

reComBat_tutorial

August 9, 2022

1 reComBat Tutorial

In this tutorial we investigate the basic functionality of the **reComBat** package and discuss best use cases. A full documentation is available [here](#).

You can install the **reComBat** via `pip install recombat`.

In this tutorial we are going to use a preprocessed version of the public *pseudomonas aeruginosa* dataset used in [recent paper](#). The preprocessing is described in our paper. If you would like to redo the preprocessing yourself, please follow the steps in `harmonizedDataCreation.py` of [this repository](#).

First, we import necessary libraries.

```
[15]: import os
import numpy as np
import pandas as pd
import scanpy as sc
import anndata as an
from reComBat import reComBat
```

Define a utility plotting function to make visualisations more convenient

```
[16]: def plot(data, metadata, type='tsne', plot_mode='all', name=None):
    adata = an.AnnData(X=data, obs=metadata)

    if plot_mode == 'all':
        to_colour_by = ['gse',
                        'strain',
                        'MediumCoarse',
                        'GrowthPhase',
                        'Oxygenation',
                        'Temperature_Coarse',
                        'Culture_Coarse',
                        'ZeroHop']
    else:
        to_colour_by = [
                        'Zero-hop cluster'
        ]
```

```

    if type == 'tsne':
        sc.tl.tsne(adata, use_rep='X')
        if name is not None:
            sc.pl.tsne(adata, color=to_colour_by, show=False, ncols=1, hspace=0.
↪25, legend_fontsize=8, save='_'+name)
        else:
            sc.pl.tsne(adata, color=to_colour_by, show=True, ncols=1, hspace=0.
↪25, legend_fontsize=8)
    elif type == 'umap':
        sc.pp.neighbors(adata, use_rep='X')
        sc.tl.umap(adata)
        if name is not None:
            sc.pl.umap(adata, color=to_colour_by, show=True, ncols=1, hspace=0.
↪25, legend_fontsize=8, save='_'+name)
        else:
            sc.pl.umap(adata, color=to_colour_by, show=True, ncols=1, hspace=0.
↪25, legend_fontsize=8)
    elif type == 'pca':
        sc.tl.pca(adata, use_highly_variable=False)
        if name is not None:
            sc.pl.pca(adata, color=to_colour_by, show=True, ncols=1, hspace=0.
↪25, legend_fontsize=8, save='_'+name)
        else:
            sc.pl.pca(adata, color=to_colour_by, show=True, ncols=1, hspace=0.
↪25, legend_fontsize=8)

```

Set the data path. In our case it is the current directory.

```
[17]: data_path = '.'
```

Load the data. This is a snapshot of the preprocessed dataset used in our publication.

```
[18]: data      = pd.read_csv(os.path.join(data_path, 'data_standardized_20211027.
↪csv'), index_col=0)
metadata = pd.read_csv(os.path.join(data_path, 'metadata_allChip_20211027.
↪csv'), index_col=0)

```

The zero hops are treated as categorical data.

```
[19]: metadata.ZeroHop = metadata.ZeroHop.astype(str)
```

Use only the coarsed information to obtain more sample overlap

```
[20]: metadata_coarse =
↪metadata[['gse', 'strain', 'MediumCoarse', 'GrowthPhase', 'Oxygenation', 'Temperature_Coarse', 'C
```

Also fill possible nan values.

```
[21]: metadata_coarse = metadata_coarse.fillna('None')
```

Only use samples which are both in the metadata and expression data. Some expression data might have been filtered out by QC.

```
[22]: valid_ids = list(set(data.index).intersection(set(metadata_coarse.index)))
      valid_ids.sort()
      data_fit = data.loc[valid_ids]
      metadata_fit = metadata_coarse.loc[valid_ids]
```

Also make sure that the samples are in correct order

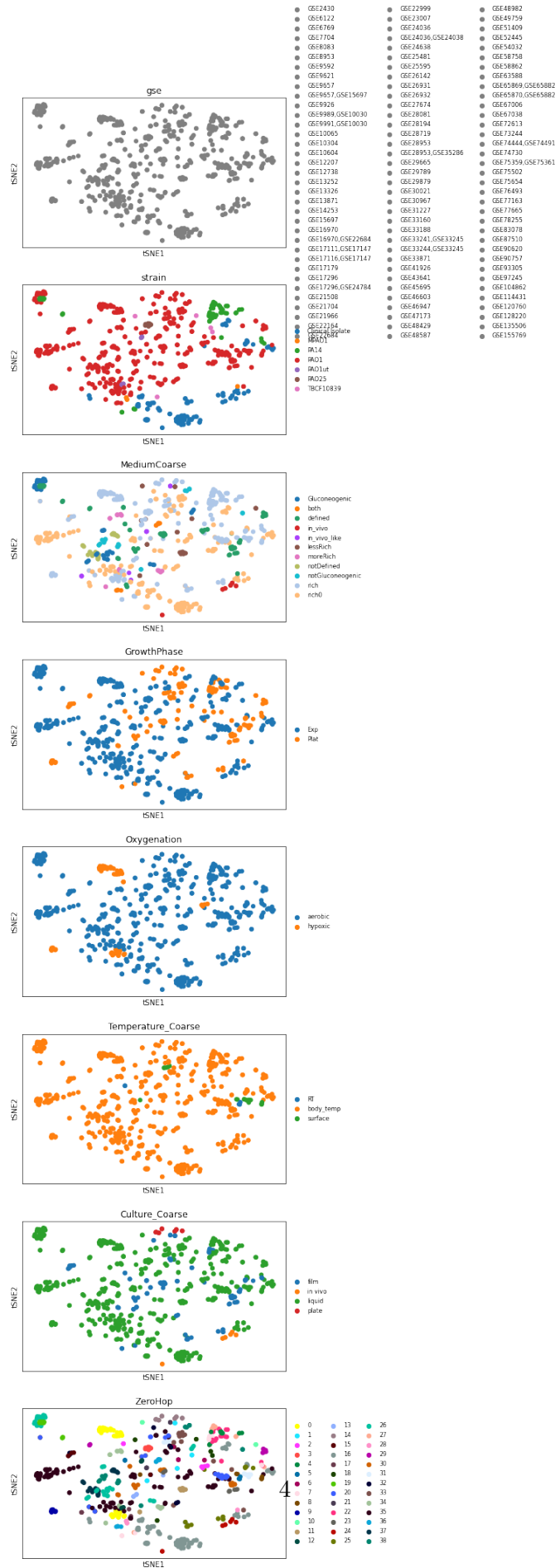
E.g. line 10 of the expression data and metadata correspond to the same sample

```
[23]: assert np.array(data_fit.index == metadata_fit.index).all()
```

Plot the initial data

```
[24]: plot(data_fit, metadata_fit, type='tsne')
```

```
... storing 'gse' as categorical
... storing 'strain' as categorical
... storing 'MediumCoarse' as categorical
... storing 'GrowthPhase' as categorical
... storing 'Oxygenation' as categorical
... storing 'Temperature_Coarse' as categorical
... storing 'Culture_Coarse' as categorical
... storing 'Antibiotic' as categorical
... storing 'ZeroHop' as categorical
```



Initialize the reComBat class

```
[25]: model = reComBat(parametric=True,          # use parametric or non-parametric
    ↪ empirical Bayes method.                    # The parametric method is
    ↪ significantly faster, whereas the          # non-parametric method is more
    ↪ flexible.                                # The regression model to be used.
    model='ridge',                            # In our experience pure ridge
    ↪ regression performs best for singular design matrices # and pure linear regression is
    ↪ best for non=singular matrices.           # Optional arguments for the
    config={'alpha':1e-12},                    # We tend to use a tiny
    ↪ regression model.                        # This has also been confirmed by
    ↪ regularisation parameter.                # The convergence criterion for
    ↪ CV.                                       # This value works well in
    conv_criterion=1e-4,                       # The maximum number of iterations
    ↪ the empirical Bayes optimisation.         # This may also be useful for
    ↪ practise.                               # This parameter is only useful in
    max_iter=1000,                             # The non-parametric optimisation
    ↪ to stop if convergence is not reached.    # Set this to the number of CPUs
    ↪ smaller convergence criteria.             # Adjust the mean of your data
    n_jobs=1,                                  # This can be useful for single
    ↪ non-parametric optimisation.              ↪ sample batches (where the variance is infinite)
    ↪ is very slow, but can be parallelised easily. # If False no empirical Bayes
    ↪ on your machine for significant speed ups. optimize_params=True,
    mean_only=False,                           # If a reference batch is present
    ↪ only (not the variance).                  # it can be set such that all data
    reference_batch=None,                       # is adjusted with respect to this reference batch.
    ↪ (e.g. a batch which is considered "batch effect free")
    ↪ is adjusted with respect to this reference batch.
    verbose=True                                # Turn log messages on or off
```

```
)
```

Fit reComBat. The input is the raw data, the batch identifiers and the design matrix

```
[26]: model.fit(data_fit,metadata_fit.gse,X=metadata_fit.  
↳drop(['gse','ZeroHop'],axis=1))
```

```
[reComBat] 2022-08-09 15:47:24,556 Starting to fot reComBat.  
[reComBat] 2022-08-09 15:47:24,576 Fit the linear model.  
[reComBat] 2022-08-09 15:47:24,782 Starting the empirical parametric  
optimisation.  
[reComBat] 2022-08-09 15:47:25,199 Optimisation finished.  
[reComBat] 2022-08-09 15:47:25,199 reComBat is fitted.
```

Not transform the data. The input stays the same

```
[27]: data_combat = model.transform(data_fit,metadata_fit.gse,X=metadata_fit.  
↳drop(['gse','ZeroHop'],axis=1))
```

```
[reComBat] 2022-08-09 15:47:25,211 Starting to transform.  
[reComBat] 2022-08-09 15:47:25,416 Transform finished.
```

Plot the output data for comparison

```
[28]: plot(data_combat,metadata_fit,type='tsne')
```

```
... storing 'gse' as categorical  
... storing 'strain' as categorical  
... storing 'MediumCoarse' as categorical  
... storing 'GrowthPhase' as categorical  
... storing 'Oxygenation' as categorical  
... storing 'Temperature_Coarse' as categorical  
... storing 'Culture_Coarse' as categorical  
... storing 'Antibiotic' as categorical  
... storing 'ZeroHop' as categorical
```

