
episcanpy Documentation

Release 0.2.0+66.g32c2282

Anna Danese

Aug 07, 2020

CONTENTS

1	Version 0.2.0 August 7, 2020	3
2	Version 0.1.8 November 5, 2019	5
3	Version 0.1.7 November 5, 2019	7
4	Version 0.1.0 May 10, 2019	9
	Bibliography	23
	Python Module Index	25
	Index	27

EpiScanpy is a toolkit to analyse single-cell open chromatin (scATAC-seq) and single-cell DNA methylation (for example scBS-seq) data. **EpiScanpy** is the epigenomic extension of the very popular scRNA-seq analysis tool **Scanpy** (*Genome Biology*, 2018) [Wolf18]. For more information, read scanpy [documentation](#).

The documentation for epiScanpy is available [here](#). If epiScanpy is useful to your research, consider citing [epiScanpy](#). Report issues and access the code on [GitHub](#).

Note: Also see the [release notes](#) of scanpy.

Also see the [release notes](#) of anndata.

VERSION 0.2.0 AUGUST 7, 2020

This release deal with the compatibility problems with the latest version of scanpy. Additionally, it contains new features to build quick custom count matrices (`bld_mtx_fly`), to convert snap into h5ad files (`snap2anndata`) or build gene activity matrices (`geneactivity`).

VERSION 0.1.8 NOVEMBER 5, 2019

Release new processing function & quality controls.

VERSION 0.1.7 NOVEMBER 5, 2019

Release for SCOG epiScanpy Hackathon in Saarbrücken.

This version is not fully compatible with previous version.

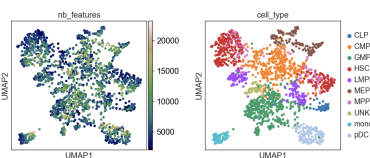
VERSION 0.1.0 MAY 10, 2019

Initial release.

4.1 Tutorials

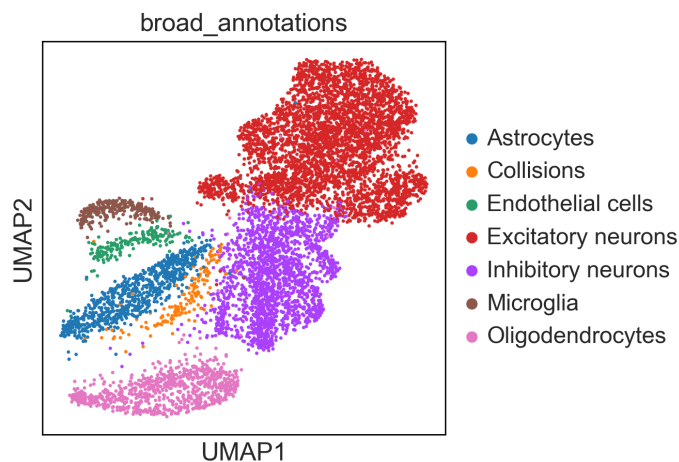
4.1.1 Single cell ATAC-seq

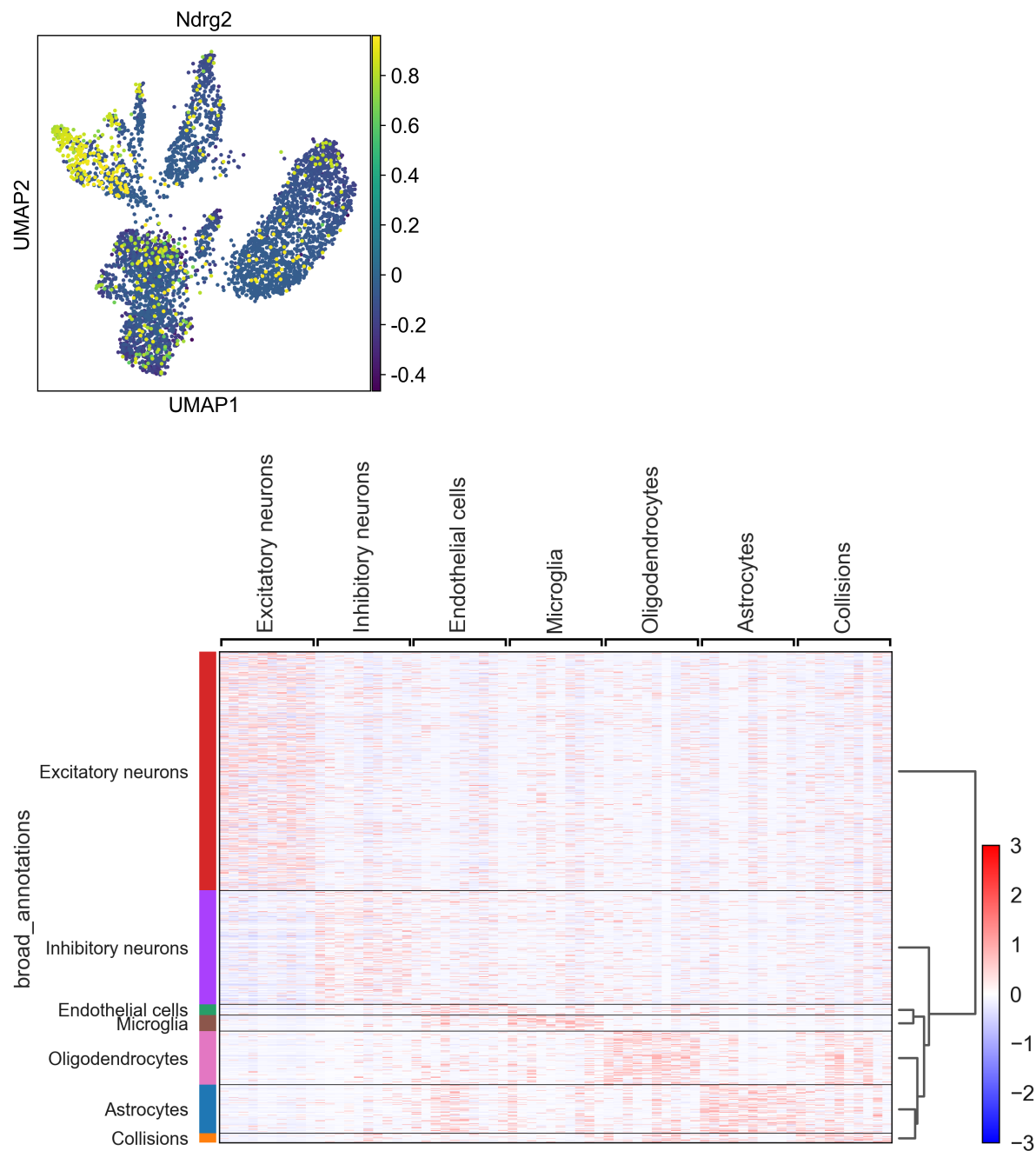
To get started, we recommend epiScanpy's analysis pipeline for scATAC-seq data from Buenrostro et al. [Buenrostro18]., the dataset consist of ~3000cells of human PBMCs. This tutorial focuses on preprocessing, clustering, identification of cell types via known marker genes and trajectory inference. The tutorial can be found [here](#).



If you want to see how to build count matrices from ATAC-seq bam files for different set of annotations (like enhancers). The tutorial can be found [here](#).

Soon available, there will be a tutorial providing a function to very quickly build custom count matrices using standard 10x single cell ATAC output.





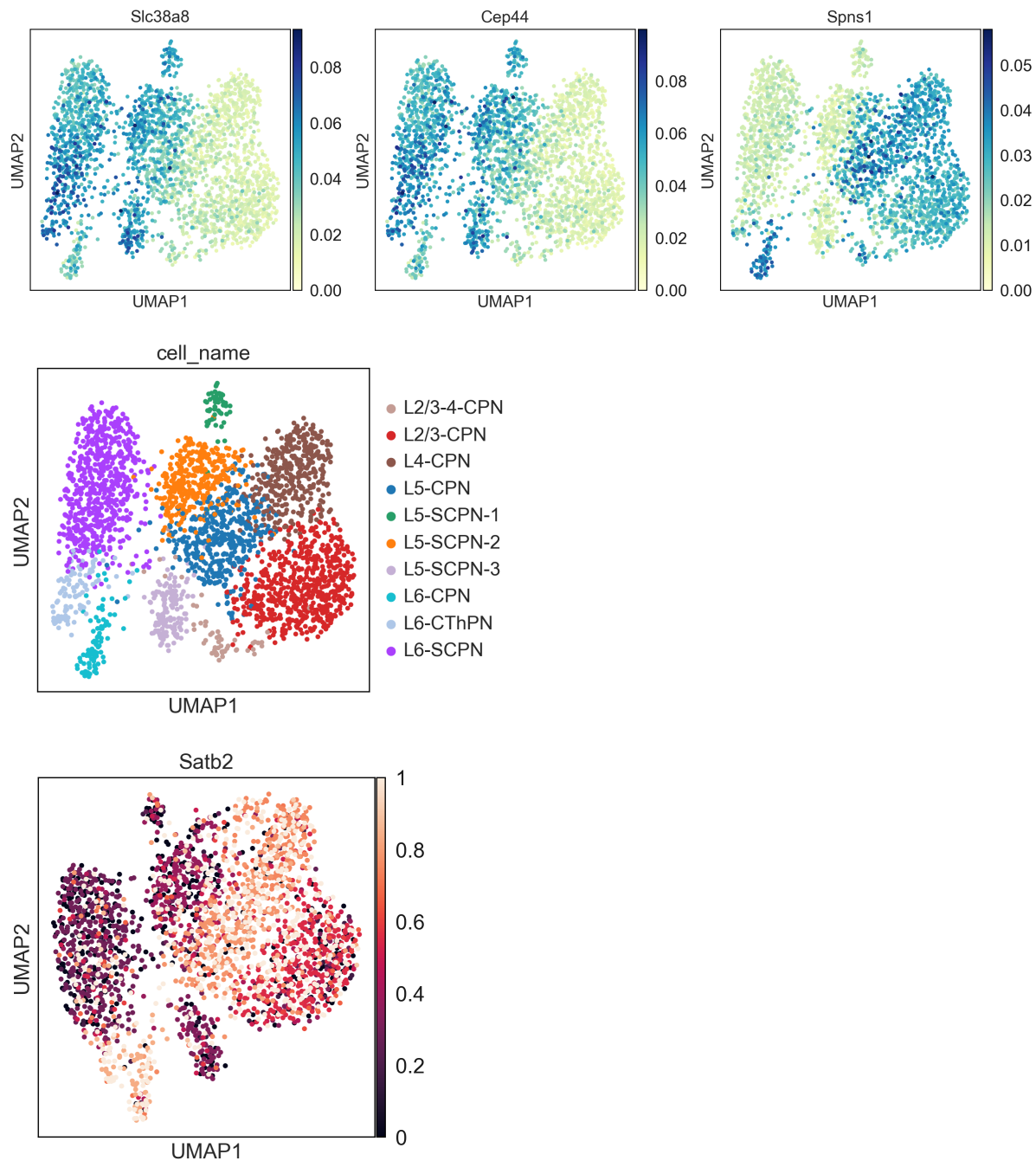
An additional tutorial on processing and clustering count matrices from the Cusanovich mouse scATAC-seq atlas [Cusanovich18].. [Here](#).

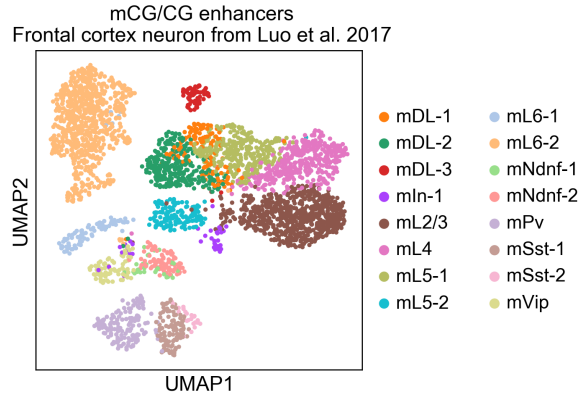
4.1.2 Single cell DNA methylation

Here you can find a tutorial for the preprocessing, clustering and identification of cell types for single-cell DNA methylation data using the publicly available data from Luo et al. [Luo17].

The first tutorial shows how to build the count matrices for the different feature spaces (windows, promoters) in different cytosine contexts. Here is the [tutorial](#).

Then, there is a second tutorial on how to use them and compare the results. The data used comes from mouse brain (frontal cortex). It will be available very soon.





4.2 Usage Principles

Import the epiScanpy API as:

```
import episcanpy.api as epi
import anndata as ad
```

4.2.1 Workflow

The first step is to build the count matrix. Because single-cell epigenomic data types have different characteristics (count data in ATAC-seq versus methylation level in DNA methylation, for example), epiScanpy implements -omic specific approaches to build the count matrix. All the functions to build the count matrices (for ATAC, methylation or other) will use `epi.ct` (ct = count).

The first step is to load an annotation and then build the count matrix that will be either methylation or ATAC-seq specific. For example using `epi.ct`, e.g.:

```
epi.ct.load_features(file_features, **tool_params) # to load annotation files
epi.ct.build_count_mtx(cell_file_names, omic="ATAC") # to build the ATAC-seq count_
↪matrix
```

If you have an already build matrix, you can load it with any additional metadata (such as cell annotations or batches).

The count matrix, either the one that has been constructed or uploaded, with any additional informations (such as cell annotations or batches) are stored as an `AnnData` object. All functions for quality control and preprocessing are called using `epi.pp` (pp = preprocessing).

To visualise how common features are and what is the coverage distribution of the count matrix features, use:

```
epi.pp.commonness_features(adata, **processing_params)
epi.pp.coverage_cells(adata, **processing_params)
```

To remove low quality cells you can use the following functions:

```
epi.pp.filter_cells(adata, min_features=10)
epi.pp.filter_features(adata, min_cells=10)
```

To reduce the feature space to the most variable features: `:: epi.pl.cal_var(adata) epi.pp.select_var_feature(adata, max_score=0.2, nb_features=50000)`

The next step, is the calculation of tSNE, UMAP, PCA etc. For that, we take advantage of the embedding into Scanpy and we use mostly Scanpy functions, which are called using `sc.tl` (tl = tool) [Wolf18]. For that, see Scanpy usage principles: <https://scanpy.readthedocs.io/en/latest/basic_usage.html>`. For example, to obtain cell-cell distance calculations or low dimensional representation we make use of the `adata` object, and store `n_obs` observations (cells) of `n_vars` variables (expression, methylation, chromatin features). For each tool, there typically is an associated plotting function in `sc.tl` and `sc.pl` (pl = plot)

```
epi.pp.pca(adata, n_comps=100, svd_solver='arpack')
epi.pp.neighbors(adata, n_neighbors=15)

epi.tl.tsne(adata, **tool_params)
epi.pl.tsne(adata, **plotting_params)
```

There are also epiScanpy specific tools and plotting functions that can be accessed using `epi.tl` and `epi.pl`

```
epi.tl.silhouette(adata, **tool_params)
epi.pl.silhouette(adata, **plotting_params)
epi.pl.prct_overlap(adata, **plotting_params)
```

4.2.2 Data structure

Similarly to Scanpy, the methylation and ATAC-seq matrices are stored as Anndata objects. For more information on the datastructure see here <<https://anndata.readthedocs.io/en/latest/>>`.

4.3 System Requirements

4.3.1 Hardware requirements

epiScanpy package requires only a standard computer with enough RAM to support the in-memory operations.

4.3.2 Software requirements

OS Requirements This package is supported for *macOS* and *Linux*. The package has been tested on the following systems: + macOS: Mojave & Catalina (10.14 to 10.15.4)

4.3.3 Python Dependencies

EpiScanpy require a working version of Python (≥ 3.6)

Additionally, this package epiScanpy depends on other Python dependencies and packages.:

```
anndata
matplotlib
numpy
pandas
pyliftOver
pysam
scanpy
scipy
scikit-learn
```

(continues on next page)

(continued from previous page)

```
seaborn
bamnostic
```

4.4 Installation

4.4.1 Anaconda

If you do not have a working Python 3.5 or 3.6 installation, consider installing Miniconda (see [Installing Miniconda](#)). Then run:

```
conda install seaborn scikit-learn statsmodels numba
conda install -c conda-forge python-igraph louvain
conda create -n scanpy python=3.6 scanpy
```

Finally, run:

```
conda install -c annadanese episcanpy
```

Pull epiScanpy from [PyPI](#) (consider using `pip3` to access Python 3):

```
pip install episcanpy
```

4.4.2 Github

you can also install epiScanpy directly from Github:

```
pip install git+https://github.com/colomemaria/epiScanpy
```

4.5 API

Import epiScanpy's high-level API as:

```
import episcanpy.api as epi
```

4.5.1 Count Matrices: CT

Loading data, loading annotations, building count matrices, filtering of lowly covered methylation variables. Filtering of lowly covered cells.

Building count matrices

Quickly build a count matrix from tsv/tbi file.

<code>ct.bld_mtx_fly(tsv_file, annotation[, ...])</code>	Building count matrix on the fly.
--	-----------------------------------

episcanpy.ct.bld_mtx_fly

Load features

In order to build a count matrix for either methylation or open chromatin data, loading the segmentation of the genome of interest or the set of features of interest is a prerequisite.

<code>ct.load_features(file_features[, ...])</code>	The function load features is here to transform a bed file into a usable set of units to measure methylation levels.
<code>ct.make_windows(size[, chromosomes, max_length])</code>	Generate windows/bins of the given size for the appropriate genome (default choice is human).
<code>ct.size_feature_norm(loaded_feature, size)</code>	If the features loaded are too smalls or of different sizes, it is possible to normalise them to a unique given size by extending the feature coordinate in both directions.
<code>ct.plot_size_features(loaded_feature[, ...])</code>	Plot the different feature sizes in an histogram.
<code>ct.name_features(loaded_features)</code>	Extract the names of the loaded features, specifying the chromosome they originated from.

episcanpy.ct.load_features

episcanpy.ct.make_windows

episcanpy.ct.size_feature_norm

episcanpy.ct.plot_size_features

episcanpy.ct.name_features

Reading methylation file

Functions to read methylation files, extract methylation and build the count matrices:

<code>ct.build_count_mtx(cells, annotation[, ...])</code>	Build methylation count matrix for a given annotation.
<code>ct.read_cyt_summary(sample_name, meth_type, ...)</code>	Read file from which you want to extract the methylation level and (assuming it is like the Ecker/Methylpy format) extract the number of methylated read and the total number of read for the cytosines covered and in the right genomic context (CG or CH) :param sample_name: name of the file to read to extract key information.
<code>ct.load_met_noimput(matrix_file[, path, save])</code>	read the raw count matrix and convert it into an AnnData object.

`episcanpy.ct.build_count_mtx`

`episcanpy.ct.read_cyt_summary`

`episcanpy.ct.load_met_noimput`

Reading open chromatin(ATAC) file

ATAC-seq specific functions to build count matrices and load data:

<code>ct.bld_atac_mtx(list_bam_files, loaded_feat)</code>	Build a count matrix one set of features at a time.
<code>ct.save_sparse_mtx(initial_matrix[, ...])</code>	Convert regular atac matrix into a sparse Anndata:

`episcanpy.ct.bld_atac_mtx`

`episcanpy.ct.save_sparse_mtx`

General functions

Functions non -omic specific:

<code>ct.save_sparse_mtx(initial_matrix[, ...])</code>	Convert regular atac matrix into a sparse Anndata:
--	--

4.5.2 Preprocessing: PP

Imputing missing data (methylation), filtering lowly covered cells or variables, correction for batch effect.

<code>pp.coverage_cells(adata[, key_added, log, ...])</code>	Histogram of the number of open features (in the case of ATAC-seq data) per cell.
<code>pp.correlation_pc(adata, variable[, pc, ...])</code>	Correlation between a given PC and a covariate.
<code>pp.coverage_features(adata[, binary, log, ...])</code>	Display how often a feature is measured as open (for ATAC-seq).
<code>pp.density_features(adata[, threshold, ...])</code>	Display how often a feature is measured as open (for ATAC-seq).
<code>pp.select_var_feature(adata[, min_score, ...])</code>	This function computes a variability score to rank the most variable features across all cells.
<code>pp.cal_var(adata[, show, color, save])</code>	Show distribution plots of cells sharing features and variability score.
<code>pp.variability_features(adata[, min_score, ...])</code>	This function computes a variability score to rank the most variable features across all cells.
<code>pp.binarize(adata[, copy])</code>	convert the count matrix into a binary matrix.
<code>pp.lazy(adata[, pp_pca, svd_solver, nb_pcs, ...])</code>	Automatically computes PCA coordinates, loadings and variance decomposition, a neighborhood graph of observations, t-distributed stochastic neighborhood embedding (tSNE) Uniform Manifold Approximation and Projection (UMAP)
<code>pp.load_metadata(adata, metadata_file[, ...])</code>	Load observational metadata in adata.obs.

continues on next page

Table 6 – continued from previous page

<code>pp.read_ATAC_10x(matrix[, cell_names, ...])</code>	Load sparse matrix (including matrices corresponding to 10x data) as AnnData objects.
<code>pp.filter_cells(adata[, min_counts, ...])</code>	Filter cell outliers based on counts and numbers of genes expressed.
<code>pp.filter_features(data[, min_counts, ...])</code>	Filter features based on number of cells or counts.
<code>pp.normalize_total(adata[, target_sum, ...])</code>	Normalize counts per cell.
<code>pp.pca(adata[, n_comps, zero_center, ...])</code>	Principal component analysis [Pedregosa11].
<code>pp.normalize_per_cell(adata[, ...])</code>	Normalize total counts per cell.
<code>pp.regress_out(adata, keys[, n_jobs, copy])</code>	Regress out unwanted sources of variation.
<code>pp.subsample(data[, fraction, n_obs, ...])</code>	Subsample to a fraction of the number of observations.
<code>pp.downsample_counts(adata[, ...])</code>	Downsample counts from count matrix.
<code>pp.neighbors(adata[, n_neighbors, n_pcs, ...])</code>	Compute a neighborhood graph of observations [McInnes18].
<code>pp.sparse(adata[, copy])</code>	Transform adata.X from a matrix or array to a csc sparse matrix.
<code>pp.sparse(adata[, copy])</code>	Transform adata.X from a matrix or array to a csc sparse matrix.

episcanpy.pp.coverage_cells

episcanpy.pp.correlation_pc

episcanpy.pp.coverage_features

episcanpy.pp.density_features

episcanpy.pp.select_var_feature

episcanpy.pp.cal_var

episcanpy.pp.variability_features

episcanpy.pp.binarize

episcanpy.pp.lazy

episcanpy.pp.load_metadata

episcanpy.pp.read_ATAC_10x

episcanpy.pp.filter_cells

episcanpy.pp.filter_features

episcanpy.pp.normalize_total

episcanpy.pp.pca

episcanpy.pp.normalize_per_cell

episcanpy.pp.regress_out

episcanpy.pp.subsample

episcanpy.pp.downsample_counts

episcanpy.pp.neighbors

episcanpy.pp.sparse

Methylation matrices

Methylation specific count matrices.

<code>pp.imputation_met(adata[,...])</code>	Impute missing values in methylation level matrices.
<code>pp.load_met_noimput(matrix_file[, path, save])</code>	read the raw count matrix and convert it into an AnnData object.
<code>pp.readandimputematrix(file_name[, min_coverage])</code>	Temporary function to load and impute methylation count matrix into an AnnData object

episcanpy.pp.imputation_met

episcanpy.pp.load_met_noimput

episcanpy.pp.readandimputematrix

4.5.3 Tools: TL

<code>tl.rank_features(adata, groupby[, omic, ...])</code>	It is a wrap-up function of scanpy <code>sc.tl.rank_genes_groups</code> function.
<code>tl.lazy(adata[, pp_pca, copy])</code>	Automatically computes PCA coordinates, loadings and variance decomposition, a neighborhood graph of observations, t-distributed stochastic neighborhood embedding (tSNE) Uniform Manifold Approximation and Projection (UMAP)
<code>tl.load_markers(path, marker_list_file)</code>	Convert list of known cell type markers from literature to a dictionary Input list of known marker genes First row is considered the header
<code>tl.identify_cluster(adata, cell_type, ...[, ...])</code>	Use markers of a given cell type to plot peak openness for peaks in promoters of the given markers Input cell type, cell type markers, peak promoter intersections
<code>tl.top_feature_genes(adata, gtf_file[, ...])</code>	Deprecated - Please use <code>epi.tl.var_features_to_genes</code> instead.
<code>tl.var_features_to_genes(adata, gtf_file[, ...])</code>	Once you called the most variable features.
<code>tl.geneactivity(adata, gtf_file[, ...])</code>	merge values of peaks/windows/features overlapping genebodies + 2kb upstream.
<code>tl.diffmap(adata[, n_comps, copy])</code>	Diffusion Maps [Coifman05] [Haghverdi15] [Wolf18].

continues on next page

Table 8 – continued from previous page

<code>tl.draw_graph(adata[, layout, init_pos, ...])</code>	Force-directed graph drawing [Islam11] [Jacomy14] [Chippada18].
<code>tl.tsne(adata[, n_pcs, use_rep, perplexity, ...])</code>	t-SNE [Maaten08] [Amir13] [Pedregosa11].
<code>tl.umap(adata[, min_dist, spread, ...])</code>	Embed the neighborhood graph using UMAP [McInnes18].
<code>tl.dpt(adata[, n_dcs, n_branchings, ...])</code>	Infer progression of cells through geodesic distance along the graph [Haghverdi16] [Wolf19].
<code>tl.louvain(adata[, resolution, ...])</code>	Cluster cells into subgroups [Blondel08] [Levine15] [Traag17].
<code>tl.leiden(adata[, resolution, restrict_to, ...])</code>	Cluster cells into subgroups [Traag18].
<code>tl.kmeans(adata, num_clusters)</code>	Compute kmeans clustering using X_pca fits.
<code>tl.hc(adata, num_clusters)</code>	Compute hierarchical clustering using X_pca fits.
<code>tl.getNClusters(adata, n_cluster[, ...])</code>	Function will test different settings of louvain to obtain the target number of clusters.
<code>tl.dendrogram(adata, groupby[, n_pcs, ...])</code>	Computes a hierarchical clustering for the given <i>groupby</i> categories.
<code>tl.ARI(adata, label_1, label_2)</code>	Compute Adjusted Rand Index.
<code>tl.AMI(adata, label_1, label_2)</code>	Compute adjusted Mutual Info.
<code>tl.homogeneity(adata, label_1, label_2)</code>	Compute homogeneity score.
<code>tl.silhouette(adata_name, cluster_annot[, ...])</code>	Compute silhouette scores.

episcanpy.tl.rank_features

episcanpy.tl.lazy

episcanpy.tl.load_markers

episcanpy.tl.identify_cluster

episcanpy.tl.top_feature_genes

episcanpy.tl.var_features_to_genes

episcanpy.tl.geneactivity

episcanpy.tl.diffmap

episcanpy.tl.draw_graph

episcanpy.tl.tsne

episcanpy.tl.umap

episcanpy.tl.dpt

episcanpy.tl.louvain

episcanpy.tl.leiden

episcanpy.tl.kmeans

[episcanpy.tl.hc](#)

[episcanpy.tl.getNClusters](#)

[episcanpy.tl.dendrogram](#)

[episcanpy.tl.ARI](#)

[episcanpy.tl.AMI](#)

[episcanpy.tl.homogeneity](#)

[episcanpy.tl.silhouette](#)

4.5.4 Plotting: PL

The plotting module `episcanpy.plotting` largely parallels the `tl.*` and a few of the `pp.*` functions. For most tools and for some preprocessing functions, you'll find a plotting function with the same name.

<code>pl.pca(adata, basis, *, color, ...)</code>	Scatter plot in PCA coordinates.
<code>pl.pca_overview(adata[, color, use_raw, ...])</code>	Plot PCA results.
<code>pl.pca_variance_ratio(adata[, n_pcs, log, ...])</code>	Plot the variance ratio.
<code>pl.tsne(adata, basis, *, color, ...)</code>	Scatter plot in tSNE basis.
<code>pl.umap(adata, basis, *, color, ...)</code>	Scatter plot in UMAP basis.
<code>pl.rank_feat_groups(adata[, groups, ...])</code>	Plot ranking of features.
<code>pl.rank_feat_groups_violin(adata[, groups, ...])</code>	Plot ranking of features for all tested comparisons.
<code>pl.rank_feat_groups_dotplot(adata[, groups, ...])</code>	Plot ranking of features using dotplot plot (see <code>dotplot()</code>)
<code>pl.rank_feat_groups_stacked_violin(adata[, ...])</code>	Plot ranking of features using stacked_violin plot (see <code>stacked_violin()</code>)
<code>pl.rank_feat_groups_matrixplot(adata[, ...])</code>	Plot ranking of features using matrixplot plot (see <code>matrixplot()</code>)
<code>pl.rank_feat_groups_heatmap(adata[, groups, ...])</code>	Plot ranking of features using heatmap plot (see <code>heatmap()</code>)
<code>pl.rank_feat_groups_tracksplot(adata[, ...])</code>	Plot ranking of features using heatmap plot (see <code>heatmap()</code>)
<code>pl.cal_var(adata[, show, color, save])</code>	Show distribution plots of cells sharing features and variability score.
<code>pl.violin(adata, keys[, groupby, log, ...])</code>	Violin plot.
<code>pl.scatter(adata[, x, y, color, use_raw, ...])</code>	Scatter plot along observations or variables axes.
<code>pl.ranking(adata, attr, keys[, dictionary, ...])</code>	Plot rankings.
<code>pl.clustermap(adata[, obs_keys, use_raw, ...])</code>	Hierarchically-clustered heatmap.
<code>pl.heatmap(adata, var_names[, groupby, ...])</code>	Heatmap of the expression values of genes.
<code>pl.dotplot(adata, var_names[, groupby, ...])</code>	Makes a <i>dot plot</i> of the expression values of <code>var_names</code> .
<code>pl.matrixplot(adata, var_names[, groupby, ...])</code>	Creates a heatmap of the mean expression values per cluster of each <code>var_names</code> If <code>groupby</code> is not given, the matrixplot assumes that all data belongs to a single category.

continues on next page

Table 9 – continued from previous page

<code>pl.tracksplot(adata, var_names, groupby[, ...])</code>	In this type of plot each var_name is plotted as a filled line plot where the y values correspond to the var_name values and x is each of the cells.
<code>pl.dendrogram(adata, groupby[, ...])</code>	Plots a dendrogram of the categories defined in <i>groupby</i> .
<code>pl.correlation_matrix(adata, groupby[, ...])</code>	Plots the correlation matrix computed as part of <i>sc.tl.dendrogram</i> .
<code>pl.prct_overlap(adata, key_1, key_2[, norm, ...])</code>	% or cell count corresponding to the overlap of different cell types between 2 set of annotations/clusters.
<code>pl.overlap_heatmap(adata, key_1, key_2[, ...])</code>	Heatmap of the cluster correspondance between 2 set of annotations.
<code>pl.cluster_composition(adata, cluster, condition)</code>	
<code>pl.silhouette(adata_name, cluster_annot[, ...])</code>	Plot the product of <i>tl.silhouette</i> as a silhouette plot
<code>pl.silhouette_tot(adata_name, cluster_annot)</code>	Both compute silhouette scores and plot it.
<code>pl.cal_var(adata[, show, color, save])</code>	Show distribution plots of cells sharing features and variability score.
<code>pl.variability_features(adata[, min_score, ...])</code>	This function computes a variability score to rank the most variable features across all cells.

episcanpy.pl.pca

episcanpy.pl.pca_overview

episcanpy.pl.pca_variance_ratio

episcanpy.pl.tsne

episcanpy.pl.umap

episcanpy.pl.rank_feat_groups

episcanpy.pl.rank_feat_groups_violin

episcanpy.pl.rank_feat_groups_dotplot

episcanpy.pl.rank_feat_groups_stacked_violin

episcanpy.pl.rank_feat_groups_matrixplot

episcanpy.pl.rank_feat_groups_heatmap

episcanpy.pl.rank_feat_groups_tracksplot

episcanpy.pl.cal_var

episcanpy.pl.violin

episcanpy.pl.scatter

episcanpy.pl.ranking

`episcanpy.pl.clustermap`

`episcanpy.pl.heatmap`

`episcanpy.pl.dotplot`

`episcanpy.pl.matrixplot`

`episcanpy.pl.tracksplot`

`episcanpy.pl.dendrogram`

`episcanpy.pl.correlation_matrix`

`episcanpy.pl.prct_overlap`

`episcanpy.pl.overlap_heatmap`

`episcanpy.pl.cluster_composition`

`episcanpy.pl.silhouette`

`episcanpy.pl.silhouette_tot`

`episcanpy.pl.variability_features`

4.6 References

BIBLIOGRAPHY

- [Angerer16] Angerer *et al.* (2016), *destiny – diffusion maps for large-scale single-cell data in R*, [Bioinformatics](#).
- [Cusanovich18] Cusanovich, D. A. et al. A Single-Cell Atlas of In Vivo Mammalian Chromatin Accessibility. *Cell* 174, 1309–1324.e18 (2018).
- [Luo17] Luo, C. et al. Single-cell methylomes identify neuronal subtypes and regulatory elements in mammalian cortex. *Science* 357, 600–604 (2017).
- [Wolf18] Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol.* 19, 15 (2018).
- [Buenrostro18] Buenrostro J. D. et al. Integrated Single-Cell Analysis Maps the Continuous Regulatory Landscape of Human Hematopoietic Differentiation. *Cell* 173 (2018)

PYTHON MODULE INDEX

e

episcanpy, [14](#)

INDEX

E

episcanpy
 module, 14

M

module
 episcanpy, 14