

# Package ‘harmony’

June 2, 2021

**Title** Fast, Sensitive, and Accurate Integration of Single Cell Data

**Version** 0.1.0

**Description** Implementation of the Harmony algorithm for single cell integration, described in Korsunsky et al <[doi:10.1038/s41592-019-0619-0](https://doi.org/10.1038/s41592-019-0619-0)>. Package includes a standalone Harmony function and interfaces to external frameworks.

**URL** [software.broadinstitute.org/harmony](https://software.broadinstitute.org/harmony)

**License** GPL-3

**Encoding** UTF-8

**RoxygenNote** 7.1.1

**Depends** R(>= 3.5.0), Rcpp

**LazyData** true

**LinkingTo** Rcpp, RcppArmadillo, RcppProgress

**Imports** dplyr, cowplot, tidyr, ggplot2, irlba, Matrix, methods, tibble, rlang

**Suggests** Seurat, testthat, knitr, rmarkdown

**VignetteBuilder** knitr

**NeedsCompilation** yes

**Author** Ilya Korsunsky [cre, aut] (<<https://orcid.org/0000-0003-4848-3948>>),  
Nghia Millard [aut] (<<https://orcid.org/0000-0002-0518-7674>>),  
Jean Fan [aut, ctb] (<<https://orcid.org/0000-0002-0212-5451>>),  
Kamil Slowikowski [aut, ctb] (<<https://orcid.org/0000-0002-2843-6370>>),  
Miles Smith [ctb],  
Soumya Raychaudhuri [aut] (<<https://orcid.org/0000-0002-1901-8265>>)

**Maintainer** Ilya Korsunsky <[ilya.korsunsky@gmail.com](mailto:ilya.korsunsky@gmail.com)>

**Repository** CRAN

**Date/Publication** 2021-06-02 08:10:02 UTC

**R topics documented:**

cell_lines . . . . .	2
cell_lines_small . . . . .	2
harmony . . . . .	3
HarmonyMatrix . . . . .	3
moe_ridge_get_betas . . . . .	5
RunHarmony . . . . .	6

<b>Index</b>	<b>9</b>
--------------	----------

---

cell_lines	<i>List of metadata table and scaled PCs matrix</i>
------------	---

---

**Description**

List of metadata table and scaled PCs matrix

**Usage**

cell\_lines

**Format**

: meta\_data: data.table of 9478 rows with defining dataset and cell\_type scaled\_pcs: data.table of 9478 rows (cells) and 20 columns (PCs)

**Source**

<https://support.10xgenomics.com/single-cell-gene-expression/datasets>

---

cell_lines_small	<i>Same as cell_lines but smaller (300 cells).</i>
------------------	--

---

**Description**

Same as cell\_lines but smaller (300 cells).

**Usage**

cell\_lines\_small

**Format**

An object of class list of length 2.

**Source**

<https://support.10xgenomics.com/single-cell-gene-expression/datasets>

---

harmony	<i>Harmony: fast, accurate, and robust single cell integration.</i>
---------	---

---

**Description**

Algorithm for single cell integration.

**Usage**

1. ?HarmonyMatrix to run Harmony on gene expression or PCA embeddings matrix.
2. ?RunHarmony to run Harmony on Seurat objects.

**Useful links**

1. Report bugs at <https://github.com/immunogenomics/harmony/issues>
2. Read the manuscript [online](#).

---

HarmonyMatrix	<i>Main Harmony interface</i>
---------------	-------------------------------

---

**Description**

Use this to run the Harmony algorithm on gene expression or PCA matrix.

**Usage**

```
HarmonyMatrix(  
  data_mat,  
  meta_data,  
  vars_use,  
  do_pca = TRUE,  
  npcs = 20,  
  theta = NULL,  
  lambda = NULL,  
  sigma = 0.1,  
  nclust = NULL,  
  tau = 0,  
  block.size = 0.05,  
  max.iter.harmony = 10,  
  max.iter.cluster = 200,  
  epsilon.cluster = 1e-05,  
  epsilon.harmony = 1e-04,  
  plot_convergence = FALSE,  
  return_object = FALSE,  
  verbose = TRUE,
```

```

    reference_values = NULL,
    cluster_prior = NULL
)

```

### Arguments

<code>data_mat</code>	Matrix of normalized gene expression (default) or PCA embeddings (see <code>do_pca</code> ). Cells can be rows or columns.
<code>meta_data</code>	Either (1) Dataframe with variables to integrate or (2) vector with labels.
<code>vars_use</code>	If <code>meta_data</code> is dataframe, this defined which variable(s) to remove (character vector).
<code>do_pca</code>	Whether to perform PCA on input matrix.
<code>npcs</code>	If doing PCA on input matrix, number of PCs to compute.
<code>theta</code>	Diversity clustering penalty parameter. Specify for each variable in <code>vars_use</code> . Default <code>theta=2</code> . <code>theta=0</code> does not encourage any diversity. Larger values of <code>theta</code> result in more diverse clusters.
<code>lambda</code>	Ridge regression penalty parameter. Specify for each variable in <code>vars_use</code> . Default <code>lambda=1</code> . <code>lambda</code> must be strictly positive. Smaller values result in more aggressive correction.
<code>sigma</code>	Width of soft kmeans clusters. Default <code>sigma=0.1</code> . <code>sigma</code> scales the distance from a cell to cluster centroids. Larger values of <code>sigma</code> result in cells assigned to more clusters. Smaller values of <code>sigma</code> make soft kmeans cluster approach hard clustering.
<code>nclust</code>	Number of clusters in model. <code>nclust=1</code> equivalent to simple linear regression.
<code>tau</code>	Protection against overclustering small datasets with large ones. <code>tau</code> is the expected number of cells per cluster.
<code>block.size</code>	What proportion of cells to update during clustering. Between 0 to 1, default 0.05. Larger values may be faster but less accurate
<code>max.iter.harmony</code>	Maximum number of rounds to run Harmony. One round of Harmony involves one clustering and one correction step.
<code>max.iter.cluster</code>	Maximum number of rounds to run clustering at each round of Harmony.
<code>epsilon.cluster</code>	Convergence tolerance for clustering round of Harmony. Set to <code>-Inf</code> to never stop early.
<code>epsilon.harmony</code>	Convergence tolerance for Harmony. Set to <code>-Inf</code> to never stop early.
<code>plot_convergence</code>	Whether to print the convergence plot of the clustering objective function. <code>TRUE</code> to plot, <code>FALSE</code> to suppress. This can be useful for debugging.
<code>return_object</code>	(Advanced Usage) Whether to return the Harmony object or only the corrected PCA embeddings.
<code>verbose</code>	Whether to print progress messages. <code>TRUE</code> to print, <code>FALSE</code> to suppress.

reference\_values (Advanced Usage) Defines reference dataset(s). Cells that have batch variables values matching reference\_values will not be moved.

cluster\_prior (Advanced Usage) Provides user defined clusters for cluster initialization. If the number of provided clusters C is less than K, Harmony will initialize K-C clusters with kmeans. C cannot exceed K.

### Value

By default, matrix with corrected PCA embeddings. If return\_object is TRUE, returns the full Harmony object (R6 reference class type).

### Examples

```
## By default, Harmony inputs a normalized gene expression matrix
## Not run:
harmony_embeddings <- HarmonyMatrix(exprs_matrix, meta_data, 'dataset')

## End(Not run)

## Harmony can also take a PCA embeddings matrix
data(cell_lines_small)
pca_matrix <- cell_lines_small$scaled_pcs
meta_data <- cell_lines_small$meta_data
harmony_embeddings <- HarmonyMatrix(pca_matrix, meta_data, 'dataset',
                                   do_pca=FALSE)

## Output is a matrix of corrected PC embeddings
dim(harmony_embeddings)
harmony_embeddings[seq_len(5), seq_len(5)]

## Finally, we can return an object with all the underlying data structures
harmony_object <- HarmonyMatrix(pca_matrix, meta_data, 'dataset',
                                do_pca=FALSE, return_object=TRUE)
dim(harmony_object$Y) ## cluster centroids
dim(harmony_object$R) ## soft cluster assignment
dim(harmony_object$Z_corr) ## corrected PCA embeddings
head(harmony_object$O) ## batch by cluster co-occurrence matrix
```

---

moe\_ridge\_get\_betas    *Get beta Utility*

---

### Description

Utility function to get ridge regression coefficients from trained Harmony object

**Usage**

```
moe_ridge_get_betas(harmonyObj)
```

**Arguments**

harmonyObj      Trained harmony object. Get this by running HarmonyMatrix function with return\_object=TRUE.

**Value**

Returns nothing, modifies object in place.

---

RunHarmony

*Harmony single cell integration*

---

**Description**

Run Harmony algorithm with Seurat and SingleCellAnalysis pipelines.

**Usage**

```
RunHarmony(object, group.by.vars, ...)
```

```
## S3 method for class 'Seurat'  
RunHarmony(  
  object,  
  group.by.vars,  
  reduction = "pca",  
  dims.use = NULL,  
  theta = NULL,  
  lambda = NULL,  
  sigma = 0.1,  
  nclust = NULL,  
  tau = 0,  
  block.size = 0.05,  
  max.iter.harmony = 10,  
  max.iter.cluster = 20,  
  epsilon.cluster = 1e-05,  
  epsilon.harmony = 1e-04,  
  plot_convergence = FALSE,  
  verbose = TRUE,  
  reference_values = NULL,  
  reduction.save = "harmony",  
  assay.use = "RNA",  
  project.dim = TRUE,  
  ...  
)
```

**Arguments**

<code>object</code>	Pipeline object. Must have PCA computed.
<code>group.by.vars</code>	Which variable(s) to remove (character vector).
<code>...</code>	other parameters
<code>reduction</code>	Name of dimension reduction to use. Default is PCA.
<code>dims.use</code>	Which PCA dimensions to use for Harmony. By default, use all
<code>theta</code>	Diversity clustering penalty parameter. Specify for each variable in <code>group.by.vars</code> . Default <code>theta=2</code> . <code>theta=0</code> does not encourage any diversity. Larger values of <code>theta</code> result in more diverse clusters.
<code>lambda</code>	Ridge regression penalty parameter. Specify for each variable in <code>group.by.vars</code> . Default <code>lambda=1</code> . <code>lambda</code> must be strictly positive. Smaller values result in more aggressive correction.
<code>sigma</code>	Width of soft kmeans clusters. Default <code>sigma=0.1</code> . <code>sigma</code> scales the distance from a cell to cluster centroids. Larger values of <code>sigma</code> result in cells assigned to more clusters. Smaller values of <code>sigma</code> make soft kmeans cluster approach hard clustering.
<code>nclust</code>	Number of clusters in model. <code>nclust=1</code> equivalent to simple linear regression.
<code>tau</code>	Protection against overclustering small datasets with large ones. <code>tau</code> is the expected number of cells per cluster.
<code>block.size</code>	What proportion of cells to update during clustering. Between 0 to 1, default 0.05. Larger values may be faster but less accurate
<code>max.iter.harmony</code>	Maximum number of rounds to run Harmony. One round of Harmony involves one clustering and one correction step.
<code>max.iter.cluster</code>	Maximum number of rounds to run clustering at each round of Harmony.
<code>epsilon.cluster</code>	Convergence tolerance for clustering round of Harmony Set to <code>-Inf</code> to never stop early.
<code>epsilon.harmony</code>	Convergence tolerance for Harmony. Set to <code>-Inf</code> to never stop early.
<code>plot_convergence</code>	Whether to print the convergence plot of the clustering objective function. <code>TRUE</code> to plot, <code>FALSE</code> to suppress. This can be useful for debugging.
<code>verbose</code>	Whether to print progress messages. <code>TRUE</code> to print, <code>FALSE</code> to suppress.
<code>reference_values</code>	(Advanced Usage) Defines reference dataset(s). Cells that have batch variables values matching <code>reference_values</code> will not be moved
<code>reduction.save</code>	Keyword to save Harmony reduction. Useful if you want to try Harmony with multiple parameters and save them as e.g. <code>'harmony_theta0'</code> , <code>'harmony_theta1'</code> , <code>'harmony_theta2'</code>
<code>assay.use</code>	(Seurat V3 only) Which assay to Harmonize with (RNA by default).
<code>project.dim</code>	Project dimension reduction loadings. Default <code>TRUE</code> .

**Value**

Seurat (version 3) object. Harmony dimensions placed into dimensional reduction object harmony. For downstream Seurat analyses, use reduction='harmony'.



# Index

## \* datasets

cell\_lines, [2](#)

cell\_lines\_small, [2](#)

cell\_lines, [2](#)

cell\_lines\_small, [2](#)

harmony, [3](#)

HarmonyMatrix, [3](#)

moe\_ridge\_get\_betas, [5](#)

RunHarmony, [6](#)