

# Package ‘bbknnR’

June 5, 2025

**Title** Perform Batch Balanced KNN in R

**Version** 2.0.1

**Date** 2025-06-05

**Description** A fast and intuitive batch effect removal tool for single-cell data. BBKNN is originally used in the 'scanpy' python package, and now can be used with 'Seurat' seamlessly.

**License** MIT + file LICENSE

**Encoding** UTF-8

**Depends** R (>= 4.1.0), methods, utils

**LinkingTo** Rcpp (>= 1.0.8), RcppEigen

**Imports** future, future.apply, glmnet, Rcpp, RcppAnnoy, RcppEigen, rlang, rndescent, Rtsne, Seurat, SeuratObject, tidytable, uwot (>= 0.2.1)

**LazyData** true

**RoxygenNote** 7.3.1

**URL** <https://github.com/ycli1995/bbknnR>,  
<https://github.com/Teichlab/bbknn>,  
<https://bbknn.readthedocs.io/en/latest/>

**BugReports** <https://github.com/ycli1995/bbknnR/issues>

**Suggests** dplyr, knitr, rmarkdown, testthat (>= 3.0.0), patchwork

**Config/testthat/edition** 3

**VignetteBuilder** knitr

**NeedsCompilation** yes

**Author** Yuchen Li [aut, cre]

**Maintainer** Yuchen Li <ycli1995@outlook.com>

**Repository** CRAN

**Date/Publication** 2025-06-05 12:10:02 UTC

## Contents

panc8_small . . . . .	2
RidgeRegression . . . . .	3
RunBBKNN . . . . .	4

<b>Index</b>	<b>8</b>
--------------	----------

---

panc8_small	<i>A small example version of the pancreas scRNA-seq dataset</i>
-------------	--

---

## Description

A subsetted version of the pancreas scRNA-seq dataset to test BBKNN

## Usage

```
panc8_small
```

## Format

A Seurat object with the following slots filled

**assays** Currently only contains one assay ("RNA" - scRNA-seq expression data)

counts - Raw expression data

- data - Normalized expression data
- scale.data - Scaled expression data
- var.features - names of the current features selected as variable
- meta.features - Assay level metadata such as mean and variance

**meta.data** Cell level metadata

**active.assay** Current default assay

**active.ident** Current default ident

**graphs** Empty

**reductions** Dimensional reductions: currently PCA

**version** Seurat version used to create the object

**commands** Command history

## Source

SeuratData <https://github.com/satijalab/seurat-data>

---

RidgeRegression	<i>Perform ridge regression on scaled expression data</i>
-----------------	---

---

**Description**

Perform ridge regression on scaled expression data, accepting both technical and biological categorical variables. The effect of the technical variables is removed while the effect of the biological variables is retained. This is a preprocessing step that can aid BBKNN integration.

**Usage**

```
RidgeRegression(object, ...)  
  
## Default S3 method:  
RidgeRegression(  
  object,  
  latent_data,  
  batch_key,  
  confounder_key,  
  lambda = 1,  
  seed = 42,  
  verbose = TRUE,  
  ...  
)  
  
## S3 method for class 'Seurat'  
RidgeRegression(  
  object,  
  batch_key,  
  confounder_key,  
  assay = NULL,  
  features = NULL,  
  lambda = 1,  
  run_pca = TRUE,  
  npcs = 50,  
  reduction.name = "pca",  
  reduction.key = "PC_",  
  replace = FALSE,  
  seed = 42,  
  verbose = TRUE,  
  ...  
)
```

**Arguments**

object	An object
...	Arguments passed to other methods

latent_data	Extra data to regress out, should be cells x latent data
batch_key	Variables to regress out as technical effects. Must be included in column names of latent_data
confounder_key	Variables to retain as biological effects. Must be included in column names of latent_data
lambda	A user supplied lambda sequence. pass to <a href="#">glmnet</a>
seed	Set a random seed. By default, sets the seed to 42. Setting NULL will not set a seed.
verbose	Whether or not to print output to the console
assay	Name of Assay ridge regression is being run on
features	Features to compute ridge regression on. If features=NULL, ridge regression will be run using the variable features for the Assay.
run_pca	Whether or not to run pca with regressed expression data (TRUE by default)
npcs	Total Number of PCs to compute and store (50 by default)
reduction.name	Dimensional reduction name (pca by default)
reduction.key	Dimensional reduction key, specifies the string before the number for the dimension names (PC by default)
replace	Whether or not to replace original scale.data with regressed expression data (TRUE by default)

### Value

Returns a Seurat object.

### References

Park, Jong-Eun, et al. "A cell atlas of human thymic development defines T cell repertoire formation." *Science* 367.6480 (2020): eaay3224.

### Examples

```
data("panc8_small")
panc8_small <- RidgeRegression(panc8_small, "tech", c("nCount_RNA"))
```

---

RunBBKNN

*Perform batch balanced KNN*

---

### Description

Batch balanced KNN, altering the KNN procedure to identify each cell's top neighbours in each batch separately instead of the entire cell pool with no accounting for batch. The nearest neighbours for each batch are then merged to create a final list of neighbours for the cell. Aligns batches in a quick and lightweight manner.

**Usage**

```
RunBBKNN(object, ...)

## S3 method for class 'matrix'
RunBBKNN(
  object,
  batch_list,
  neighbors_within_batch = 3,
  n_pcs = 50,
  method = c("annoy", "nndescent"),
  metric = "euclidean",
  n_trees = 10L,
  k_build_nndescent = 30,
  trim = NULL,
  set_op_mix_ratio = 1,
  local_connectivity = 1,
  seed = 42,
  verbose = TRUE,
  ...
)

## S3 method for class 'Seurat'
RunBBKNN(
  object,
  batch_key,
  assay = NULL,
  reduction = "pca",
  n_pcs = 50L,
  graph_name = "bbknn",
  set_op_mix_ratio = 1,
  local_connectivity = 1,
  run_TSNE = TRUE,
  TSNE_name = "tsne",
  TSNE_key = "tsNE_",
  run_UMAP = TRUE,
  UMAP_name = "umap",
  UMAP_key = "UMAP_",
  return.umap.model = FALSE,
  min_dist = 0.3,
  spread = 1,
  seed = 42,
  verbose = TRUE,
  ...
)
```

**Arguments**

object            An object

...	Arguments passed to other methods
batch_list	A character vector with the same length as nrow(pca)
neighbors_within_batch	How many top neighbours to report for each batch; total number of neighbours in the initial k-nearest-neighbours computation will be this number times the number of batches. This then serves as the basis for the construction of a symmetrical matrix of connectivities.
n_pcs	Number of dimensions to use. Default is 50.
method	Method to find k nearest neighbors (kNNs). One of "annoy" and "nndescent".
metric	Metric to calculate the distances. The options depend on the choice of kNN method. The following metrics are supported in both annoy and nndescent: <ul style="list-style-type: none"> <li>• 'euclidean' (the default)</li> <li>• 'manhattan'</li> <li>• 'hamming'</li> </ul> <p>The following metrics are only supported in nndescent:</p> <ul style="list-style-type: none"> <li>• 'sqeuclidean'</li> <li>• 'chebyshev'</li> <li>• 'canberra'</li> <li>• 'braycurtis'</li> <li>• 'cosine'</li> <li>• 'correlation'</li> <li>• 'jaccard'</li> <li>• 'dice'</li> <li>• 'matching'</li> <li>• 'russellrao'</li> <li>• 'kulsinski'</li> <li>• 'rogerstanimoto'</li> <li>• 'sokalmichener'</li> <li>• 'sokalsneath'</li> <li>• 'tsss'</li> <li>• 'yule'</li> <li>• 'hellinger'</li> </ul>
n_trees	The number of trees to use in the random projection forest. More trees give higher precision when querying, at the cost of increased run time and resource intensity.
k_build_nndescent	Used with nndescent neighbour identification. The number of neighbours to include when building the approximate nearest neighbors index and neighbor graph. More neighbours give higher precision when querying, at the cost of increased run time and resource intensity.
trim	Trim the neighbours of each cell to these many top connectivities. May help with population independence and improve the tidiness of clustering. The lower the value the more independent the individual populations, at the cost of more conserved batch effect. Default is 10 times neighbors_within_batch times the number of batches. Set to 0 to skip.

set_op_mix_ratio	Pass to 'set_op_mix_ratio' parameter for <a href="#">umap</a>
local_connectivity	Pass to 'local_connectivity' parameter for <a href="#">umap</a>
seed	Set a random seed. By default, sets the seed to 42. Setting NULL will not set a seed.
verbose	Whether or not to print output to the console
batch_key	Column name in meta.data discriminating between your batches.
assay	Used to construct Graph.
reduction	Which dimensional reduction to use for the BBKNN input. Default is PCA
graph_name	Name of the generated BBKNN graph. Default is "bbknn".
run_TSNE	Whether or not to run t-SNE based on BBKNN results.
TSNE_name	Name to store t-SNE dimensional reduction.
TSNE_key	Specifies the string before the number of the t-SNE dimension names. tSNE by default.
run_UMAP	Whether or not to run UMAP based on BBKNN results.
UMAP_name	Name to store UMAP dimensional reduction.
UMAP_key	Specifies the string before the number of the UMAP dimension names. tSNE by default.
return.umap.model	Whether UMAP will return the uwot model.
min_dist	Pass to 'min_dist' parameter for <a href="#">umap</a>
spread	Pass to 'spread' parameter for <a href="#">umap</a>

**Value**

Returns a Seurat object containing a new BBKNN Graph and Neighbor data. If run t-SNE or UMAP, will also return corresponded reduction objects.

**References**

Polański, Krzysztof, et al. "BBKNN: fast batch alignment of single cell transcriptomes." *Bioinformatics* 36.3 (2020): 964-965.

**Examples**

```
data("panc8_small")
panc8_small <- RunBBKNN(panc8_small, "tech")
```

# Index

## \* datasets

panc8\_small, [2](#)

glmnet, [4](#)

panc8\_small, [2](#)

RidgeRegression, [3](#)

RunBKNN, [4](#)

umap, [7](#)