

---

# **Hypercluster**

***Release 0.0.2***

**Ruggleslab**

**Feb 25, 2020**



# CONTENTS

<b>1</b>	<b>hypercluster package</b>	<b>1</b>
1.1	hypercluster.classes module . . . . .	5
1.2	hypercluster.utilities module . . . . .	8
1.3	hypercluster.visualize module . . . . .	10
1.4	hypercluster.constants module . . . . .	12
1.5	hypercluster.additional_clusterers module . . . . .	13
1.6	hypercluster.additional_metrics module . . . . .	15
<b>2</b>	<b>hypercluster SnakeMake pipeline</b>	<b>17</b>
2.1	Line-by-line explanation of config.yml . . . . .	18
2.2	config.yml example from scRNA-seq workflow . . . . .	19
<b>3</b>	<b>Indices and tables</b>	<b>21</b>
<b>4</b>	<b>Installation and logistics</b>	<b>23</b>
4.1	Installation . . . . .	23
4.2	Quick reference for clustering and evaluation . . . . .	24
4.3	Quickstart and examples . . . . .	24
	<b>Python Module Index</b>	<b>27</b>
	<b>Index</b>	<b>29</b>



---

CHAPTER  
ONE

---

## HYPERCLUSTER PACKAGE

```
class hypercluster.AutoClusterer(clusterer_name:          Optional[str]      =      'KMeans',
                                 params_to_optimize:    Optional[dict]    =  None,   ran-
                                 dom_search:           bool        =  False,   random_search_fraction:
                                 Optional[float]        =  0.5,   param_weights: dict     =  {}, 
                                 clus_kwargs:          Optional[dict]  =  None,   labels_:    Optional[pandas.core.frame.DataFrame] =  None, 
                                 evaluation_:          Optional[pandas.core.frame.DataFrame] =  None,
                                 data:                 Optional[pandas.core.frame.DataFrame] =  None,
                                 labels_df:            Optional[pandas.core.frame.DataFrame] =  None,
                                 evaluation_df:        Optional[pandas.core.frame.DataFrame] =  None)
Bases: hypercluster.classes.Clusterer
```

Main hypercluster object.

**clusterer\_name**

String name of clusterer.

**Type** str

**params\_to\_optimize**

Dictionary with possibilities for different parameters. Ex format - {‘parameter\_name’:[1, 2, 3, 4, 5]}. If None, will optimize default selection, given in hypercluster.constants.variables\_to\_optimize. Default None.

**Type** dict

**random\_search**

Whether to search a random selection of possible parameters or all possibilities. Default True.

**Type** bool

**random\_search\_fraction**

If random\_search is True, what fraction of the possible parameters to search. Default 0.5.

**Type** float

**param\_weights**

Dictionary of str: dictionaries. Ex format - { ‘parameter\_name’: {‘param\_option\_1’:0.5, ‘param\_option\_2’:0.5} }.

**Type** dict

**clus\_kwargs**

Additional kwargs to pass into given clusterer, but not to be optimized. Default None.

**Type** dict

**labels\_**

If already fit, labels DataFrame fit to data.

**Type** Optional[DataFrame]

**evaluation\_**

If already fit and evalute, evaluations per label.

**Type** Optional[DataFrame]

**data**

Data to fit, will not fit by default even if passed data.

**Type** Optional[DataFrame]

**evaluate** (*methods*: Optional[Iterable[str]] = None, *metric\_kwargs*: Optional[dict] = None,

*gold\_standard*: Optional[Iterable[T\_co]] = None)

Evaluate labels with given metrics.

**Parameters**

- **methods** (Optional[Iterable[str]]) – List of evaluation methods to use.
- **metric\_kwargs** (Optional[dict]) – Additional kwargs per evaluation metric. Structure of {‘metric\_name’:{‘param1’:value, ‘param2’:val2}}.
- **gold\_standard** (Optional[Iterable]) – Gold standard labels, if available. Only needed if using a metric that needs ground truth.

**Returns (AutoClusterer):** self with attribute .evaluation\_; a DataFrame with all eval values per labels.

**fit** (*data*: pandas.core.frame.DataFrame)

Fits clusterer to data with each parameter set.

**Parameters** **data** (DataFrame) – DataFrame with elements to cluster as index and features as columns.

**Returns (AutoClusterer):** self

**generate\_param\_sets** ()

Uses info from init to make a Dataframe of all parameter sets that will be tried.

**Returns (AutoClusterer):** self

```
class hypercluster.MultiAutoClusterer(algorithm_names: Union[Iterable[T_co], str, None]
= None, algorithm_parameters: Optional[Dict[str, dict]] = None, random_search: bool = False,
random_search_fraction: Optional[float] = 0.5, algorithm_param_weights: Optional[dict] = None,
algorithm_clus_kwargs: Optional[dict] = None, data: Optional[pandas.core.frame.DataFrame] = None,
evaluation_methods: Optional[List[str]] = None, metric_kwargs: Optional[Dict[str, dict]] = None,
gold_standard: Optional[Iterable[T_co]] = None, autoclusterers: Iterable[hypercluster.classes.AutoClusterer] = None,
labels_: Dict[str, hypercluster.classes.AutoClusterer] = None, evaluation_: Dict[str, hypercluster.classes.AutoClusterer] = None,
labels_df: Optional[pandas.core.frame.DataFrame] = None, evaluation_df: Optional[pandas.core.frame.DataFrame] = None)
```

Bases: `hypercluster.classes.Clusterer`

Object for training multiple clustering algorithms.

#### **algorithm\_names**

List of algorithm names to test OR name of category of clusterers from `hypercluster.constants.categories`, OR None. If None, default is `hypercluster.constants.variables_to_optimize.keys()`.

**Type** `Optional[Union[Iterable, str]]`

#### **algorithm\_parameters**

Dictionary of hyperparameters to optimize. Example format: {‘clusterer\_name1’:{‘hyperparam1’:[val1, val2]} }.

**Type** `Optional[Dict[str, dict]]`

#### **random\_search**

Whether to search a random subsample of possible conditions.

**Type** `bool`

#### **random\_search\_fraction**

If `random_search`, what fraction of conditions to search.

**Type** `float`

#### **algorithm\_param\_weights**

If `random_search`, and you want to give probability weights to certain parameters, dictionary of probability weights. Example format: {‘clusterer1’: {‘hyperparam1’:{val1:probability1, val2:probability2}} }.

**Type** `Dict[str, Dict[str, dict]]`

#### **algorithm\_clus\_kwargs**

Dictionary of additional keyword args for any clusterer. Example format: {‘clusterer1’: {‘param1’:val1}}.

**Type** `Dict[str, dict]`

#### **data**

Optional, data to fit. Will not fit even if passed, need to call fit method.

**Type** `Optional[DataFrame]`

**evaluation\_methods**

List of metrics with which to evaluate. If None, will use hypercluster.constants.inherent\_metrics. Default is None.

**Type** Optional[List[str]]

**metric\_kwargs**

Additional keyword args for any metric function. Example format: {‘metric1’: {‘param1’: value}}.

**Type** Optional[Dict[str, dict]]

**gold\_standard**

If using methods that need ground truth, vector of correct labels. Can also pass in during evaluate.

**Type** Optional[Iterable]

**autoclusterers**

If building from initialized AutoClusterer objects, can give a list of them here. If these are given, it will override anything

**Type** Iterable[AutoClusterer]

**passed to labels\_ and evaluation\_**

**labels\_**

Dictionary of label DataFrames per clusterer, if already fit. Example format: {‘clusterer1’: labels\_df}.

**Type** Optional[Dict[str, DataFrame]]

**evaluation\_**

Dictionary of evaluation DataFrames per clusterer, if already fit and evaluated. Example format: {‘clusterer1’: evaluation\_df}.

**Type** Optional[Dict[str, DataFrame]]

**labels\_df**

Combined DataFrame of all labeling results.

**Type** Optional[DataFrame]

**evaluation\_df**

Combined DataFrame of all evaluation results.

**Type** Optional[DataFrame]

**evaluate** (evaluation\_methods: Optional[list] = None, metric\_kwargs: Optional[dict] = None, gold\_standard: Optional[Iterable[T\_co]] = None)

**fit** (data: Optional[pandas.core.frame.DataFrame] = None)

## 1.1 hypercluster.classes module

```
class hypercluster.classes.AutoClusterer(clusterer_name: Optional[str] = 'KMeans',
                                         params_to_optimize: Optional[dict] = None,
                                         random_search: bool = False,
                                         random_search_fraction: Optional[float] = 0.5,
                                         param_weights: dict = {}, clus_kwarg: Optional[dict] = None,
                                         labels_: Optional[pandas.core.frame.DataFrame] = None,
                                         evaluation_: Optional[pandas.core.frame.DataFrame] = None,
                                         data: Optional[pandas.core.frame.DataFrame] = None,
                                         labels_df: Optional[pandas.core.frame.DataFrame] = None,
                                         evaluation_df: Optional[pandas.core.frame.DataFrame] = None)
```

Bases: `hypercluster.classes.Clusterer`

Main hypercluster object.

**clusterer\_name**

String name of clusterer.

**Type** str

**params\_to\_optimize**

Dictionary with possibilities for different parameters. Ex format - {‘parameter\_name’:[1, 2, 3, 4, 5]}. If None, will optimize default selection, given in `hypercluster.constants.variables_to_optimize`. Default None.

**Type** dict

**random\_search**

Whether to search a random selection of possible parameters or all possibilities. Default True.

**Type** bool

**random\_search\_fraction**

If random\_search is True, what fraction of the possible parameters to search. Default 0.5.

**Type** float

**param\_weights**

Dictionary of str: dictionaries. Ex format - { ‘parameter\_name’: {‘param\_option\_1’:0.5, ‘param\_option\_2’:0.5} }.

**Type** dict

**clus\_kwarg**

Additional kwarg to pass into given clusterer, but not to be optimized. Default None.

**Type** dict

**labels\_**

If already fit, labels DataFrame fit to data.

**Type** `Optional[DataFrame]`

**evaluation\_**

If already fit and evaluate, evaluations per label.

**Type** `Optional[DataFrame]`

**data**

Data to fit, will not fit by default even if passed data.

**Type** Optional[DataFrame]

**evaluate** (*methods*: Optional[Iterable[str]] = None, *metric\_kwargs*: Optional[dict] = None, *gold\_standard*: Optional[Iterable[T\_co]] = None)

Evaluate labels with given metrics.

**Parameters**

- **methods** (Optional[Iterable[str]]) – List of evaluation methods to use.
- **metric\_kwargs** (Optional[dict]) – Additional kwargs per evaluation metric. Structure of {'metric\_name': {'param1': value, 'param2': val2}}.
- **gold\_standard** (Optional[Iterable]) – Gold standard labels, if available. Only needed if using a metric that needs ground truth.

**Returns (AutoClusterer):** self with attribute .evaluation\_; a DataFrame with all eval values per labels.

**fit** (*data*: pandas.core.frame.DataFrame)

Fits clusterer to data with each parameter set.

**Parameters** **data** (DataFrame) – DataFrame with elements to cluster as index and features as columns.

**Returns (AutoClusterer):** self

**generate\_param\_sets()**

Uses info from init to make a Dataframe of all parameter sets that will be tried.

**Returns (AutoClusterer):** self

**class** hypercluster.classes.Clusterer  
Bases: object

Meta class for shared methods for both AutoClusterer and MultiAutoClusterer.

**fit\_predict** (*data*: Optional[pandas.core.frame.DataFrame], *parameter\_set\_name*, *method*, *min\_of\_max*)

**pick\_best\_labels** (*method*: Optional[str] = None, *min\_or\_max*: Optional[str] = None)

**visualize\_evaluations** (*savefig*: bool = False, *output\_prefix*: str = 'evaluations', \*\*heatmap\_kws) → List[matplotlib.axes.\_axes.Axes]

**visualize\_for\_picking\_labels** (*method*: Optional[str] = None, *savefig\_prefix*: Optional[str] = None)

**visualize\_label\_agreement** (*method*: Optional[str] = None, *savefig*: bool = False, *output\_prefix*: Optional[str] = None, \*\*heatmap\_kws) → List[matplotlib.axes.\_axes.Axes]

**visualize\_sample\_label\_consistency** (*savefig*: bool = False, *output\_prefix*: Optional[str] = None, \*\*heatmap\_kws) → List[matplotlib.axes.\_axes.Axes]

```
class hypercluster.classes.MultiAutoClusterer(algorithm_names: Union[Iterable[T_co],
    str, None] = None, algorithm_parameters:
    Optional[Dict[str, dict]] = None, random_search: bool =
    False, random_search_fraction:
    Optional[float] = 0.5, algorithm_param_weights:
    Optional[dict] = None, algorithm_clus_kwargs:
    Optional[dict] = None, data: Optional[pandas.core.frame.DataFrame]
    = None, evaluation_methods: Optional[List[str]] = None, metric_kwargs:
    Optional[Dict[str, dict]] = None, gold_standard: Optional[Iterable[T_co]]
    = None, autoclusterers: Iterable[hypercluster.classes.AutoClusterer]
    = None, labels_: Dict[str, hypercluster.classes.AutoClusterer] =
    None, evaluation_: Dict[str, hypercluster.classes.AutoClusterer]
    = None, labels_df: Optional[pandas.core.frame.DataFrame]
    = None, evaluation_df: Optional[pandas.core.frame.DataFrame] =
    None)
```

Bases: `hypercluster.classes.Clusterer`

Object for training multiple clustering algorithms.

#### `algorithm_names`

List of algorithm names to test OR name of category of clusterers from `hypercluster.constants.categories`, OR None. If None, default is `hypercluster.constants.variables_to_optimize.keys()`.

**Type** `Optional[Union[Iterable, str]]`

#### `algorithm_parameters`

Dictionary of hyperparameters to optimize. Example format: `{‘clusterer_name1’: {‘hyperparam1’: [val1, val2]} }`.

**Type** `Optional[Dict[str, dict]]`

#### `random_search`

Whether to search a random subsample of possible conditions.

**Type** `bool`

#### `random_search_fraction`

If `random_search`, what fraction of conditions to search.

**Type** `float`

#### `algorithm_param_weights`

If `random_search`, and you want to give probability weights to certain parameters, dictionary of probability weights. Example format: `{‘clusterer1’: {‘hyperparam1’: {val1:probability1, val2:probability2}}}`.

**Type** `Dict[str, Dict[str, dict]]`

#### `algorithm_clus_kwargs`

Dictionary of additional keyword args for any clusterer. Example format: `{‘clusterer1’: {‘param1’:val1}}`.

**Type** `Dict[str, dict]`

```
data
    Optional, data to fit. Will not fit even if passed, need to call fit method.

    Type Optional[DataFrame]

evaluation_methods
    List of metrics with which to evaluate. If None, will use hypercluster.constants.inherent_metrics. Default is None.

    Type Optional[List[str]]

metric_kwargs
    Additional keyword args for any metric function. Example format: {'metric1': {'param1': value}}.

    Type Optional[Dict[str, dict]]

gold_standard
    If using methods that need ground truth, vector of correct labels. Can also pass in during evaluate.

    Type Optional[Iterable]

autoclusterers
    If building from initialized AutoClusterer objects, can give a list of them here. If these are given, it will override anything.

    Type Iterable[AutoClusterer]

passed to labels_ and evaluation_.

labels_
    Dictionary of label DataFrames per clusterer, if already fit. Example format: {'clusterer1': labels_df}.

    Type Optional[Dict[str, DataFrame]]

evaluation_
    Dictionary of evaluation DataFrames per clusterer, if already fit and evaluated. Example format: {'clusterer1': evaluation_df}.

    Type Optional[Dict[str, DataFrame]]

labels_df
    Combined DataFrame of all labeling results.

    Type Optional[DataFrame]

evaluation_df
    Combined DataFrame of all evaluation results.

    Type Optional[DataFrame]

evaluate (evaluation_methods: Optional[list] = None, metric_kwargs: Optional[dict] = None,
           gold_standard: Optional[Iterable[T_co]] = None)

fit (data: Optional[pandas.core.frame.DataFrame] = None)
```

## 1.2 hypercluster.utilities module

```
hypercluster.utilities.calculate_row_weights (row: Iterable[T_co], param_weights: dict,
                                             vars_to_optimize: dict) → float
```

Used to select random rows of parameter combinations using individual parameter weights.

### Parameters

- **row** (*Iterable*) – Series of parameters, with parameter names as index.

- **param\_weights** (*dict*) – Dictionary of str: dictionaries. Ex format - {‘parameter\_name’:{‘param\_option\_1’:0.5, ‘param\_option\_2’:0.5}}.
- **vars\_to\_optimize** (*Iterable*) – Dictionary with possibilities for different parameters. Ex format - {‘parameter\_name’:[1, 2, 3, 4, 5]}.

**Returns (float):** Float representing the probability of seeing that combination of parameters, given their individual weights.

```
hypercluster.utilities.cluster(clusterer_name: str, data: pandas.core.frame.DataFrame,
                               params: dict = {})
```

Runs a given clusterer with a given set of parameters.

#### Parameters

- **clusterer\_name** (*str*) – String name of clusterer.
- **data** (*DataFrame*) – Dataframe with elements to cluster as index and examples as columns.
- **params** (*dict*) – Dictionary of parameter names and values to feed into clusterer. Default {}

**Returns** Instance of the clusterer fit with the data provided.

```
hypercluster.utilities.convert_to_multiind(key: str, df: pandas.core.frame.DataFrame)
                                         → pandas.core.frame.DataFrame
```

Takes columns from a single clusterer from Clusterer.labels\_df or .evaluation\_df and converts to a multiindexed rather than collapsed into string. Equivalent to grabbing Clusterer.labels[clusterer] or .evaluations[clusterer]. Opposite of generate\_flattened\_df.

#### Parameters

- **key** (*str*) – Name of clusterer, must match beginning of columns to convert.
- **df** (*DataFrame*) – Dataframe to grab chunk from.

**Returns** Subset DataFrame with multiindex.

```
hypercluster.utilities.evaluate_one(labels: Iterable[T_co], method: str = 'silhouette_score',
                                    data: Optional[pandas.core.frame.DataFrame] = None,
                                    gold_standard: Optional[Iterable[T_co]] = None, metric_kwargs: Optional[dict] = None) → dict
```

Uses a given metric to evaluate clustering results.

#### Parameters

- **labels** (*Iterable*) – Series of labels.
- **method** (*str*) – Str of name of evaluation to use. Default is silhouette.
- **data** (*DataFrame*) – If using an inherent metric, must provide DataFrame with which to calculate the metric.
- **gold\_standard** (*Iterable*) – If using a metric that compares to ground truth, must provide a set of gold standard labels.
- **metric\_kwargs** (*dict*) – Additional kwargs to use in evaluation.

**Returns (float):** Metric value

```
hypercluster.utilities.generate_flattened_df (df_dict: Dict[str, das.core.frame.DataFrame]) → pandas.core.frame.DataFrame
Takes dictionary of results from many clusterers and makes 1 DataFrame. Opposite of convert_to_multiind.
```

**Parameters** `df_dict` (`Dict[str, DataFrame]`) – Dictionary of dataframes to flatten. Can be `.labels_` or `.evaluations_` from MultiAutoClusterer.

**Returns** Flattened DataFrame with all data.

```
hypercluster.utilities.pick_best_labels (evaluation_results_df: pandas.core.frame.DataFrame, clustering_labels_df: pandas.core.frame.DataFrame, method: Optional[str] = None, min_or_max: Optional[str] = None) → Iterable[T_co]
From evaluations and a metric to minimize or maximize, return all labels with top pick.
```

**Parameters**

- `evaluation_results_df` (`DataFrame`) – Evaluations DataFrame from optimize\_clustering.
- `clustering_labels_df` (`DataFrame`) – Labels DataFrame from optimize\_clustering.
- `method` (`str`) – Method with which to choose the best labels.
- `min_or_max` (`str`) – Whether to minimize or maximize the metric. Must be ‘min’ or ‘max’.

**Returns (DataFrame):** DataFrame of all top labels.

## 1.3 hypercluster.visualize module

```
hypercluster.visualize.compute_order (df, dist_method: str = 'euclidean', cluster_method: str = 'average')
```

Gives hierarchical clustering order for the rows of a DataFrame

**Parameters**

- `df` (`DataFrame`) – DataFrame with rows to order.
- `dist_method` (`str`) – Distance method to pass to `scipy.cluster.hierarchy.linkage`.
- `cluster_method` (`str`) – Clustering method to pass to `scipy.spatial.distance.pdist`.

**Returns (pandas.Index):** Ordered row index.

```
hypercluster.visualize.visualize_evaluations (evaluations_df: pandas.core.frame.DataFrame, savefig: bool = False, output_prefix: str = 'evaluations', **heatmap_kws) → List[matplotlib.axes.Axes]
```

Makes a z-scored visualization of all evaluations.

**Parameters**

- `evaluations_df` (`DataFrame`) – Evaluations dataframe from `clustering.optimize_clustering`
- `output_prefix` (`str`) – If saving a figure, file prefix to use.

- **savefig** (*bool*) – Whether to save a pdf
- **\*\*heatmap\_kws** – Additional keyword arguments to pass to seaborn.heatmap.

**Returns (List[matplotlib.axes.Axes]):** List of all matplotlib axes.

```
hypercluster.visualize.visualize_for_picking_labels (evaluation_df: pandas.core.frame.DataFrame, method: Optional[str] = None, savefig_prefix: Optional[str] = None)
```

Generates graphs similar to a scree graph for PCA for each parameter and each clusterer.

#### Parameters

- **evaluation\_df** (*DataFrame*) – DataFrame of evaluations to visualize. Clusterer.evaluation\_df.
- **method** (*str*) – Which metric to visualize.
- **savefig\_prefix** (*str*) – If not None, save a figure with give prefix.

**Returns** matplotlib axes.

```
hypercluster.visualize.visualize_label_agreement (labels: pandas.core.frame.DataFrame, method: Optional[str] = None, savefig: bool = False, output_prefix: Optional[str] = None, **heatmap_kws) → List[matplotlib.axes._axes.Axes]
```

Visualize similarity between clustering results given an evaluation metric.

#### Parameters

- **labels** (*DataFrame*) – Labels DataFrame, e.g. from optimize\_clustering or **AutoClusterer.labels**
- **method** (*str*) – Method with which to compare labels. Must be a metric like the ones in constants.need\_ground\_truth, which takes two sets of labels.
- **savefig** (*bool*) – Whether to save a pdf.
- **output\_prefix** (*str*) – If saving a pdf, file prefix to use.
- **\*\*heatmap\_kws** – Additional keywords to pass to seaborn.heatmap

**Returns (List[matplotlib.axes.Axes]):** List of matplotlib axes

```
hypercluster.visualize.visualize_pairwise (df: pandas.core.frame.DataFrame, savefig: bool = False, output_prefix: Optional[str] = None, method: Optional[str] = None, **heatmap_kws) → List[matplotlib.axes._axes.Axes]
```

Visualize symmetrical square DataFrames.

#### Parameters

- **df** (*DataFrame*) – DataFrame to visualize.
- **savefig** (*bool*) – Whether to save a pdf.
- **output\_prefix** (*str*) – If saving a pdf, file prefix to use.

- **method** (*str*) – Label for cbar, if relevant.
- **\*\*heatmap\_kws** – Additional keywords to pass to `seaborn.heatmap`

**Returns (List[matplotlib.axes.Axes]):** List of matplotlib axes for figure.

```
hypercluster.visualize.visualize_sample_label_consistency(labels:          pan-
                                                          das.core.frame.DataFrame,
                                                          savefig: bool = False,
                                                          output_prefix: Optional[str] = None,
                                                          **heatmap_kws) →
                                                          List[matplotlib.axes._axes.Axes]
```

Visualize how often two samples are labeled in the same group across conditions. Interpret with care—if you use more conditions for some type of clusterers, e.g. more n\_clusters for KMeans, those cluster more similarly across conditions than between clusterers. This means that more agreement in labeling could be due to the choice of clusterers rather than true similarity between samples.

#### Parameters

- **labels** (*DataFrame*) – Labels DataFrame, e.g. from `optimize_clustering` or `AutoClusterer.labels_`
- **savefig** (*bool*) – Whether to save a pdf.
- **output\_prefix** (*str*) – If saving a pdf, file prefix to use.
- **\*\*heatmap\_kws** – Additional keywords to pass to `seaborn.heatmap`

**Returns (List[matplotlib.axes.Axes]):** List of matplotlib axes

```
hypercluster.visualize.zscore(df)
Row zscores a DataFrame, ignores np.nan
```

**Parameters df** (*DataFrame*) – DataFrame to z-score

**Returns (DataFrame):** Row-zscored DataFrame.

## 1.4 hypercluster.constants module

```
hypercluster.constants.param_delim
```

delimiter between hyperparameters for snakemake file labels and labels DataFrame columns.

```
hypercluster.constants.val_delim
```

delimiter between hyperparameter label and value for snakemake file labels and labels DataFrame columns.

```
hypercluster.constants.categories
```

Convenient groups of clusterers to use. If all samples need to be clustered, ‘partitioners’ is a good choice. If there are millions of samples, ‘fastest’ might be a good choice.

```
hypercluster.constants.variables_to_optimize
```

Some default hyperparameters to optimize and value ranges for a selection of commonly used clustering algorithms from sklearn. Used as defaults for `clustering.AutoClusterer` and `clustering.optimize_clustering`.

```
hypercluster.constants.need_ground_truth
```

list of sklearn metrics that need ground truth labeling. “adjusted\_rand\_score”, “adjusted\_mutual\_info\_score”, “homogeneity\_score”, “completeness\_score”, “fowlkes\_mallows\_score”, “mutual\_info\_score”, “v\_measure\_score”

```
hypercluster.constants.inherent_metrics
```

list of sklearn metrics that need original data for calculation. “silhouette\_score”, “calinski\_harabasz\_score”, “davies\_bouldin\_score”, “smallest\_largest\_clusters\_ratio”, “number\_of\_clusters”, “smallest\_cluster\_size”, “largest\_cluster\_size”

```
hypercluster.constants.min_or_max
```

establishing whether each sklearn metric is better when minimized or maximized for clustering.pick\_labels.

## 1.5 hypercluster.additional\_clusterers module

Additonal clustering classes can be added here, as long as they have a ‘fit’ method.

```
hypercluster.additional_clusterers.HDBSCAN
```

See [hdbscan](#)

**Type** clustering class

```
class hypercluster.additional_clusterers.LeidenCluster(adjacency_method: str
= 'MNN', k: int = 20,
resolution: float = 0.8,
adjacency_kwargs: Optional[dict] = None, partition_type: str = 'RBConfigurationVertexPartition',
**leiden_kwargs)
```

Bases: object

Leidein clustering on graph derived from an adjacency matrix. See [reference](#) for more info

### Parameters

- **adjacency\_method** – Method to use to construct adjacency matrix, which is used to construct graph that will be clustered. Valid methods are any metric valid in `scipy.spatial.distance.pdist`, or MNN, for mutual nearest neighbors and CNN for common nearest neighbors. Both use `sklearn.neighbors.NearestNeighbors` at a given `k` to calculate NNs. MNN then uses whether points `i` and `j` are each others NNs as edge weights. CNN uses the count of how many NNs `i` and `j` have in common as the edge weight.
- **k** – If using CNN or MNN, `k` to use to construct the `NearestNeighbors` matrix.
- **resolution** – If using ‘RBConfigurationVertexPartition’, ‘CPMVertexPartition’ which resolution to use. If using other partitioners, this is ignored but any other kwargs for those partitioners can be passed too.
- **adjacency\_kwargs** – Additional keyword arguments to pass to `sklearn.neighbors.NearestNeighbors` or `scipy.spatial.distance.pdist` to construct the adjacency matrix.
- **partition\_type** – Which partition to use for leiden clustering, see [leidenalg](#) for more info.
- **\*\*leiden\_kwargs** – Additional kwargs to be passed to ‘[find\\_partition](#)’

**fit** (`data: pandas.core.frame.DataFrame`)

```
class hypercluster.additional_clusterers.LouvainCluster (adjacency_method: str  
= 'MNN', k: int = 20,  
resolution: float = 0.8,  
adjacency_kwargs: Optional[dict] = None,  
partition_type: str = 'RB-  
ConfigurationVertexPartition', **louvain_kwargs)
```

Bases: object

Louvain clustering on graph derived from an adjacency matrix.

#### Parameters

- **adjacency\_method** – Method to use to construct adjacency matrix, which is used to construct graph that will be clustered. Valid methods are any metric valid in `scipy.spatial.distance.pdist`, or MNN, for mutual nearest neighbors and CNN for common nearest neighbors. Both use `sklearn.neighbors.NearestNeighbors` at a given `k` to calculate NNs. MNN then uses whether points `i` and `j` are each others NNs as edge weights. CNN uses the count of how many NNs `i` and `j` have in common as the edge weight.
- **k** – If using CNN or MNN, `k` to use to construct the `NearestNeighbors` matrix.
- **resolution** – If using ‘RBConfigurationVertexPartition’, ‘CPMVertexPartition’ which resolution to use. If using other partitioners, this is ignored but any other kwargs for those partitioners can be passed too.
- **adjacency\_kwargs** – Additional keyword arguments to pass to `sklearn.neighbors.NearestNeighbors` or `scipy.spatial.distance.pdist` to construct the adjacency matrix.
- **partition\_type** – Which partition to use for louvain clustering, see `louvain-igraph` for more info.
- **\*\*louvain\_kwargs** – Additional kwargs to be passed to ‘`find_partition`’

**fit** (data: `pandas.core.frame.DataFrame`)

```
class hypercluster.additional_clusterers.NMFCluster (n_clusters: int = 8,  
**nmf_kwargs)
```

Bases: object

Uses non-negative factorization from `sklearn` to assign clusters to samples, based on the maximum membership score of the sample per component.

#### Parameters

- **n\_clusters** – The number of clusters to find. Used as `n_components` when fitting.
- **\*\*nmf\_kwargs** –

**fit** (data)

If negative numbers are present, creates one data matrix with all negative numbers zeroed. Create another data matrix with all positive numbers zeroed and the signs of all negative numbers reversed. Concatenate both matrices resulting in a data matrix twice as large as the original, but with positive values only and zeros and hence appropriate for NMF. Uses decomposed matrix `H`, which is `nxk` (with `n`=number of samples and `k`=number of components) to assign cluster membership. Each sample is assigned to the cluster for which it has the highest membership score. See `sklearn.decomposition.NMF`

**Parameters** `data` (`DataFrame`) – Data to fit with samples as rows and features as columns.

**Returns** self with `labels_` attribute.

## 1.6 `hypercluster.additional_metrics` module

More functions for evaluating clustering results. Additional metric evaluations can be added here, as long as the second argument is the labels to evaluate

```
hypercluster.additional_metrics.largest_cluster_size(_, labels: Iterable[T_co]) → float
    Number in largest cluster
```

### Parameters

- `_` – Dummy, pass anything or None
- **labels** (`Iterable`) – Vector of sample labels.

**Returns (int):** Number of samples in largest cluster.

```
hypercluster.additional_metrics.number_clustered(_, labels: Iterable[T_co]) → float
    Returns the number of clustered samples.
```

### Parameters

- `_` – Dummy, pass anything or None.
- **labels** (`Iterable`) – Vector of sample labels.

**Returns (int):** The number of clustered labels.

```
hypercluster.additional_metrics.number_of_clusters(_, labels: Iterable[T_co]) → float
    Number of total clusters.
```

### Parameters

- `_` – Dummy, pass anything or None
- **labels** (`Iterable`) – Vector of sample labels.

**Returns (int):** Number of clusters.

```
hypercluster.additional_metrics.smallest_cluster_ratio(_, labels: Iterable[T_co]) → float
    Number in the smallest cluster over the total samples.
```

### Parameters

- `_` – Dummy, pass anything or None.
- **labels** (`Iterable`) – Vector of sample labels.

**Returns (float):** Ratio of number of members in smallest over all samples.

```
hypercluster.additional_metrics.smallest_cluster_size(_, labels: Iterable[T_co]) → float
    Number in smallest cluster
```

### Parameters

- `_` – Dummy, pass anything or None
- **labels** (`Iterable`) – Vector of sample labels.

**Returns (int):** Number of samples in smallest cluster.

```
hypercluster.additional_metrics.smallest_largest_clusters_ratio(_ labels: Iterable[T_co])  
→ float
```

Number in the smallest cluster over the number in the largest cluster.

#### Parameters

- `_` – Dummy, pass anything or None.
- `labels (Iterable)` – Vector of sample labels.

**Returns (float):** Ratio of number of members in smallest over largest cluster.



---

 CHAPTER  
 TWO
 

---

## HYPERCLUSTER SNAKEMAKE PIPELINE

### 2.1 Line-by-line explanation of config.yml

Table 1: Explanation for config.yml

config.yml parameter	Explanation	Example from scRNA-seq workflow
input_data_folder	Path to folder in which input data can be found. No / at the end.	/input_data
input_data_files	List of prefixes of data files. Exclude extension, .csv, .tsv and .txt allowed.	['input_data1', 'input_data2']
gold_standard_file	File name of gold_standard_file. Must have same pandas.read_csv kwargs as the corresponding input file. Must be in input_data_folder.	{'input_data': 'gold_standard_file.txt'}
read_csv_kwargs	Per input data file, keyword args to put into pandas.read_csv. <b>If specifying multiindex, also put the same in output_kwargs['labels']</b>	{'test_input': {'index_col': [0]}}
output_folder	Path to folder in which results will be written. No / at the end.	/hypercluster_results
intermediates_folder	Name of the folder within the output_folder to put intermediate results, such as labels and evaluations per condition. No need to change this usually.	clustering_intermediates
clustering_results		clustering
18	Chapter 2. Name of the folder within the output_folder to put final results. No need to change this usually.	hypercluster SnakeMake pipeline

\*\*Note: Formatting of lists and dictionaries can be in python syntax (like above) or yaml syntax, or a mixture, like below. \*\*

## 2.2 config.yml example from scRNA-seq workflow

```

input_data_folder: '.'
input_data_files:
  - sc_data
gold_standards:
  test_input: 'gold_standard.csv'
read_csv_kwargs:
  test_input: {'index_col': [0]}

output_folder: 'results'
intermediates_folder: 'clustering_intermediates'
clustering_results: 'clustering'

clusterer_kwargs: {}
generate_parameters_addtl_kwargs: {}

evaluations:
  - silhouette_score
  - calinski_harabasz_score
  - davies_bouldin_score
  - number_clustered
  - smallest_largest_clusters_ratio
  - smallest_cluster_ratio
eval_kwargs: {}

metric_to_choose_best: silhouette_score
metric_to_compare_labels: adjusted_rand_score
compare_samples: true

output_kwargs:
  evaluations:
    index_col: [0]
  labels:
    index_col: [0]
heatmap_kwargs: {}

optimization_parameters:
  HDBSCAN:
    min_cluster_size: &id002
    - 2
    - 3
    - 4
    - 5
  KMeans:
    n_clusters: &id001
    - 5
    - 6
    - 7
  MiniBatchKMeans:
    n_clusters: *id001
  OPTICS:
    min_samples: *id002

```

(continues on next page)

(continued from previous page)

```
NMFCluster:  
    n_clusters: *id001  
LouvainCluster: &id003  
    resolution:  
        - 0.2  
        - 0.4  
        - 0.6  
        - 0.8  
        - 1.0  
        - 1.2  
        - 1.4  
        - 1.6  
    k:  
        - 10  
        - 15  
        - 20  
        - 40  
        - 80  
        - 120  
LeidenCluster: *id003
```

---

**CHAPTER  
THREE**

---

**INDICES AND TABLES**

- genindex
- modindex
- search



## INSTALLATION AND LOGISTICS

### 4.1 Installation

Available via pip:

```
pip install hypercluster
```

Or bioconda:

```
conda install hypercluster
# or
conda install -c conda-forge -c bioconda hypercluster
```

If you are having problems installing with conda, try changing your channel priority. Priority of conda-forge > bioconda > defaults is recommended.

To check channel priority: `conda config --get channels`

It should look like:

```
--add channels 'defaults'    # lowest priority
--add channels 'bioconda'
--add channels 'conda-forge'   # highest priority
```

If it doesn't look like that, try:

```
conda config --add channels bioconda
conda config --add channels conda-forge
```

## 4.2 Quick reference for clustering and evaluation

Table 1: Clustering algorithms

Clusterer	Type
KMeans/MiniBatch KMeans	Partitioner
Affinity Propagation	Partitioner
Mean Shift	Partitioner
DBSCAN	Clusterer
OPTICS	Clusterer
Birch	Partitioner
OPTICS	Clusterer
HDBSCAN	Clusterer
NMF	Partitioner
LouvainCluster	Partitioner
LeidenCluster	Partitioner

Table 2: Evaluations

Metric	Type
adjusted_rand_score	Needs ground truth
adjusted_mutual_info_score	Needs ground truth
homogeneity_score	Needs ground truth
completeness_score	Needs ground truth
fowlkes_mallows_score	Needs ground truth
mutual_info_score	Needs ground truth
v_measure_score	Needs ground truth
silhouette_score	Inherent metric
calinski_harabasz_score	Inherent metric
davies_bouldin_score	Inherent metric
smallest_largest_clusters_ratio	Inherent metric
number_of_clusters	Inherent metric
smallest_cluster_size	Inherent metric
largest_cluster_size	Inherent metric

## 4.3 Quickstart and examples

### 4.3.1 With snakemake:

```
snakemake -s hypercluster.smk --configfile config.yml --config input_data_files=test_
→data input_data_folder=.
```

### 4.3.2 With python:

```
import pandas as pd
from sklearn.datasets import make_blobs
import hypercluster
```

(continues on next page)

(continued from previous page)

```
data, labels = make_blobs()
data = pd.DataFrame(data)
labels = pd.Series(labels, index=data.index, name='labels')

# With a single clustering algorithm
clusterer = hypercluster.AutoClusterer()
clusterer.fit(data).evaluate(
    methods = hypercluster.constants.need_ground_truth+hypercluster.constants.inherent_
    ↵metrics,
    gold_standard = labels
)

clusterer.visualize_evaluations()

# With a range of algorithms

clusterer = hypercluster.MultiAutoClusterer()
clusterer.fit(data).evaluate(
    methods = hypercluster.constants.need_ground_truth+hypercluster.constants.inherent_
    ↵metrics,
    gold_standard = labels
)

clusterer.visualize_evaluations()
```

Example work flows for both python and snakemake are [here](#)

Source code is available [here](#)



## PYTHON MODULE INDEX

### h

hypercluster, 1  
hypercluster.additional\_clusterers, 13  
hypercluster.additional\_metrics, 15  
hypercluster.classes, 5  
hypercluster.constants, 12  
hypercluster.utilities, 8  
hypercluster.visualize, 10



# INDEX

## A

algorithm\_clus\_kwarg  
    `ter.classes.MultiAutoClusterer`, 7  
algorithm\_clus\_kwarg  
    `ter.MultiAutoClusterer attribute`, 3  
algorithm\_names  
    `ter.classes.MultiAutoClusterer`, 7  
algorithm\_names (`hypercluster.MultiAutoClusterer attribute`), 3  
algorithm\_param\_weights  
    `ter.classes.MultiAutoClusterer`, 7  
algorithm\_param\_weights  
    `ter.MultiAutoClusterer attribute`, 3  
algorithm\_parameters  
    `ter.classes.MultiAutoClusterer`, 7  
algorithm\_parameters  
    `ter.MultiAutoClusterer attribute`, 3  
AutoClusterer (class in `hypercluster`), 1  
AutoClusterer (class in `hypercluster.classes`), 5  
autoclusterers  
    `hypercluster.classes.MultiAutoClusterer` (attribute), 8  
autoclusterers (`hypercluster.MultiAutoClusterer attribute`), 4

## C

calculate\_row\_weights () (in module `hypercluster.utilities`), 8  
categories (in module `hypercluster.constants`), 12  
clus\_kwarg (`hypercluster.AutoClusterer attribute`), 1  
clus\_kwarg (`hypercluster.classes.AutoClusterer attribute`), 5  
cluster () (in module `hypercluster.utilities`), 9  
Clusterer (class in `hypercluster.classes`), 6  
clusterer\_name (`hypercluster.AutoClusterer attribute`), 1  
clusterer\_name (`hypercluster.classes.AutoClusterer attribute`), 5

compute\_order () (in module `hypercluster.visualize`), 10  
convert\_to\_multiind () (in module `hypercluster.utilities`), 9

## D

data (`hypercluster.AutoClusterer attribute`), 2  
data (`hypercluster.classes.AutoClusterer attribute`), 5  
data (`hypercluster.classes.MultiAutoClusterer attribute`), 8  
data (`hypercluster.MultiAutoClusterer attribute`), 3

## E

evaluate () (`hypercluster.AutoClusterer method`), 2  
evaluate () (`hypercluster.classes.AutoClusterer method`), 6  
evaluate () (`hypercluster.classes.MultiAutoClusterer method`), 8  
evaluate () (`hypercluster.MultiAutoClusterer method`), 4  
evaluate\_one () (in module `hypercluster.utilities`), 9  
evaluation\_ (`hypercluster.AutoClusterer attribute`), 2  
evaluation\_ (`hypercluster.classes.AutoClusterer attribute`), 5  
evaluation\_ (`hypercluster.classes.MultiAutoClusterer attribute`), 8

evaluation\_ (`hypercluster.MultiAutoClusterer attribute`), 4  
evaluation\_df  
    `hypercluster.classes.MultiAutoClusterer` (attribute), 8  
evaluation\_df (`hypercluster.MultiAutoClusterer attribute`), 4  
evaluation\_methods  
    `hypercluster.classes.MultiAutoClusterer` (attribute), 8  
evaluation\_methods  
    `hypercluster.classes.MultiAutoClusterer attribute`, 3

## F

fit () (`hypercluster.additional_clusterers.LeidenCluster`

```

        method), 13
fit() (hypercluster.additional_clusterers.LouvainCluster
       method), 14
fit() (hypercluster.additional_clusterers.NMFCluster
       method), 14
fit() (hypercluster.AutoClusterer method), 2
fit() (hypercluster.classes.AutoClusterer method), 6
fit() (hypercluster.classes.MultiAutoClusterer
       method), 8
fit() (hypercluster.MultiAutoClusterer method), 4
fit_predict() (hypercluster.classes.Clusterer
       method), 6

```

## G

```

generate_flattened_df() (in module hyperclus-
ter.utilities), 9
generate_param_sets() (hyperclus-
ter.AutoClusterer method), 2
generate_param_sets() (hyperclus-
ter.classes.AutoClusterer method), 6
gold_standard (hyperclus-
ter.classes.MultiAutoClusterer
       attribute), 8
gold_standard (hypercluster.MultiAutoClusterer at-
tribute), 4

```

## H

```

HDBSCAN (in module hyperclus-
ter.additional_clusterers), 13
hypercluster (module), 1
hypercluster.additional_clusterers (mod-
ule), 13
hypercluster.additional_metrics (module),
15
hypercluster.classes (module), 5
hypercluster.constants (module), 12
hypercluster.utilities (module), 8
hypercluster.visualize (module), 10

```

## I

```

inherent_metrics (in module hyperclus-
ter.constants), 12

```

## L

```

labels_ (hypercluster.AutoClusterer attribute), 1
labels_ (hypercluster.classes.AutoClusterer attribute),
5
labels_ (hypercluster.classes.MultiAutoClusterer at-
tribute), 8
labels_ (hypercluster.MultiAutoClusterer attribute), 4
labels_df (hypercluster.classes.MultiAutoClusterer
       attribute), 8
labels_df (hypercluster.MultiAutoClusterer at-
tribute), 4

```

## M

```

largest_cluster_size() (in module hyperclus-
ter.additional_metrics), 15
LeidenCluster (class in hyperclus-
ter.additional_clusterers), 13
LouvainCluster (class in hyperclus-
ter.additional_clusterers), 13

```

## N

```

metric_kwargs (hypercluster.classes.MultiAutoClusterer
       attribute), 8
metric_kwargs (hypercluster.MultiAutoClusterer at-
tribute), 4
min_or_max (in module hypercluster.constants), 13
MultiAutoClusterer (class in hypercluster), 2
MultiAutoClusterer (class in hyperclus-
ter.classes), 6

```

## P

```

need_ground_truth (in module hyperclus-
ter.constants), 12
NMFCluster (class in hyperclus-
ter.additional_clusterers), 14
number_clustered() (in module hyperclus-
ter.additional_metrics), 15
number_of_clusters() (in module hyperclus-
ter.additional_metrics), 15

```

## R

```

param_delim (in module hypercluster.constants), 12
param_weights (hypercluster.AutoClusterer at-
tribute), 1
param_weights (hypercluster.classes.AutoClusterer
       attribute), 5
params_to_optimize (hypercluster.AutoClusterer
       attribute), 1
params_to_optimize (hyperclus-
ter.classes.AutoClusterer attribute), 5
pick_best_labels() (hyperclus-
ter.classes.Clusterer method), 6
pick_best_labels() (in module hyperclus-
ter.utilities), 10

```

## R

```

random_search (hypercluster.AutoClusterer at-
tribute), 1
random_search (hypercluster.classes.AutoClusterer
       attribute), 5
random_search (hyperclus-
ter.classes.MultiAutoClusterer
       attribute), 7
random_search (hypercluster.MultiAutoClusterer at-
tribute), 3

```

random\_search\_fraction (hypercluster.AutoClusterer attribute), 1  
random\_search\_fraction (hypercluster.classes.AutoClusterer attribute), 5  
random\_search\_fraction (hypercluster.classes.MultiAutoClusterer attribute), 7  
random\_search\_fraction (hypercluster.MultiAutoClusterer attribute), 3

## S

smallest\_cluster\_ratio() (in module hypercluster.additional\_metrics), 15  
smallest\_cluster\_size() (in module hypercluster.additional\_metrics), 15  
smallest\_largest\_clusters\_ratio() (in module hypercluster.additional\_metrics), 15

## V

val\_delim (in module hypercluster.constants), 12  
variables\_to\_optimize (in module hypercluster.constants), 12  
visualize\_evaluations() (hypercluster.classes.Clusterer method), 6  
visualize\_evaluations() (in module hypercluster.visualize), 10  
visualize\_for\_picking\_labels() (hypercluster.classes.Clusterer method), 6  
visualize\_for\_picking\_labels() (in module hypercluster.visualize), 11  
visualize\_label\_agreement() (hypercluster.classes.Clusterer method), 6  
visualize\_label\_agreement() (in module hypercluster.visualize), 11  
visualize\_pairwise() (in module hypercluster.visualize), 11  
visualize\_sample\_label\_consistency() (hypercluster.classes.Clusterer method), 6  
visualize\_sample\_label\_consistency() (in module hypercluster.visualize), 12

## Z

zscore() (in module hypercluster.visualize), 12