

# ReFACTor v1.0

Reference-Free Adjustment for Cell-Type composition (ReFACTor) is an unsupervised method for the correction of cell-type heterogeneity in epigenome-wide association studies (EWAS), which is based on a variant of principal component analysis (PCA). ReFACTor is described in the following [paper](#).

**An updated implementation of ReFACTor is now available in [GLINT](#), a user-friendly command-line tool for fast analysis of genome-wide DNA methylation data.**

ReFACTor is also available in standalone implementations in both R and Python (described in details below).

## Download

1. Download the latest release from [here](#).
2. To use the Python version:
  - i. Install [Anaconda Python version 2.7](#) (automatically includes all required dependencies).
  - ii. Install ReFACTor by running the provided install.py file (run `cd python; python install.py`). For more details see "Dependencies (Python version)".
3. Make sure the provided demo works (see below).

## Running ReFACTor

### Python

Execute the refactor.py script found in the 'python' folder from the command line as follows:

```
python refactor.py --datafile <datafile> --k <k>
```

Additional optional arguments can be included (for more details see "Input arguments"):

```
python refactor.py --datafile <datafile> --k <k> --covarfile [covarfile] --t [t] --numcomp [numcomp] --stdth [stdth] --out [out]
```

### Demo

The following demo computes the ReFACTor components of a simulated example dataset and performs an EWAS. The demo shows that while a standard PCA cannot adjust the data well, using ReFACTor can adjust the data similarly to using the true cell proportions. From the command line run:

```
python demo.py
```

### R

The refactor.R function in the "R" folder can be executed directly from R. For example:

```
# <R code>
source("refactor.R")
k = 5
datafile = "../demo_files/demo_datafile.txt"
results <- refactor(datafile,k)
RC <- results$refactor_components # Extract the ReFACTor components
ranked_list <- results$ranked_list # Extract the list of sites ranked by ReFACTor

# Can also provide one or more of the optional arguments
results <- refactor(datafile,k,covarfile="../demo_files/demo_covariates.txt",t=500,numcomp=10,stdth=0.01,out="demo_results")
```

## Demo

The following demo computes the ReFACTor components of a simulated example dataset and performs an EWAS. The demo shows that while a standard PCA cannot adjust the data well, using ReFACTor can adjust the data similarly to using the true cell proportions. From the command line run:

```
Rscript demo.R
```

## Input arguments

ReFACTor gets the following arguments as an input:

Required:

- **datafile** - path to a sites by samples matrix file of tab-delimited beta-normalized methylation levels; the first row should contain the sample IDs and the first column should contain the CpG IDs (see "demo\_files/demo\_datafile.txt" for example). Important data preparation instructions are described below under "Data preparation".
- **k** - the number of assumed cell types; guidelines for selecting k are described below under "Parameters selection".

Optional:

- **covarfile** - path to a samples by covariates matrix file of tab-delimited covariates; the first column should contain the sample IDs ordered as in the first row of the data file (see "demo\_files/demo\_covariates.txt" for example). If provided, the data are adjusted for the covariates before running ReFACTor. For more details see "Data preparation".
- **t** - the number of sites to use for computing the ReFACTor components (default is 500); guidelines for selecting t are described below under "Parameters selection".
- **numcomp** - the number of ReFACTor components to output (default is the same as k).
- **stdth** - standard deviation (std) threshold for excluding low variance sites; all sites with  $\text{std} < \text{stdth}$  will be excluded before running ReFACTor (default is 0.02). For more details see "Data preparation".
- **out** - prefix of the output files (default is "refactor").

## Output

The software outputs two files:

1. refactor.out.components.txt - a matrix with the first numcomp ReFACTor components for each individual
2. refactor.out.rankedlist.txt - a ranked list of the methylation sites; from the most informative to the least informative

Note that the default prefix of these files ("refactor") can be changed using the "out" argument.

## Data preparation

### Preprocessing raw data

ReFACTor is designed to handle Beta normalized methylation levels (although it may perform well on M-value normalized data as well). Prior to running ReFACTor, the data should be adjusted for known technical artifacts of the technology used for probing the methylation levels as well as adjusted for known technical covariates (such as batches). For a comprehensive comparison between methods for preprocessing raw data collected by the Illumina 27K/450K platforms see [Lenhe et al. \(2015\)](#). In order to best fit to the assumptions of ReFACTor, any normalization applied should keep the data approximately normal.

### Preparing data for ReFACTor

For best results, we suggest to take the following steps when preparing the data for ReFACTor:

- **Exclude problematic probes** - remove non-autosomal probes, cross-hybridized probes and probes with SNPs. Note that once the ReFACTor components are computed, any of the excluded probes can be rejoined to the data for the rest of the analysis.
- **Exclude outlier samples** - outliers can be revealed using dimensionality reduction methods such as PCA or multidimensional scaling (MDS).
- **Adjust the data for covariates** ("covarfile" argument) - adjusting the methylation levels, before running ReFACTor, for known technical covariates such as batch information can be crucial in some cases. In addition, we observe that adjusting the methylation levels for genome-wide affecting factors, such as gender, smoking status and global ancestry, improves the performance of

ReFACTor. However, we do not suggest to adjust the data for covariates that are correlated with the cell type composition, such as age, before running ReFACTor (these covariates should be accounted for after running ReFACTor). The "covarfile" optional argument allows to adjust the data for covariates before running ReFACTor.

- **Remove sites with very low variance** ("stdth" argument) - Many sites in the Illumina 27K/450K platforms are constant or nearly-constant. We observe that removing those sites improves the performance of ReFACTor (defined by the "stdth" argument; the default value should be sufficient in most cases).
- **Handle missing values** - The current version of ReFACTor does not handle missing values. If missing values exist in the data they should be assigned with values before running ReFACTor (e.g. for each site its missings values can be assigned with the mean value of the site - across all samples with no missing values).

## Parameters selection

The manuscript describing ReFACTor demonstrates that the algorithm is robust to the selection of the parameters  $k$  and  $t$  in simulated and real data. However, sometimes even an approximation of  $k$  is not available, and the default value of  $t$  ( $t=500$ ) may not be adequate in some cases. Therefore, we provide the following tools for guiding the parameters selection. These tools are available in the Python version only.

### Selecting $k$ (the number of assumed cell types)

The `estimate_k.py` script (under the "python" folder) computes a score for each of the first several eigenvalues of the empirical covariance matrix of the input data. The score of the  $i$ -th eigenvalue is defined to be  $-\log$  of the ratio between the  $i$ -th eigenvalue to the  $(i-1)$ -th eigenvalue, thus a high score of a specific eigenvalue suggests its eigenvector as a substantial variance component in the data (compared with the previous one). The ratio between adjacent eigenvalues, as well as other test statistics of the eigenvalues, is described by [Peres-Neto et al. \(2004\)](#) as a method for determining the number of non-trivial axes of variance in data.

For plotting the scores of the first several eigenvalues (starting from the second eigenvalue), run:

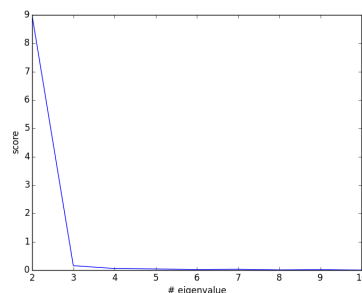
```
python estimate_k.py --datafile <datafile>
```

The maximal number of eigenvalues in the plot can be changed:

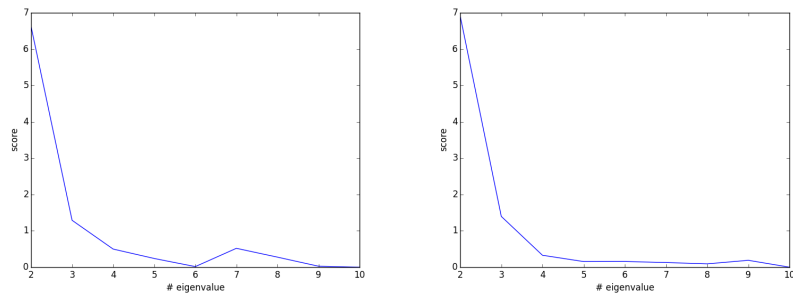
```
python estimate_k.py --datafile <datafile> --max_k <max_k>
```

$k$  should be selected to be the number of high score eigenvalues, before reaching a right tail of flat scores (the scores of the last several eigenvalues). Below are examples of plots generated by `estimate_k.py`:

- `estimate_k_results_simulated.png` - the result of applying the script on the example dataset provided here under the "demo\_files" folder. The plot suggests 4 or 5 to be a reasonable choice of  $k$ . These simulated data were in fact generated using  $k=5$ .



- `estimate_k_results_ra.png`, `estimate_k_results_gala.png` - the result of applying the script on the Rheumatoid arthritis and GALA II datasets presented in the manuscript describing ReFACTor. Both plots suggest 5 or 6 to be a reasonable choice of  $k$ .



### Selecting $t$ (the number of sites to use for computing ReFACTOR's components)

The `estimate_t.py` script (under the "python" folder) provides a qualitative tool for assessing the number of features in the data that are highly informative in terms of the main structure in the data. The script first follows the ReFACTOR algorithm in order to find the distance of each site from its  $k$ -rank approximation. Then, the sites are sorted by their distance, and a score for site  $i$  in the sorted list is defined to be the difference between the distance of the  $i$ -th site and the distance of the  $(i-1)$ -th site. Finally, the scores of the first several thousands of sites are plotted (using moving average for smoothing the signal), thus providing a qualitative way to get the number of highly informative sites.

Execution:

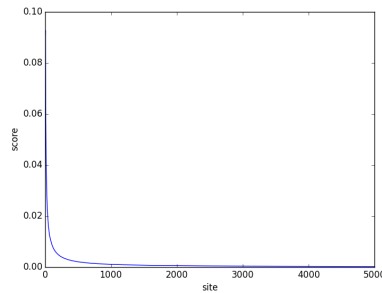
```
python estimate_t.py --datafile <datafile> --k <k>
```

The number of sites in the plot can be changed:

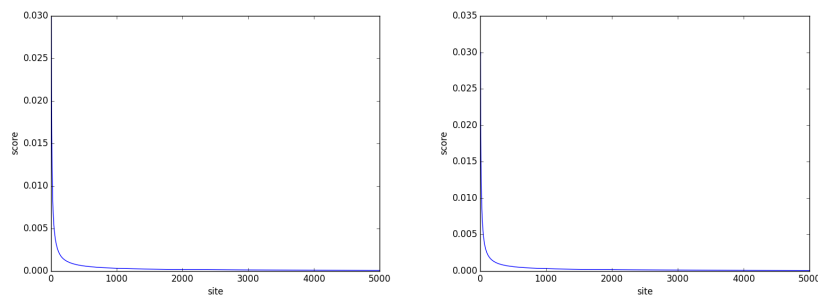
```
python estimate_t.py --datafile <datafile> --k <k> --numsites <num_sites>
```

$t$  should be selected to be the number of sites after which the signal dramatically decays. Below are examples of plots generated by `estimate_t.py`:

- `estimate_t_results_simulated.png` - the result of applying the script on the example dataset provided here under the "demo\_files" folder. The plot suggests the range between 500 and 1000 to be a reasonable choice of  $t$ .



- `estimate_t_results_ra.png`, `estimate_t_results_gala.png` - the result of applying the script on the Rheumatoid arthritis and GALA II datasets presented in the manuscript describing ReFACTOR. Both plots suggest the range between 500 and 1000 to be a reasonable choice of  $t$ .



### Dependencies (Python version)

This release of ReFACTor was implemented for Python 2.7 and has the following dependencies:

```
numpy
scipy
sklearn
matplotlib (required only for demo.py)
statsmodels (required only for demo.py)
```

We recommend installing [Anaconda Python version 2.7](#), which already includes all necessary dependencies.

If you already have Python installed and do not want to install Anaconda Python, run "install.py" script (found in the "python" folder):

```
python install.py
```

The script automatically installs missing dependencies that are required for ReFACTor. It also adds "refactor.py" script to the path of your operating system (not implemented on Windows). Note that in some environments the script may fail to install some of the dependencies, in which case you will need to manually install them.

## Citing ReFACTor

If you use ReFACTor in any published work, please cite the manuscript describing the method:

Elior Rahmani, Noah Zaitlen, Yael Baran, Celeste Eng, Donglei Hu, Joshua Galanter, Sam Oh, Esteban G Burchard, Eleazar Eskin, James Zou and Eran Halperin. "Sparse PCA corrects for cell type heterogeneity in epigenome-wide association studies". Nature Methods (2016).

Remark: The data contributed to GEO ([GSE77716](#)) is the full GALA 2 dataset. In the [ReFACTor paper](#) we used only a subset of the data that was available for us at the time of the analysis. As described in the paper, Figure 2 was generated using 78 samples (for which we had cell counts at the time of the analysis). The identifiers of these individuals can be found [here](#).

## Authors

This software was developed by Reut Yedidim, Noah Zaitlen and Elior Rahmani.

For any question and for reporting bugs please send an email to Elior Rahmani at: [elior.rahmani@gmail.com](mailto:elior.rahmani@gmail.com)