

---

# **singlet Documentation**

***Release 0.1***

**Fabio Zanini**

**Aug 27, 2017**



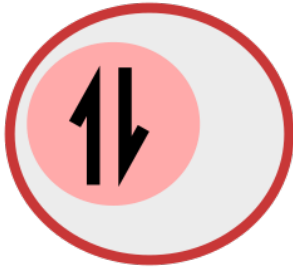
---

## Contents

---

<b>1</b>	<b>Requirements</b>	<b>3</b>
<b>2</b>	<b>Install</b>	<b>5</b>
<b>3</b>	<b>Usage example</b>	<b>7</b>
<b>4</b>	<b>Contents</b>	<b>9</b>
4.1	API reference . . . . .	9
4.2	Examples . . . . .	17
<b>5</b>	<b>Indices and tables</b>	<b>19</b>
	<b>Python Module Index</b>	<b>21</b>





Single cell RNA-Seq analysis.



# CHAPTER 1

---

## Requirements

---

Python 3.4+ is required. Moreover, you will need:

- pyyaml
- numpy
- scipy
- pandas
- scikit-learn
- matplotlib
- seaborn
- bhtsne (for t-SNE dimensionality reduction)

Get those from pip or conda.





## CHAPTER 2

---

### Install

---

TODO: upload on pypi.

For the time being, you can clone the git repo and then call:

```
python3 setup.py install
```



## CHAPTER 3

---

### Usage example

---

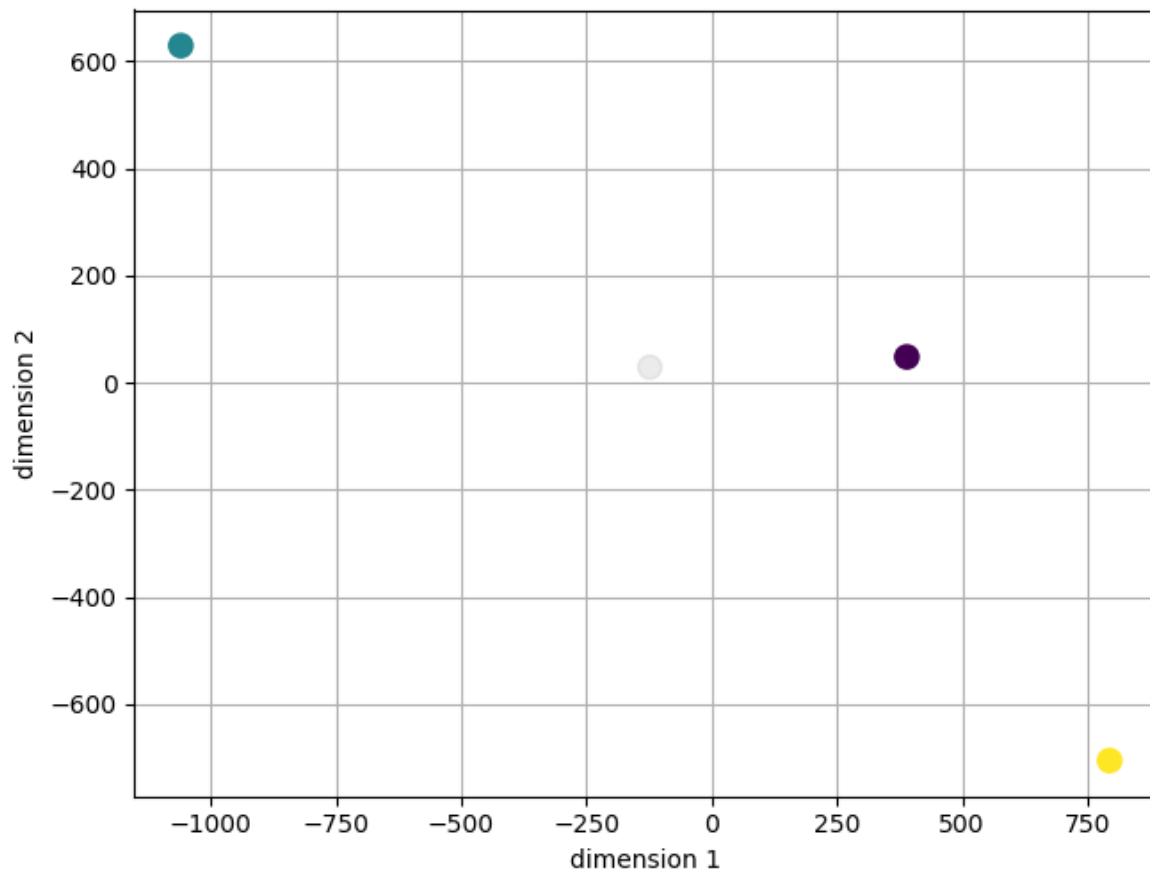
You can have a look inside the *test* folder for examples. To start using the example dataset:

- Set the environment variable *SINGLET\_CONFIG\_FILENAME* to the location of the example YAML file
- Open a Python/IPython shell and type:

```
from singlet.dataset import Dataset
ds = Dataset(
    samplesheet='example_sheet_tsv',
    counts_table='example_table_tsv')

ds.counts = ds.counts.iloc[:200]
vs = ds.dimensionality.tsne(
    n_dims=2,
    transform='log10',
    theta=0.5,
    perplexity=0.8)
ax = ds.plot.scatter_reduced_samples(
    vs,
    color_by='quantitative_phenotype_1_[A.U.]')
plt.show()
```

This will calculate a t-SNE embedding of the first 200 features and then show your samples in the reduced space. It should look like this:



### API reference

Singlet analysis is centered around the *Dataset* class, which describes a set of samples (usually single cells). Each *Dataset* has two main properties:

- a *CountsTable* with the counts of genomic features, typically transcripts
- a *SampleSheet* with the metadata and phenotypic information.

Moreover, a *Dataset* has a number of “action properties” that perform operations on the data:

- *Dataset.correlations*: correlate feature expressions and phenotypes
- *Dataset.dimensionality*: reduce dimensionality of the data including phenotypes
- *Dataset.cluster*: cluster samples, features, and phenotypes
- *Dataset.plot*: plot the results of various analyses

Supporting modules are useful for particular purposes or internal use only:

- *config*
- *utils*
- *io*

### singlet.counts\_table module

**class** singlet.counts\_table.**CountsTable** (*data=None, index=None, columns=None, dtype=None, copy=False*)

Bases: pandas.core.frame.DataFrame

Table of gene expression counts

- Rows are features, e.g. genes.
- Columns are samples.

**exclude\_features** (*spikeins=True, other=True, inplace=False, errors='raise'*)

Get a slice that excludes secondary features.

**Parameters**

- **spikeins** (*bool*) – Whether to exclude spike-ins
- **other** (*bool*) – Whether to exclude other features, e.g. unmapped reads
- **inplace** (*bool*) – Whether to drop those features in place.
- **errors** (*string*) – Whether to raise an exception if the features to be excluded are already not present.

**Returns** a slice of self without those features.

**Return type** *CountsTable*

**classmethod from\_tablename** (*tablename*)

Instantiate a CountsTable from its name in the config file.

**Parameters** **tablename** (*string*) – name of the counts table in the config file.

**Returns** the counts table.

**Return type** *CountsTable*

**get\_other\_features** ()

Get other features

**Returns** a slice of self with only other features (e.g. unmapped).

**Return type** *CountsTable*

**get\_spikeins** ()

Get spike-in features

**Returns** a slice of self with only spike-ins.

**Return type** *CountsTable*

**get\_statistics** (*metrics=('mean', 'cv')*)

Get statistics of the counts.

**Parameters** **metrics** (*sequence of strings*) – any of 'mean', 'var', 'std', 'cv', 'fano', 'min', 'max'.

**Returns** pandas.DataFrame with features as rows and metrics as columns.

**normalize** (*method='counts\_per\_million', include\_spikeins=False, inplace=False, \*\*kwargs*)

Normalize counts and return new CountsTable.

**Parameters**

- **method** (*string or function*) – The method to use for normalization. One of 'counts\_per\_million', 'counts\_per\_thousand\_spikeins', 'counts\_per\_thousand\_features'. If this argument is a function, it must take the CountsTable as input and return the normalized one as output.
- **include\_spikeins** (*bool*) – Whether to include spike-ins in the normalization and result.
- **inplace** (*bool*) – Whether to modify the CountsTable in place or return a new one.

**Returns** If *inplace* is False, a new, normalized CountsTable.

**pseudocount** = 0.1

## singlet.samplesheet module

```
class singlet.samplesheet.SampleSheet (data=None, index=None, columns=None, dtype=None,
                                       copy=False)
    Bases: pandas.core.frame.DataFrame
    classmethod from_sheetname (sheetname)
```

## singlet.dataset module

```
class singlet.dataset.Dataset (samplesheet, counts_table)
    Bases: object
    Collection of cells, with feature counts and metadata
    copy ()
        Copy of the Dataset including a new SampleSheet and CountsTable
```

**counts**  
Matrix of gene expression counts.  
Rows are features, columns are samples.

**featurenames**  
pandas.Index of feature names

**metadatanames**  
pandas.Index of metadata column names

**n\_features**  
Number of features

**n\_samples**  
Number of samples

**query\_features** (*expression*, *inplace=False*)  
Select features based on their expression.

### Parameters

- **expression** (*string*) – An expression compatible with pandas.DataFrame.query.
- **inplace** (*bool*) – Whether to change the Dataset in place or return a new one.

**Returns** If *inplace* is True, None. Else, a Dataset.

**query\_samples\_by\_counts** (*expression*, *inplace=False*)  
Select samples based on gene expression.

### Parameters

- **expression** (*string*) – An expression compatible with pandas.DataFrame.query.
- **inplace** (*bool*) – Whether to change the Dataset in place or return a new one.

**Returns** If *inplace* is True, None. Else, a Dataset.

**samplenames**  
pandas.Index of sample names

**samplesheet**  
Matrix of metadata.  
Rows are samples, columns are metadata (e.g. phenotypes).

**split** (*phenotypes*, *copy=True*)

Split Dataset based on one or more categorical phenotypes

**Parameters** **phenotypes** (*string or list of strings*) – one or more phenotypes to use for the split. Unique values of combinations of these determine the split Datasets.

**Returns** the keys are either unique values of the phenotype chosen or, if more than one, tuples of unique combinations.

**Return type** dict of Datasets

## Dataset action properties

### singlet.dataset.correlations module

**class** singlet.dataset.correlations.**Correlation** (*dataset*)

Bases: object

Correlate gene expression and phenotype in single cells

**correlate\_features\_features** (*features='all'*, *features2=None*, *method='spearman'*)

Correlate feature expression with one or more phenotypes.

#### Parameters

- **features** (*list or string*) – list of features to correlate. Use a string for a single feature. The special string 'all' (default) uses all features.
- **features** – list of features to correlate with. Use a string for a single feature. The special string 'all' uses all features. None (default) takes the same list as features, returning a square matrix.
- **method** (*string*) – type of correlation. Must be one of 'pearson' or 'spearman'.

**Returns** pandas.DataFrame with the correlation coefficients. If either features or features2 is a single string, the function returns a pandas.Series. If both are a string, it returns a single correlation coefficient.

**correlate\_features\_phenotypes** (*phenotypes*, *features='all'*, *method='spearman'*, *fillna=None*)

Correlate feature expression with one or more phenotypes.

#### Parameters

- **phenotypes** (*list of string*) – list of phenotypes, i.e. columns of the samplesheet. Use a string for a single phenotype.
- **features** (*list or string*) – list of features to correlate. Use a string for a single feature. The special string 'all' (default) uses all features.
- **method** (*string*) – type of correlation. Must be one of 'pearson' or 'spearman'.
- **fillna** (*dict, int, or None*) – a dictionary with phenotypes as keys and numbers to fill for NaNs as values. None will do nothing, potentially yielding NaN as correlation coefficients.

**Returns** pandas.DataFrame with the correlation coefficients. If either phenotypes or features is a single string, the function returns a pandas.Series. If both are a string, it returns a single correlation coefficient.



**correlate\_phenotypes\_phenotypes** (*phenotypes*, *phenotypes2=None*, *method='spearman'*, *fillna=None*, *fillna2=None*)

Correlate feature expression with one or more phenotypes.

#### Parameters

- **phenotypes** (*list of string*) – list of phenotypes, i.e. columns of the samplesheet. Use a string for a single phenotype.
- **phenotypes2** (*list of string*) – list of phenotypes, i.e. columns of the samplesheet. Use a string for a single phenotype. None (default) uses the same as phenotypes.
- **method** (*string*) – type of correlation. Must be one of 'pearson' or 'spearman'.
- **fillna** (*dict, int, or None*) – a dictionary with phenotypes as keys and numbers to fill for NaNs as values. None will do nothing, potentially yielding NaN as correlation coefficients.
- **fillna2** (*dict, int, or None*) – as fillna, but for phenotypes2.

**Returns** pandas.DataFrame with the correlation coefficients. If either phenotypes or features is a single string, the function returns a pandas.Series. If both are a string, it returns a single correlation coefficient.

## singlet.dataset.dimensionality module

**class** singlet.dataset.dimensionality.**DimensionalityReduction** (*dataset*)

Bases: object

Reduce dimensionality of gene expression and phenotype in single cells

**pca** (*n\_dims=2*, *transform='log10'*, *robust=True*, *random\_state=None*)

Principal component analysis

#### Parameters

- **n\_dims** (*int*) – Number of dimensions (2+).
- **transform** (*string or None*) – Whether to preprocess the data.
- **robust** (*bool*) – Whether to use Principal Component Pursuit to exclude outliers.

**Returns** dict of the left eigenvectors (vs), right eigenvectors (us) of the singular value decomposition, eigenvalues (lambdas), the transform, and the whiten function (for plotting).

**tsne** (*n\_dims=2*, *transform='log10'*, *perplexity=30*, *theta=0.5*, *rand\_seed=0*, *\*\*kwargs*)

t-SNE algorithm.

#### Parameters

- **n\_dims** (*int*) – Number of dimensions to use.
- **perplexity** (*float*) – Perplexity of the algorithm.
- **theta** (*float*) – A number between 0 and 1. Higher is faster but less accurate (via the Barnes-Hut approximation).
- **rand\_seed** (*int*) – Random seed. -1 randomizes each run.
- **\*\*kwargs** – Named arguments passed to the t-SNE algorithm.

Returns:

## singlet.dataset.cluster module

**class** singlet.dataset.cluster.**Cluster**(dataset)

Bases: object

Cluster samples, features, and phenotypes

**hierarchical** (axis, phenotypes=(), metric='correlation', method='average', log\_features=True, optimal\_ordering=False, \*\*kwargs)

Hierarchical clustering.

### Parameters

- **axis** (*string*) – It must be 'samples' or 'features'. The Dataset.counts matrix is used and either samples or features are clustered.
- **phenotypes** (*iterable of strings*) – Phenotypes to add to the features for joint clustering.
- **metric** (*string*) – Metric to calculate the distance matrix. Should be a string accepted by `scipy.spatial.distance.pdist`.
- **method** (*string*) – Clustering method. Must be a string accepted by `scipy.cluster.hierarchy.linkage`.
- **log\_features** (*bool*) – Whether to add pseudocounts and take a log of the feature counts before calculating distances.
- **optimal\_ordering** (*bool*) – Whether to resort the linkage so that nearest neighbours have shortest distance. This may take longer than the clustering itself.

**Returns** dict with the linkage, distance matrix, and ordering.

## singlet.dataset.plot module

**class** singlet.dataset.plot.**Plot**(dataset)

Bases: object

Plot gene expression and phenotype in single cells

**clustermap** (cluster\_samples=False, cluster\_features=False, phenotypes\_cluster\_samples=(), phenotypes\_cluster\_features=(), subtract\_mean=False, divide\_std=False, orientation='horizontal', legend=False, \*\*kwargs)

Samples versus features / phenotypes.

### Parameters

- **cluster\_samples** (*bool or linkage*) – Whether to cluster samples and show the dendrogram. Can be either, False, True, or a linkage from `scipy.cluster.hierarchy.linkage`.
- **cluster\_features** (*bool or linkage*) – Whether to cluster features and show the dendrogram. Can be either, False, True, or a linkage from `scipy.cluster.hierarchy.linkage`.
- **phenotypes\_cluster\_samples** (*iterable of strings*) – Phenotypes to add to the features for joint clustering of the samples. If the clustering has been precomputed including phenotypes and the linkage matrix is explicitly set as `cluster_samples`, the *same* phenotypes must be specified here, in the same order.
- **phenotypes\_cluster\_features** (*iterable of strings*) – Phenotypes to add to the features for joint clustering of the features and phenotypes. If the clustering

has been precomputed including phenotypes and the linkage matrix is explicitly set as `cluster_features`, the *same* phenotypes must be specified here, in the same order.

- **orientation** (*string*) – Whether the samples are on the abscissa ('horizontal') or on the ordinate ('vertical').
- **tight\_layout** (*bool or dict*) – Whether to call `matplotlib.pyplot.tight_layout` at the end of the plotting. If it is a dict, pass it unpacked to that function.
- **legend** (*bool or dict*) – If True, call `ax.legend()`. If a dict, pass as **\*\*kwargs** to `ax.legend`.
- **\*\*kwargs** – named arguments passed to the plot function.

**Returns** A seaborn ClusterGrid instance.

**gate\_features\_from\_statistics** (*features='mapped', x='mean', y='cv', \*\*kwargs*)

Select features for downstream analysis with a gate.

Usage: Click with the left mouse button to set the vertices of a polygon. Double left-click closes the shape. Right click resets the plot.

#### Parameters

- **features** (*list or string*) – List of features to plot. The string 'mapped' means everything excluding spikeins and other, 'all' means everything including spikeins and other.
- **x** (*string*) – Statistics to plot on the x axis.
- **y** (*string*) – Statistics to plot on the y axis.
- **\*\*kwargs** – named arguments passed to the plot function.

**Returns** `pd.Index` of features within the gate.

**plot\_coverage** (*features='total', kind='cumulative', ax=None, tight\_layout=True, legend=False, \*\*kwargs*)

Plot number of reads for each sample

#### Parameters

- **features** (*list or string*) – Features to sum over. The string 'total' means all features including spikeins and other, 'mapped' means all features excluding spikeins and other, 'spikeins' means only spikeins, and 'other' means only 'other' features.
- **kind** (*string*) – Kind of plot (default: cumulative distribution).
- **ax** (*matplotlib.axes.Axes*) – The axes to plot into. If None (default), a new figure with one axes is created. `ax` must not strictly be a matplotlib class, but it must have common methods such as 'plot' and 'set'.
- **tight\_layout** (*bool or dict*) – Whether to call `matplotlib.pyplot.tight_layout` at the end of the plotting. If it is a dict, pass it unpacked to that function.
- **legend** (*bool or dict*) – If True, call `ax.legend()`. If a dict, pass as **\*\*kwargs** to `ax.legend`.
- **\*\*kwargs** – named arguments passed to the plot function.

**Returns** `matplotlib.axes.Axes` with the axes containing the plot.

**plot\_distributions** (*features, kind='violin', ax=None, tight\_layout=True, legend=False, orientation='vertical', sort=False, bottom=0, grid=None, \*\*kwargs*)

Plot distribution of spike-in controls

**Parameters**

- **features** (*list or string*) – List of features to plot. If it is the string ‘spikeins’, plot all spikeins, if the string ‘other’, plot other features.
- **kind** (*string*) – Kind of plot, one of ‘violin’ (default), ‘box’, ‘swarm’.
- **ax** (*matplotlib.axes.Axes*) – Axes to plot into. If None (default), create a new figure and axes.
- **tight\_layout** (*bool or dict*) – Whether to call `matplotlib.pyplot.tight_layout` at the end of the plotting. If it is a dict, pass it unpacked to that function.
- **legend** (*bool or dict*) – If True, call `ax.legend()`. If a dict, pass as **\*\*kwargs** to `ax.legend`.
- **orientation** (*string*) – ‘horizontal’ or ‘vertical’.
- **sort** (*bool or string*) – True or ‘ascending’ sorts the features by median, ‘descending’ uses the reverse order.
- **bottom** (*float or string*) – The value of zero-count features. If you are using a log axis, you may want to set this to 0.1 or any other small positive number. If a string, it must be ‘pseudocount’, then the `CountsTable.pseudocount` will be used.
- **grid** (*bool or None*) – Whether to add a grid to the plot. None defaults to your existing settings.
- **\*\*kwargs** – named arguments passed to the plot function.

**Returns** The axes with the plot.

**Return type** `matplotlib.axes.Axes`

**scatter\_reduced\_samples** (*vectors\_reduced, color\_by=None, color\_log=None, cmap='viridis', ax=None, tight\_layout=True, legend=False, \*\*kwargs*)

Scatter samples after dimensionality reduction.

**Parameters**

- **vectors\_reduced** (*pandas.DataFrame*) – matrix of coordinates of the samples after dimensionality reduction. Rows are samples, columns (typically 2 or 3) are the component in the low-dimensional embedding.
- **color\_by** (*string or None*) – color sample dots by phenotype or expression of a certain feature.
- **color\_log** (*bool or None*) – use log of phenotype/expression in the colormap. Default None only logs expression, but not phenotypes.
- **cmap** (*string or matplotlib colormap*) – color map to use for the sample dots.
- **ax** (*matplotlib.axes.Axes*) – The axes to plot into. If None (default), a new figure with one axes is created. `ax` must not strictly be a matplotlib class, but it must have common methods such as ‘plot’ and ‘set’.
- **tight\_layout** (*bool or dict*) – Whether to call `matplotlib.pyplot.tight_layout` at the end of the plotting. If it is a dict, pass it unpacked to that function.
- **legend** (*bool or dict*) – If True, call `ax.legend()`. If a dict, pass as **\*\*kwargs** to `ax.legend`.
- **\*\*kwargs** – named arguments passed to the plot function.

**Returns** `matplotlib.axes.Axes` with the axes containing the plot.

**scatter\_statistics** (*features='mapped', x='mean', y='cv', ax=None, tight\_layout=True, legend=False, grid=None, \*\*kwargs*)

Scatter plot statistics of features.

#### Parameters

- **features** (*list or string*) – List of features to plot. The string ‘mapped’ means everything excluding spikeins and other, ‘all’ means everything including spikeins and other.
- **x** (*string*) – Statistics to plot on the x axis.
- **y** (*string*) – Statistics to plot on the y axis.
- **ax** (*matplotlib.axes.Axes*) – The axes to plot into. If None (default), a new figure with one axes is created. ax must not strictly be a matplotlib class, but it must have common methods such as ‘plot’ and ‘set’.
- **tight\_layout** (*bool or dict*) – Whether to call matplotlib.pyplot.tight\_layout at the end of the plotting. If it is a dict, pass it unpacked to that function.
- **legend** (*bool or dict*) – If True, call ax.legend(). If a dict, pass as **\*\*kwargs** to ax.legend.
- **grid** (*bool or None*) – Whether to add a grid to the plot. None defaults to your existing settings.
- **\*\*kwargs** – named arguments passed to the plot function.

**Returns** matplotlib.axes.Axes with the axes containing the plot.

## Examples

- examples/quality\_control
- examples/pca
- examples/tsne



## CHAPTER 5

---

### Indices and tables

---

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





### S

`singlet.counts_table`, [9](#)  
`singlet.dataset.cluster`, [14](#)  
`singlet.samplesheet`, [11](#)



## C

Cluster (class in singlet.dataset.cluster), 14  
 clustermap() (singlet.dataset.plot.Plot method), 14  
 copy() (singlet.dataset.Dataset method), 11  
 correlate\_features\_features() (singlet.dataset.correlations.Correlation method), 12  
 correlate\_features\_phenotypes() (singlet.dataset.correlations.Correlation method), 12  
 correlate\_phenotypes\_phenotypes() (singlet.dataset.correlations.Correlation method), 12  
 Correlation (class in singlet.dataset.correlations), 12  
 counts (singlet.dataset.Dataset attribute), 11  
 CountsTable (class in singlet.counts\_table), 9

## D

Dataset (class in singlet.dataset), 11  
 DimensionalityReduction (class in singlet.dataset.dimensionality), 13

## E

exclude\_features() (singlet.counts\_table.CountsTable method), 10

## F

featurenames (singlet.dataset.Dataset attribute), 11  
 from\_sheetname() (singlet.samplesheet.SampleSheet class method), 11  
 from\_tablename() (singlet.counts\_table.CountsTable class method), 10

## G

gate\_features\_from\_statistics() (singlet.dataset.plot.Plot method), 15  
 get\_other\_features() (singlet.counts\_table.CountsTable method), 10

get\_spikeins() (singlet.counts\_table.CountsTable method), 10  
 get\_statistics() (singlet.counts\_table.CountsTable method), 10

## H

hierarchical() (singlet.dataset.cluster.Cluster method), 14

## M

metadatanames (singlet.dataset.Dataset attribute), 11

## N

n\_features (singlet.dataset.Dataset attribute), 11  
 n\_samples (singlet.dataset.Dataset attribute), 11  
 normalize() (singlet.counts\_table.CountsTable method), 10

## P

pca() (singlet.dataset.dimensionality.DimensionalityReduction method), 13  
 Plot (class in singlet.dataset.plot), 14  
 plot\_coverage() (singlet.dataset.plot.Plot method), 15  
 plot\_distributions() (singlet.dataset.plot.Plot method), 15  
 pseudocount (singlet.counts\_table.CountsTable attribute), 10

## Q

query\_features() (singlet.dataset.Dataset method), 11  
 query\_samples\_by\_counts() (singlet.dataset.Dataset method), 11

## S

samplenames (singlet.dataset.Dataset attribute), 11  
 SampleSheet (class in singlet.samplesheet), 11  
 samplesheet (singlet.dataset.Dataset attribute), 11  
 scatter\_reduced\_samples() (singlet.dataset.plot.Plot method), 16  
 scatter\_statistics() (singlet.dataset.plot.Plot method), 17  
 singlet.counts\_table (module), 9

`singlet.dataset` (module), [11](#)  
`singlet.dataset.cluster` (module), [14](#)  
`singlet.dataset.correlations` (module), [12](#)  
`singlet.dataset.dimensionality` (module), [13](#)  
`singlet.dataset.plot` (module), [14](#)  
`singlet.samplesheet` (module), [11](#)  
`split()` (`singlet.dataset.Dataset` method), [11](#)

## T

`tsne()` (`singlet.dataset.dimensionality.DimensionalityReduction`  
method), [13](#)