

---

**topicpy**  
*Release 0.2.0*

May 19, 2022



---

## Contents

---

<b>1 Installation</b>	<b>1</b>
<b>2 Resources</b>	<b>3</b>
<b>3 Indices and tables</b>	<b>5</b>
3.1 converter . . . . .	5
3.2 geneontology . . . . .	5
3.3 gtex . . . . .	6
3.4 hsbmpy . . . . .	6
3.5 hypergeom . . . . .	7
3.6 lda . . . . .	8
3.7 tableanalyser . . . . .	13
3.8 tacos_plot . . . . .	13
3.9 TCGA_files . . . . .	13
<b>Python Module Index</b>	<b>15</b>
<b>Index</b>	<b>17</b>



# CHAPTER 1

---

## Installation

---

Install topicpy by running:

```
pip install topicpy
```

This package consists of multiple modules:

- converter: handles gene name conversions
- geneontology: uses GSEA to perform gene ontologies
- gtex: handle GTEx data
- hsbmpy: handle output of hierarchical stochastic block model
- hypergeom: perform hypergeometric tests
- lda: Run sklearn LatentDirichletAllocation with some more parameters
- tableanalyser: study distributions of a RNA-Seq data
- TCGA\_files: handle TCGA metadata



# CHAPTER 2

---

## Resources

---

- Online documentation: <https://topicpy.readthedocs.io/en/latest/>
- Paper: in preparation



# CHAPTER 3

---

## Indices and tables

---

- genindex
- modindex
- search

### 3.1 converter

```
topicpy.converter.converter
    alias of topicpy.converter.converter

topicpy.converter.convert_list_to_sybmols(ensgs: list) → list
    it converts a list of ENSG to gene names
```

**Parameters** `ensgs` – list of ENSG  
convert\_list\_to\_sybmols(["ENSG00000159763"])

```
topicpy.converter.convert_list_to_ensg(symbols: list) → list
    it converts a list of gene names to ENSG identifiers
```

**Parameters** `symbols` – list of gene symbols  
convert\_list\_to\_sybmols(["PIP"])

### 3.2 geneontology

```
topicpy.geneontology.geneontology
    alias of topicpy.geneontology.geneontology
```

---

```
topicpy.geneontology.topic_analysis(directory, l, algorithm='topsbm', background=None, verbose=True, save_Pvalues=True)
```

**Parameters**

- **directory** – where to find files
- **l** – level of the analysis
- **algorithm** – name of folder and files containing topics table (e.g. path/to/topsbm/topsbm\_level\_3\_topics.csv)
- **background** – List of background genes
- **verbose** – verbosity
- **save\_Pvalues** – save data for P-values plot

```
topicpy.geneontology.save_plot_Pvalues(df_topics, l, directory, algorithm)
```

**Parameters**

- **df\_topics** – Topics DataFrame
- **l** – level of the analysis
- **directory** –

```
topicpy.geneontology.get_ontology_df(topic, cutoff=0.05, threshhold=0.5, gene_sets=['GO_Molecular_Function_2018', 'GO_Biological_Process_2018', 'GO_Cellular_Component_2018', 'Human_Phenotype_Ontology', 'GTEx_Tissue_Sample_Gene_Expression_Profiles_up', 'GTEx_Tissue_Sample_Gene_Expression_Profiles_down', 'Tissue_Protein_Expression_from_Human_Proteome_Map', 'KEGG_2019_Human', 'NCI-60_Cancer_Cell_Lines'], background=None) → pandas.core.frame.DataFrame
```

**Parameters**

- **topic** – list of genes
- **background** – enrichment test background
- **cutoff** – Enrichments cutoff
- **threshhold** – threshold on Adjusted P-value

**Returns** DataFrame with terms and P-vals

### 3.3 gtex

```
topicpy.gtex.gtex
alias of topicpy.gtex.gtex
```

### 3.4 hsbmpy

```
class topicpy.hsbmpy.painter
```

`topicipy.hsbmpy.hsbmpy`

alias of `topicipy.hsbmpy.hsbmpy`

`topicipy.hsbmpy.clusteranalysis(directory, labels, algorithm='topsbt') → None`

Perform analyses of an algorithm output

#### Parameters

- **directory** – where to search the data
- **labels** – ground truth label to search. This should be in a file called directory/files.dat
- **algorithm** – name of the folder in which data are stored

## 3.5 hypergeom

`topicipy.hypergeom.hypergeom`

`topicipy.hypergeom.parameters_for_hypergeometric(list_1: pandas.core.series.Series, list_2: pandas.core.series.Series) → (<class 'float'>, <class 'float'>, <class 'float'>, <class 'float'>, (<class 'list'>, <class 'list'>))`

#### Parameters

- **list\_1** – series
- **list\_2** – series

lists of elements

#### Returns

- x num of successes
- M population size
- k successes in population
- N sample size
- (list\_1, list\_2) tuple of original lists

Example:

```
l1 = pd.Series(index=["ENSG00000000123", "ENSG00000000456", "ENSG00000000789", "ENSG00000000XXX"], data=["c1", "c1", "c1", "c2"], dtype=str)
l2 = pd.Series(index=["ENSG00000000123", "ENSG00000000456", "ENSG00000000789"], data=["c1", "c1", "c1"], dtype=str)
x, M, k, N, _ = parameters_for_hypergeometric(l1, l2)
```

```
>>> x
c1
c1    3
c2    0
>>> M
3
>>> k
{'c1': 3}
>>> N
{'c1': 3, 'c2': 1}
```

```
topicpy.hypergeom.build_map(num_successes, population_size, pop_successes, sample_sizes, lists,
                             last_name=None)
topicpy.hypergeom.plot_map(df_cmap, first_name='topsbt', last_name='lda', *args, **kwargs)
```

## 3.6 lda

```
class topicpy lda(learning_method='online',      max_doc_update_iter=5,      max_iter=5,
                   topic_word_prior=1, doc_topic_prior=1, random_state=42, **params)

    _approx_bound(X, doc_topic_distr, sub_sampling)
        Estimate the variational bound.
```

Estimate the variational bound over “all documents” using only the documents passed in as X. Since log-likelihood of each word cannot be computed directly, we use this bound to estimate it.

### Parameters

- **X** (*{array-like, sparse matrix} of shape (n\_samples, n\_features)*) – Document word matrix.
- **doc\_topic\_distr** (*ndarray of shape (n\_samples, n\_components)*) – Document topic distribution. In the literature, this is called gamma.
- **sub\_sampling** (*bool, default=False*) – Compensate for subsampling of documents. It is used in calculate bound in online learning.

### Returns score

**Return type** float

```
_check_feature_names(X, *, reset)
    Set or check the feature_names_in_ attribute.
```

New in version 1.0.

### Parameters

- **X** (*{ndarray, dataframe} of shape (n\_samples, n\_features)*) – The input samples.
- **reset** (*bool*) – Whether to reset the feature\_names\_in\_ attribute. If False, the input will be checked for consistency with feature names of data provided when reset was last True.  
.. note:

It is recommended to call `reset=True` in `fit` and in the first call to `partial\_fit`. All other methods that validate `X` should set `reset=False`.

```
_check_n_features(X, reset)
    Set the n_features_in_ attribute, or check against it.
```

### Parameters

- **X** (*{ndarray, sparse matrix} of shape (n\_samples, n\_features)*) – The input samples.
- **reset** (*bool*) – If True, the n\_features\_in\_ attribute is set to *X.shape[1]*. If False and the attribute exists, then check that it is equal to *X.shape[1]*. If False and the attribute does not exist, then the check is skipped. .. note:

It is recommended to call `reset=True` in `fit` and in the first call to `partial_fit`. All other methods that validate `X` should set `reset=False`.

### `_check_non_neg_array(X, reset_n_features, whom)`

check X format

check X format and make sure no negative value in X.

**Parameters** `x` (*array-like or sparse matrix*) –

### `_check_params()`

Check model parameters.

### `_e_step(X, cal_sstats, random_init, parallel=None)`

E-step in EM update.

**Parameters**

- `x` (*{array-like, sparse matrix} of shape (n\_samples, n\_features)*) – Document word matrix.
- `cal_sstats` (`bool`) – Parameter that indicate whether to calculate sufficient statistics or not. Set `cal_sstats` to `True` when we need to run M-step.
- `random_init` (`bool`) – Parameter that indicate whether to initialize document topic distribution randomly in the E-step. Set it to `True` in training steps.
- `parallel` (`joblib.Parallel, default=None`) – Pre-initialized instance of `joblib.Parallel`.

**Returns** `doc_topic_distr` is unnormalized topic distribution for each document. In the literature, this is called *gamma*. `suff_stats` is expected sufficient statistics for the M-step. When `cal_sstats == False`, it will be `None`.

**Return type** (`doc_topic_distr, suff_stats`)

### `_em_step(X, total_samples, batch_update, parallel=None)`

EM update for 1 iteration.

update `_component` by batch VB or online VB.

**Parameters**

- `x` (*{array-like, sparse matrix} of shape (n\_samples, n\_features)*) – Document word matrix.
- `total_samples` (`int`) – Total number of documents. It is only used when `batch_update` is `False`.
- `batch_update` (`bool`) – Parameter that controls updating method. `True` for batch learning, `False` for online learning.
- `parallel` (`joblib.Parallel, default=None`) – Pre-initialized instance of `joblib.Parallel`.

**Returns** `doc_topic_distr` – Unnormalized document topic distribution.

**Return type** ndarray of shape (n\_samples, n\_components)

### `classemethod _get_param_names()`

Get parameter names for the estimator

### `_init_latent_vars(n_features)`

Initialize latent variables.

**\_perplexity\_precomp\_distr** (*X, doc\_topic\_distr=None, sub\_sampling=False*)  
Calculate approximate perplexity for data *X* with ability to accept precomputed doc\_topic\_distr  
Perplexity is defined as  $\exp(-1. * \text{log-likelihood per word})$

**Parameters**

- **x** ({array-like, sparse matrix} of shape (n\_samples, n\_features)) – Document word matrix.
- **doc\_topic\_distr** (ndarray of shape (n\_samples, n\_components), default=None) – Document topic distribution. If it is None, it will be generated by applying transform on *X*.

**Returns score** – Perplexity score.

**Return type** float

**\_repr\_html\_**  
HTML representation of estimator.  
This is redundant with the logic of `_repr_mimebundle_`. The latter should be favored in the long term, `_repr_html_` is only implemented for consumers who do not interpret `_repr_mimbundle_`.

**\_repr\_html\_inner()**  
This function is returned by the @property `_repr_html_` to make `hasattr(estimator, “_repr_html_”)` return ‘True’ or ‘False’ depending on `get_config()[“display”]`.

**\_repr\_mimebundle\_(\*\*kwargs)**  
Mime bundle used by jupyter kernels to display estimator

**\_unnormalized\_transform(X)**  
Transform data *X* according to fitted model.

**Parameters x** ({array-like, sparse matrix} of shape (n\_samples, n\_features)) – Document word matrix.

**Returns doc\_topic\_distr** – Document topic distribution for *X*.

**Return type** ndarray of shape (n\_samples, n\_components)

**\_validate\_data** (*X='no\_validation'*, *y='no\_validation'*, *reset=True*, *validate\_separately=False*, *\*\*check\_params*)  
Validate input data and set or check the *n\_features\_in\_* attribute.

**Parameters**

- **x** ({array-like, sparse matrix, dataframe} of shape (n\_samples, n\_features), default='no validation') – The input samples. If ‘no\_validation’, no validation is performed on *X*. This is useful for meta-estimator which can delegate input validation to their underlying estimator(s). In that case *y* must be passed and the only accepted *check\_params* are *multi\_output* and *y\_numeric*.
- **y** (array-like of shape (n\_samples,), default='no\_validation')
  - The targets.
  - If *None*, *check\_array* is called on *X*. If the estimator’s *requires\_y* tag is True, then an error will be raised.
  - If ‘no\_validation’, *check\_array* is called on *X* and the estimator’s *requires\_y* tag is ignored. This is a default placeholder and is never meant to be explicitly set. In that case *X* must be passed.
  - Otherwise, only *y* with *check\_y* or both *X* and *y* are checked with either *check\_array* or *check\_X\_y* depending on *validate\_separately*.

- **reset** (*bool, default=True*) – Whether to reset the *n\_features\_in\_* attribute. If False, the input will be checked for consistency with data provided when reset was last True. .. note:

It is recommended to call `reset=True` in `fit` and in the first call to `partial\_fit`. All other methods that validate `X` should set `reset=False`.

- **validate\_separately** (*False or tuple of dicts, default=False*) – Only used if *y* is not None. If False, call `validate_X_y()`. Else, it must be a tuple of kwargs to be used for calling `check_array()` on *X* and *y* respectively.
- **\*\*check\_params** (*kwargs*) – Parameters passed to `sklearn.utils.check_array()` or `sklearn.utils.check_X_y()`. Ignored if `validate_separately` is not False.

**Returns** *out* – The validated input. A tuple is returned if both *X* and *y* are validated.

**Return type** {ndarray, sparse matrix} or tuple of these

### **fit** (*X, y=None*)

Learn model for the data *X* with variational Bayes method.

When *learning\_method* is ‘online’, use mini-batch update. Otherwise, use batch update.

#### Parameters

- **x** ({array-like, sparse matrix} of shape *(n\_samples, n\_features)*) – Document word matrix.
- **y** (*Ignored*) – Not used, present here for API consistency by convention.

**Returns** Fitted estimator.

**Return type** self

### **fit\_transform** (*X, y=None, \*\*fit\_params*)

Fit to data, then transform it.

Fits transformer to *X* and *y* with optional parameters *fit\_params* and returns a transformed version of *X*.

#### Parameters

- **x** (array-like of shape *(n\_samples, n\_features)*) – Input samples.
- **y** (array-like of shape *(n\_samples,)* or *(n\_samples, n\_outputs)*, *default=None*) – Target values (None for unsupervised transformations).
- **\*\*fit\_params** (*dict*) – Additional fit parameters.

**Returns** *X\_new* – Transformed array.

**Return type** ndarray array of shape *(n\_samples, n\_features\_new)*

### **full\_analysis** (*directory, xl, tl=None, label='primary\_site', logarithmise=False, round\_data=False, \*args, \*\*kwargs*) → None

#### Parameters

- **df** –
- **directory** –
- **xl** –
- **tl** –

- **kwargs** – arguments to LatentDirichletAllocation().fit\_transform

**get\_params (deep=True)**

Get parameters for this estimator.

**Parameters** **deep** (*bool, default=True*) – If True, will return the parameters for this estimator and contained subobjects that are estimators.

**Returns** **params** – Parameter names mapped to their values.

**Return type** dict

**partial\_fit (X, y=None)**

Online VB with Mini-Batch update.

**Parameters**

- **X** ({array-like, sparse matrix} of shape *(n\_samples, n\_features)*) – Document word matrix.
- **y** (*Ignored*) – Not used, present here for API consistency by convention.

**Returns** Partially fitted estimator.

**Return type** self

**perplexity (X, sub\_sampling=False)**

Calculate approximate perplexity for data X.

Perplexity is defined as  $\exp(-1. * \text{log-likelihood per word})$

Changed in version 0.19: *doc\_topic\_distr* argument has been deprecated and is ignored because user no longer has access to unnormalized distribution

**Parameters**

- **X** ({array-like, sparse matrix} of shape *(n\_samples, n\_features)*) – Document word matrix.
- **sub\_sampling** (*bool*) – Do sub-sampling or not.

**Returns** **score** – Perplexity score.

**Return type** float

**score (X, y=None)**

Calculate approximate log-likelihood as score.

**Parameters**

- **X** ({array-like, sparse matrix} of shape *(n\_samples, n\_features)*) – Document word matrix.
- **y** (*Ignored*) – Not used, present here for API consistency by convention.

**Returns** **score** – Use approximate bound as score.

**Return type** float

**set\_params (\*\*params)**

Set the parameters of this estimator.

The method works on simple estimators as well as on nested objects (such as Pipeline). The latter have parameters of the form <component>\_\_<parameter> so that it's possible to update each component of a nested object.

**Parameters** **\*\*params** (*dict*) – Estimator parameters.

**Returns** `self` – Estimator instance.

**Return type** estimator instance

**transform**(*X*)  
Transform data *X* according to the fitted model.  
Changed in version 0.18: *doc\_topic\_distr* is now normalized

**Parameters** `X` ({array-like, sparse matrix} of shape (*n\_samples*, *n\_features*)) – Document word matrix.

**Returns** `doc_topic_distr` – Document topic distribution for *X*.

**Return type** ndarray of shape (*n\_samples*, *n\_components*)

## 3.7 tableanalyser

```
topiccpy.tableanalyser.tableanalyser
alias of topiccpy.tableanalyser.tableanalyser

topiccpy.tableanalyser.scalinglawsandoverexpressed(working_dir,           normal-
                                                    isation_str='counts',
                                                    method='sampling',
                                                    how_many_sigmas=3,      dis-
                                                    distance=10)
```

## 3.8 tacos\_plot

```
topiccpy.tacos_plot.tacos_plot
alias of topiccpy.tacos_plot.tacos_plot

topiccpy.tacos_plot.scatterdense(x, y, ax=None, nbins=80, colorbar=False, c_title='density',
**kwargs) → None  
Plot a scatter plot with density
```

### Parameters

- `x` – list of first variable
- `y` – list of second variable
- `ax` – matplotlib.pyplot.Axes add here or create a new one
- `nbins` – number of bins
- `colorbar` – plot colorbar
- `c_title` – color bar title
- `**kwargs` – arguments to be passed to matplotlib.pyplot.scatter

## 3.9 TCGA\_files

```
topiccpy.TCGA_files.TCGA_files
alias of topiccpy.TCGA_files.TCGA_files
```

```
topicpy.TCGA_files.get_tcga_tissue(sample)  
    Get primary_site of tcga sample
```

**Parameters** **sample** (*str*) – sample id

```
topicpy.TCGA_files.queryFiles(files)  
    Get infor for a list of files
```

**Parameters** **files** (*list*) – list of TCGA-ids

```
topicpy.TCGA_files.queryFile(idFile)  
    Get information for file
```

**Parameters** **idFile** (*str*) – file TCGA-id

---

## Python Module Index

---

### t

topicpy.converter, 5  
topicpy.geneontology, 5  
topicpy.gtex, 6  
topicpy.hsbmpy, 6  
topicpy.hypergeom, 7  
topicpy.lda, 8  
topicpy.tableanalyser, 13  
topicpy.tacos\_plot, 13  
topicpy.TCGA\_files, 13



### Symbols

\_approx\_bound() (*topiccpy.lda.lda method*), 8  
\_check\_feature\_names() (*topiccpy.lda.lda method*), 8  
\_check\_n\_features() (*topiccpy.lda.lda method*), 8  
\_check\_non\_neg\_array() (*topiccpy.lda.lda method*), 9  
\_check\_params() (*topiccpy.lda.lda method*), 9  
\_e\_step() (*topiccpy.lda.lda method*), 9  
\_em\_step() (*topiccpy.lda.lda method*), 9  
\_get\_param\_names() (*topiccpy.lda.lda class method*), 9  
\_init\_latent\_vars() (*topiccpy.lda.lda method*), 9  
\_perplexity\_precomp\_distr() (*topiccpy.lda.lda method*), 9  
\_repr\_html\_() (*topiccpy.lda.lda attribute*), 10  
\_repr\_html\_inner() (*topiccpy.lda.lda method*), 10  
\_repr\_mimebundle\_() (*topiccpy.lda.lda method*), 10  
\_unnormalized\_transform() (*topiccpy.lda.lda method*), 10  
\_validate\_data() (*topiccpy.lda.lda method*), 10

### B

build\_map() (*in module topiccpy.hypergeom*), 7

### C

clusteranalysis() (*in module topiccpy.hsbmpy*), 7  
convert\_list\_to\_ensg() (*in module topiccpy.converter*), 5  
convert\_list\_to\_sybmols() (*in module topiccpy.converter*), 5  
converter (*in module topiccpy.converter*), 5

### F

fit() (*topiccpy.lda.lda method*), 11  
fit\_transform() (*topiccpy.lda.lda method*), 11  
full\_analysis() (*topiccpy.lda.lda method*), 11

### G

geneontology (*in module topiccpy.geneontology*), 5

get\_ontology\_df() (*in module topicypy.geneontology*), 6  
get\_params() (*topiccpy.lda.lda method*), 12  
get\_tcga\_tissue() (*in module topicypy.TCGA\_files*), 13  
gtex (*in module topicypy.gtex*), 6

### H

hsbmpy (*in module topicypy.hsbmpy*), 6  
hypergeom (*in module topicypy.hypergeom*), 7

### L

lda (*class in topicypy lda*), 8

### P

painter (*class in topicypy.hsbmpy*), 6  
parameters\_for\_hypergeometric() (*in module topicypy.hypergeom*), 7  
partial\_fit() (*topicypy.lda.lda method*), 12  
perplexity() (*topicypy.lda.lda method*), 12  
plot\_map() (*in module topicypy.hypergeom*), 8

### Q

queryFile() (*in module topicypy.TCGA\_files*), 14  
queryFiles() (*in module topicypy.TCGA\_files*), 14

### S

save\_plot\_Pvalues() (*in module topicypy.geneontology*), 6  
scalinglawsandoverexpressed() (*in module topicypy.tableanalyser*), 13  
scatterdense() (*in module topicypy.tacos\_plot*), 13  
score() (*topicypy.lda.lda method*), 12  
set\_params() (*topicypy.lda.lda method*), 12

### T

tableanalyser (*in module topicypy.tableanalyser*), 13  
tacos\_plot (*in module topicypy.tacos\_plot*), 13  
TCGA\_files (*in module topicypy.TCGA\_files*), 13

```
topic_analysis()      (in      module      top-
topicpy.geneontology), 5
topicpy.converter (module), 5
topicpy.geneontology (module), 5
topicpy.gtex (module), 6
topicpy.hsbmpy (module), 6
topicpy.hypergeom (module), 7
topicpy.lda (module), 8
topicpy.tableanalyser (module), 13
topicpy.tacos_plot (module), 13
topicpy.TCGA_files (module), 13
transform() (topicpy.lda.lda method), 13
```