

```
[25]: cife_unsup.autoselect(5)
```

If you ran the example above, this will load faster because the t_SNE coordinates for genes and cells have already been computed.

```
[26]: im_unsup = pr.markers.interactive.InteractiveMarkerSelection(adata, cife_unsup, ["tsne
      ↪"])
```

```
Running tsne on cells...
```

```
[27]: im_unsup.show()
```

```
Output()
```

1.1.3 Binary Feature Selection

We can also perform feature selection specifically for individual class labels (e.g., clusters). This is done by changing the `SparseInformationSet`'s `y` array. In the example below, we will target the class label "2Ery". Notice that the features selected by MIM (MIM doesn't consider redundancy) are only those that are informative about "2Ery" in particular.

Binary (i.e., not multiclass) feature selection can be performed with any information-theoretic feature selection algorithm (e.g., CIFE, JMI, MIM).

```
[28]: # since we are changing y anyway, the value of include_y (True in the line below)
↳doesn't matter
infoset2 = pr.markers.makeinfoset(adata, True)
infoset2.set_y((adata.obs['clust'] == '2Ery').astype(int).values)
```

```
[29]: mim2 = pr.markers.mutualinformation.iterative.MIM(infoset2)
```

```
[30]: im2 = pr.markers.interactive.InteractiveMarkerSelection(adata, mim2, ["violin"])
```

```
[31]: im2.show()
```

```
Output()
```

```
[ ]:
```

1.2 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*, *name*='clust', *copy*=False)

Process cluster labels from an obs column

This copies *adata.obs[name]* into *adata.obs["clust"]* and precomputes cluster indices, number of clusters, etc for use by various functions in PicturedRocks.

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*=' ', *name*='clust', *header*=True, *copy*=False)

Read cluster labels from a csv into an obs column

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
- **name** (*str*, *optional*) – destination for label is *adata.obs[name]*
- **header** (*bool*) – determines whether csv has a header line. If false, it is assumed that data begins at the first line of csv
- **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`)

1.3 Preprocessing

The preprocessing module provides basic preprocessing tools. To avoid reinventing the wheel, we won't repeat methods already in `scanpy` 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"]`, and `varm["PCs"]` set

Return type `anndata.AnnData`

1.4 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.pcafigure` (*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 *generiplot*, 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

`picturedrocks.plot.umapfigure` (*adata*, ****scatterkwargs**)

1.5 Selecting Markers

argmax

PicturedRocks can be used to implement numerous information-theoretic feature selection methods. In both the supervised and unsupervised cases, it is computationally intractable to optimize the true objective functions. We want to find a small set of genes S , with $|S| < n$ such that, in the supervised case $I(S; y)$ is maximized and in the unsupervised case $H(S)$ is maximized. (For more information, see our paper)

Because these are computationally intractable, we have implemented the following approximations to $x_i I(S \cup \{x_i\}; y)$ (the ideal supervised objective function)

- The **Mutual Information Maximization (MIM)** is a univariate feature selection method—it does not consider the interaction

$$J_{\text{mim}}(x_i) = I(x_i; y)$$

- The **Joint Mutual Information (JMI)** algorithm proposed by Yang and Moody uses a diminishing penalty on redundancy,

$$J_{\text{jmi}}(x_i) = I(x_i; y) - \frac{1}{|S|} \sum_{x_j \in S} I(x_i; x_j; y)$$

- The **(CIFE)** algorithm proposed by Lin and Tang (and independently by others) does not diminish its redundancy penalty

$$J_{\text{cife}}(x_i) = I(x_i; y) - \sum_{x_j \in S} I(x_i; x_j; y)$$

1.5.1 Mutual information

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 `InformationSet`, but by default, we suggest using `makeinfoset()` to generate such an object. The `makeinfoset()` function uses the recursive quantile transform `quantile_discretize()`.

Iterative Feature Selection

All information-theoretic feature selection methods in PicturedRocks are greedy algorithms. In general, they implement the abstract class `IterativeFeatureSelection` class. See *Supervised Feature Selection* and *Unsupervised Feature Selection* below for specific algorithms.

```
class picturedrocks.markers.IterativeFeatureSelection (infoset)
```

Abstract Class for Iterative Feature Selection

```
add (ind)
```

Select specified feature

Parameters `ind` (*int*) – Index of feature to select

```
autoselect (n_feats)
```

Auto select features

This automatically selects `n_feats` features greedily by selecting the feature with the highest score at each iteration.

Parameters `n_feats` (*int*) – The number of features to select

```
remove (ind)
```

Remove specified feature

Parameters `ind` (*int*) – Index of feature to remove

Supervised Feature Selection

```
class picturedrocks.markers.MIM (infoset)
```

```
class picturedrocks.markers.CIFE (infoset)
```

```
class picturedrocks.markers.JMI (infoset)
```

Unsupervised Feature Selection

```
class picturedrocks.markers.UniEntropy (infoset)
```

```
class picturedrocks.markers.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

```
class picturedrocks.markers.SparseInformationSet (X, y=None)
```

Stores sparse discrete gene expression matrix

Parameters

- **X** (*scipy.sparse.csc_matrix*) – a (num_obs, num_vars) shape matrix with dtype *np.integer*
- **y** (*Union[NoneType, numpy.ndarray]*) – if there is a target label column then this should be the target label, otherwise pass *None*.

```
entropy (cols)
```

Entropy of an ensemble of columns

Parameters **cols** (*numpy.ndarray*) – a 1-d array (of dtype *int64*) with indices of columns to compute entropy over. To include the *y* column, use index *-1* as the first entry of *cols*. All entries in *cols* except for the first entry must be non-negative.

Returns the Shannon entropy of *cols*

Return type *numpy.int64*

```
entropy_wrt (cols)
```

Compute multiple entropies at once

This method computes the entropy of *cols + [i]* iterating over all possible values of *i* and returns an array of entropies (one for each column)

Parameters **cols** (*numpy.ndarray*) – a 1-d array (of dtype *int64*) with indices of columns to compute entropy with respect to. To include the *y* column, use index *-1* as the first entry of *cols*. All entries in *cols* except for the first entry must be non-negative.

Note: To compute the entropy of all columns, you need to pass an empty numpy array of dtype *int64*. A quick way to do so is *np.arange(0)*.

Returns a 1-d array of entropies (where entry *i* corresponds to the entropy of columns *cols* together with column *i*)

Return type *numpy.ndarray*

`picturedrocks.markers.makeinfoset` (*adata*, *include_y*, *k=5*)
Discretize data and make a `SparseInformationSet` object

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 `SparseInformationSet`

`picturedrocks.markers.mutualinformation.infoset.quantile_discretize` (*X*, *k=5*)
Discretize data matrix with a recursive quantile transform

Parameters

- **x** (*Union[numpy.ndarray, scipy.sparse.spmatrix]*) – The input data matrix to transform.
- **k** (*int*) – The number of bins to use in the discretization.

Returns The discretized data matrix

Return type `np.ndarray`

1.6 Interactive Marker Selection

PicturedRocks provides an interactive user interface for selecting markers that allows the practitioner to use both their domain knowledge and the information-theoretic methods to select a small set of genes.

The user can specify the tabs they want available through the `visuals` argument. It is also possible to write your own interface for a tab using any plotting library available to you. See below for more details.

```
class picturedrocks.markers.interactive.InteractiveMarkerSelection (adata,  
                                                                    fea-  
                                                                    ture_  
                                                                    selection,  
                                                                    visu-  
                                                                    als=None,  
                                                                    disp_genes=10,  
                                                                    con-  
                                                                    nected=True)
```

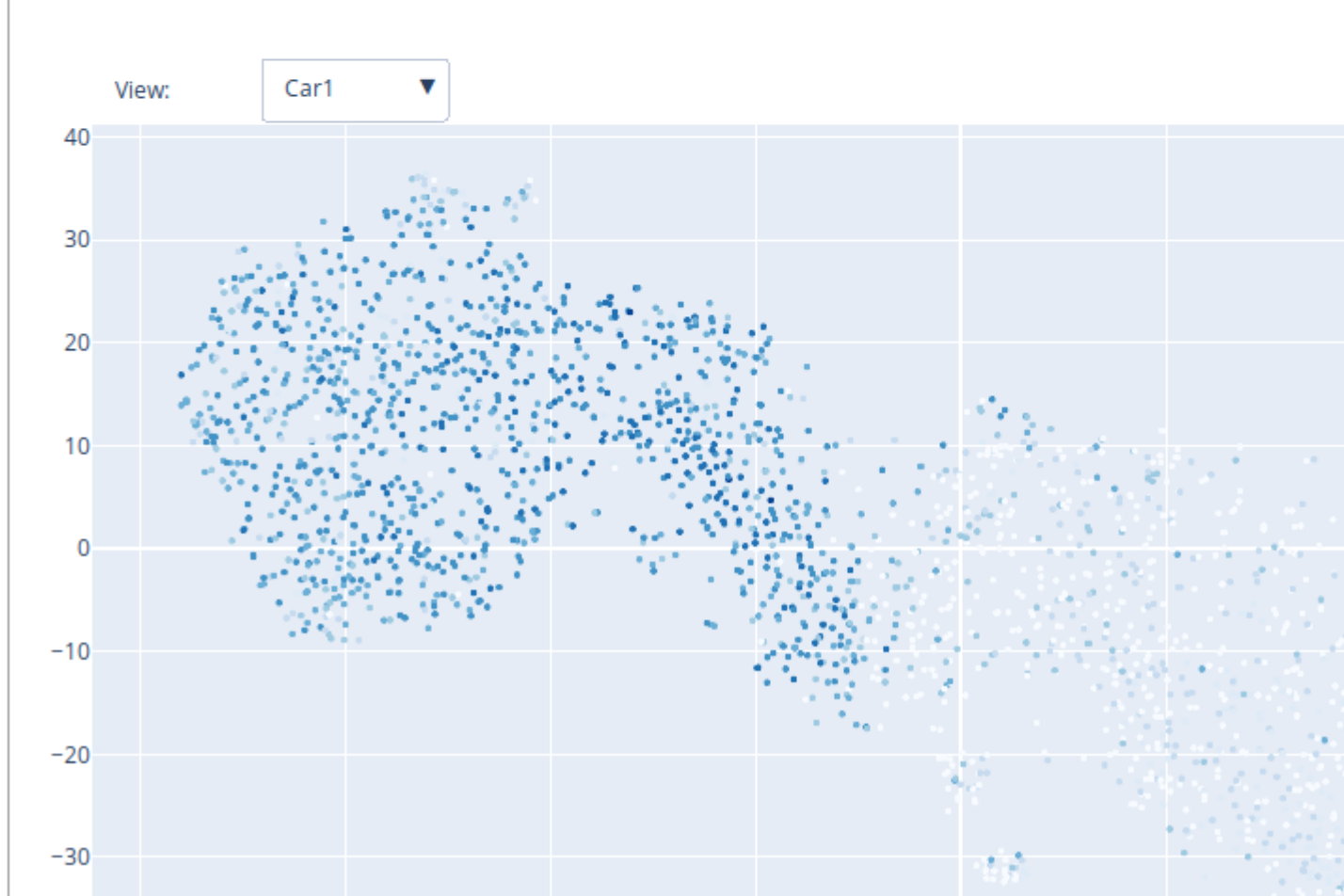
Run an interactive marker selection GUI inside a jupyter notebook

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 interactive feature selection algorithm class that corresponds to `adata` (i.e., the column indices in `feature_selection` should correspond to the column indices in `adata`)
- **visuals** (*list*) – List of visualizations to display. These can either be shorthands for built-in visualizations (currently “tsne”, “umap”, and “violin”), or an instance of `InteractiveVisualization` (see `GeneHeatmap` and `ViolinPlot` for example implementations).

Candidate Next Gene	Currently selected genes
+ Myb (score: 1.0000)	Mpo <input type="checkbox"/>
+ Ifngr1 (score: 0.9457)	Car1 <input type="checkbox"/>
+ Srgn (score: 0.9387)	Actb <input type="checkbox"/>
+ Gata2 (score: 0.9264)	H2afy <input type="checkbox"/>
+ Tmsb10 (score: 0.9244)	Hsp90ab1 <input type="checkbox"/>
+ Tmbim6 (score: 0.9236)	Gpr56 <input type="checkbox"/>
+ Apoe (score: 0.9154)	Ly6e <input type="checkbox"/>
+ Hn1 (score: 0.9116)	
+ Cbfa2t3 (score: 0.9085)	
+ Cyba (score: 0.9039)	
+ <input type="text" value="Gene Name"/> (?)	

t-SNE Heatmap Violin Plots



1.6. Interactive Marker Selection

Fig. 1: A screenshot of the Interactive Marker Selection user interface being used on the Paul dataset. The “candidate” genes are updated after any addition to or removal from the “current” list of genes. The user may also type in the name of a gene at any point.

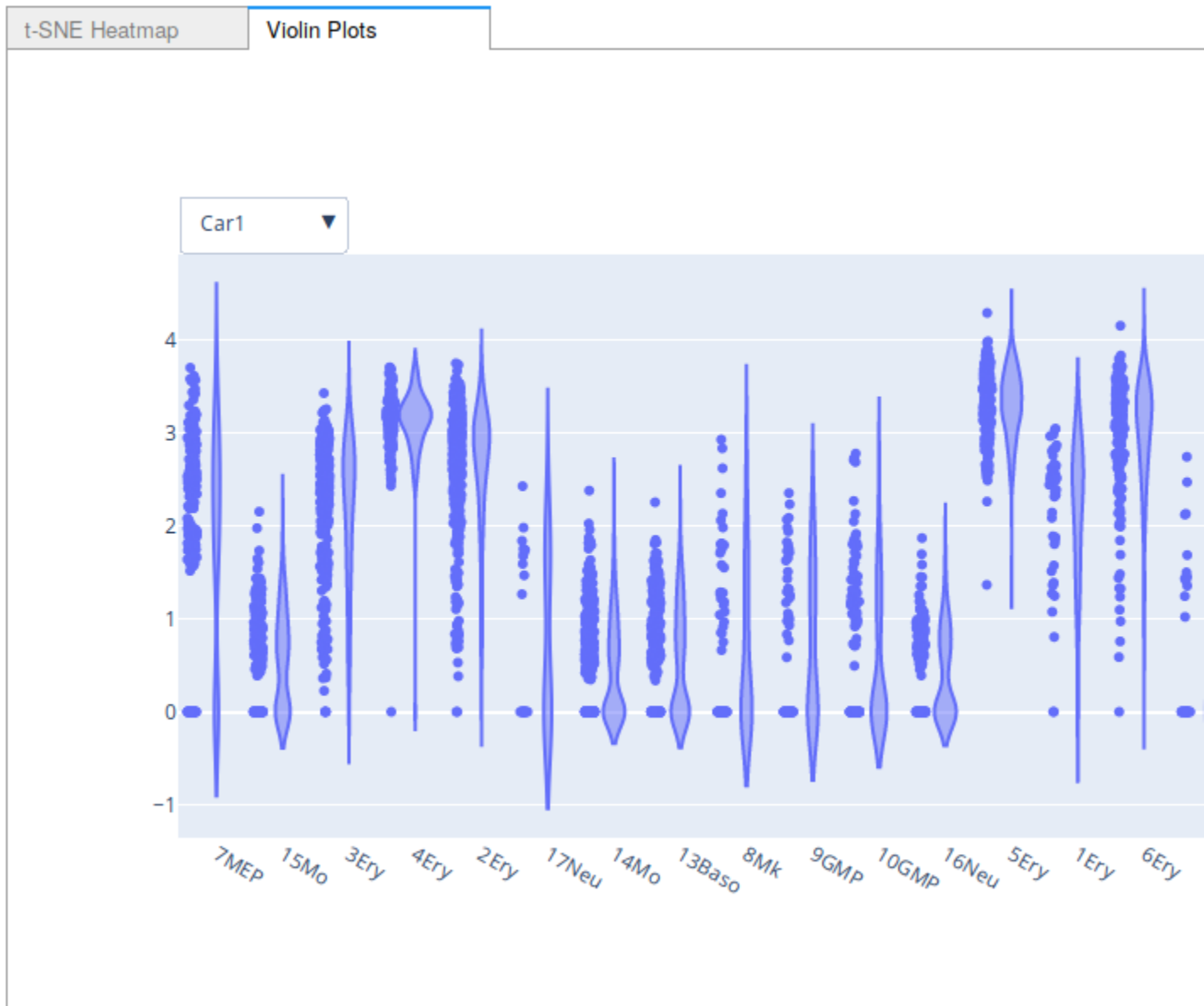


Fig. 2: A screenshot of the violin plot option for Interactive Marker Selection.

- **disp_genes** (*int*) – Number of genes to display as options (by default, number of genes plotted on the tSNE plot is $3 * disp_genes$, 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`.

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.

redraw ()

Redraw jupyter widgets after a change

This is called internally and there should usually be no need for the user to explicitly call this method.

show ()

Display the jupyter widgets

class picturedrocks.markers.interactive.**GeneHeatmap** (*dim_red='tsne', n_pcs=30*)
GeneHeatmap for Interactive Marker Selection

Parameters

- **dim_red** (*str*) – Dimensionality reduction algorithm to use. Currently available options are “umap” and “tsne”
- **n_pcs** (*int*) – The number of principal components to map to before running dimensionality reduction

class picturedrocks.markers.interactive.**ViolinPlot**
Violin Plots for each class label

1.6.1 Using Your Own Visualizations

To write your own visualization, simply extend the *InteractiveVisualization*. Use the source code for *GeneHeatmap* and *ViolinPlot* as a reference for implementation.

class picturedrocks.markers.interactive.**InteractiveVisualization**

Abstract base class for interactive visualizations

Extend this class and pass an instance of it to *InteractiveMarkerSelection* to use your own visualization. It is recommended that you begin your implementation of `__init__` with:

```
super().__init__()
```

You are welcome to add parameters specific to your visualization in the `__init__` method.

prepare (*adata, out*)

Prepare for visualization

This method is called when *InteractiveMarkerSelection* is initialized. It is recommended that you begin your implementation with

```
super().prepare(adata, out)
```

This stores *adata* and *out* in `self.adata` and `self.out` respectively.

redraw (*next_gene_inds*, *cur_gene_inds*)

Draw the visualization

You must implement this method. To display the plots in the appropriate widget, use:

```
with self.out:
    fig.show() # or your plotting library's equivalent
```

title

Title of the visualization

This should be a Python property, using the `@property` decorator.

1.7 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.

Note: The *classifier* should not attempt to modify data in-place. Any preprocessing should be done on a copy.

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

Note: The *select_function* should not attempt to modify data in-place. Any preprocessing should be done on a copy.

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 *ipplot*, 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

`picturedrocks.performance.merge_markers (ft, n_markers)`

`picturedrocks.performance.truncatemarkers (ft, n_markers)`

CHAPTER 2

Indices and tables

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

p

`picturedrocks.markers`, 9
`picturedrocks.markers.interactive`, 12
`picturedrocks.performance`, 16
`picturedrocks.plot`, 8
`picturedrocks.preprocessing`, 8
`picturedrocks.read`, 7

A

`add()` (*picturedrocks.markers.IterativeFeatureSelection* method), 10

`autoselect()` (*picturedrocks.markers.IterativeFeatureSelection* method), 10

C

`CIFE` (*class in picturedrocks.markers*), 10

`CIFEUnsup` (*class in picturedrocks.markers*), 11

`classify()` (*picturedrocks.performance.FoldTester* method), 16

`confusionmatrixfigure()` (*picturedrocks.performance.PerformanceReport* method), 17

E

`entropy()` (*picturedrocks.markers.SparseInformationSet* method), 11

`entropy_wrt()` (*picturedrocks.markers.SparseInformationSet* method), 11

F

`FoldTester` (*class in picturedrocks.performance*), 16

G

`GeneHeatmap` (*class in picturedrocks.markers.interactive*), 15

`genericplot()` (*in module picturedrocks.plot*), 8

`genericwrongplot()` (*in module picturedrocks.plot*), 8

`getconfusionmatrix()` (*picturedrocks.performance.PerformanceReport* method), 17

I

`InformationSet` (*class in picturedrocks.markers*), 11

`InteractiveMarkerSelection` (*class in picturedrocks.markers.interactive*), 12

`InteractiveVisualization` (*class in picturedrocks.markers.interactive*), 15

`IterativeFeatureSelection` (*class in picturedrocks.markers*), 10

J

`JMI` (*class in picturedrocks.markers*), 10

K

`kfoldindices()` (*in module picturedrocks.performance*), 18

L

`loadfolds()` (*picturedrocks.performance.FoldTester* method), 16

`loadfoldsandmarkers()` (*picturedrocks.performance.FoldTester* method), 16

M

`makefolds()` (*picturedrocks.performance.FoldTester* method), 16

`makeinfoset()` (*in module picturedrocks.markers*), 11

`merge_markers()` (*in module picturedrocks.performance*), 18

`MIM` (*class in picturedrocks.markers*), 10

N

`NearestCentroidClassifier` (*class in picturedrocks.performance*), 17

P

`pca()` (*in module picturedrocks.preprocessing*), 8

`pcafigure()` (*in module picturedrocks.plot*), 9

`pcawrongplot()` (*in module picturedrocks.plot*), 9

PerformanceReport (class in picture-
drocks.performance), 17

picturedrocks.markers (module), 9

picturedrocks.markers.interactive (mod-
ule), 12

picturedrocks.performance (module), 16

picturedrocks.plot (module), 8

picturedrocks.preprocessing (module), 8

picturedrocks.read (module), 7

plotgeneheat () (in module picturedrocks.plot), 9

prepare () (picturedrocks.markers.interactive.Interactive
method), 15

printscore () (picture-
drocks.performance.PerformanceReport
method), 18

process_clusts () (in module picturedrocks.read), 7

Q

quantile_discretize () (in module picture-
drocks.markers.mutualinformation.infoset),
12

R

read_clusts () (in module picturedrocks.read), 7

redraw () (picturedrocks.markers.interactive.InteractiveMarkerSelection
method), 15

redraw () (picturedrocks.markers.interactive.InteractiveVisualization
method), 15

remove () (picturedrocks.markers.IterativeFeatureSelection
method), 10

S

savefolds () (picturedrocks.performance.FoldTester
method), 17

savefoldsandmarkers () (picture-
drocks.performance.FoldTester method),
17

selectmarkers () (picture-
drocks.performance.FoldTester method),
17

show () (picturedrocks.markers.interactive.InteractiveMarkerSelection
method), 15

show () (picturedrocks.performance.PerformanceReport
method), 18

SparseInformationSet (class in picture-
drocks.markers), 11

T

test () (picturedrocks.performance.NearestCentroidClassifier
method), 17

title (picturedrocks.markers.interactive.InteractiveVisualization
attribute), 16

train () (picturedrocks.performance.NearestCentroidClassifier
method), 17

truncatemarkers () (in module picture-
drocks.performance), 18

U

umapfigure () (in module picturedrocks.plot), 9

UniEntropy (class in picturedrocks.markers), 11

V

validatefolds () (picture-
drocks.performance.FoldTester method),
17

ViolinPlot (class in picture-
drocks.markers.interactive), 15

W

wrong () (picturedrocks.performance.PerformanceReport
method), 18