

Package ‘KMEANS.KNN’

January 20, 2025

Title KMeans and KNN Clustering Package

Version 0.1.0

Description Implementation of Kmeans clustering algorithm and a supervised KNN (K Nearest Neighbors) learning method. It allows users to perform unsupervised clustering and supervised classification on their datasets. Additional features include data normalization, imputation of missing values, and the choice of distance metric. The package also provides functions to determine the optimal number of clusters for Kmeans and the best k-value for KNN: `knn_Function()`, `find_Knn_best_k()`, `KMEANS_FUNCTION()`, and `find_Kmeans_best_k()`.

License GPL-3

Encoding UTF-8

RoxygenNote 7.3.1

Imports factoextra, cluster, ggplot2, stats, assertthat, class, caret, grDevices

Suggests knitr, rmarkdown, testthat (>= 3.0.0)

Config/testthat/edition 3

VignetteBuilder knitr

NeedsCompilation no

Author LALLOGO Lassané [aut, cre] (<<https://orcid.org/0009-0004-1637-3511>>)

Maintainer LALLOGO Lassané <lallogo2002@gmail.com>

Repository CRAN

Date/Publication 2024-05-17 09:20:12 UTC

Contents

find_Kmeans_best_k	2
find_Knn_best_k	2
KMEANS_FUNCTION	3
knn_Function	4
Index	6

`find_Kmeans_best_k` *find_Kmeans_best_k*

Description

`find_Kmeans_best_k`

Usage

```
find_Kmeans_best_k(data, max_k = 10, Method = "coude", verbose = FALSE)
```

Arguments

<code>data</code>	The dataset for which K-means clustering will be performed.
<code>max_k</code>	The maximum number of clusters to consider. It defaults to 10.
<code>Method</code>	The method used to determine the optimal number of clusters. Acceptable values are "coude" (elbow method), "silhouette" (silhouette method), or "gap" (gap statistics).
<code>verbose</code>	Logical. If TRUE, additional output is provided.

Value

This function does not return a value but prints the optimal number of clusters based on the chosen method and plots the corresponding graph.

Examples

```
data(iris)
find_Kmeans_best_k(iris[,-5],9,Method = "coude")
```

`find_Knn_best_k` *find_Knn_best_k*

Description

This function finds the best k-value for KNN based on the provided data.

Usage

```
find_Knn_best_k(data, target_column, k_values, Prop_train = 0.8)
```

Arguments

<code>data</code>	A dataframe containing the dataset to be used.
<code>target_column</code>	A string specifying the name of the target column in the dataset.
<code>k_values</code>	A numeric vector containing the different k-values to be tested.
<code>Prop_train</code>	A numeric value between 0 and 1 indicating the proportion of the dataset to be used for training.

Value

A list containing a dataframe with k-values and their corresponding accuracies, and the best k-value with its accuracy.

Examples

```
data(iris)
find_Knn_best_k(iris,"Species",1:10,Prop_train=0.8)
```

KMEANS_FUNCTION	<i>KMEANS_FUNCTION</i>
-----------------	------------------------

Description

This function implements the K-Means algorithm for data clustering. It provides options for data preprocessing, such as normalization and imputation of missing values.

Usage

```
KMEANS_FUNCTION(
  data,
  k,
  max_iter = 100,
  nstart = 25,
  distance_metric = "euclidean",
  scale_data = FALSE,
  impute_data = "mean"
)
```

Arguments

<code>data</code>	A dataframe containing the numerical data to be clustered.
<code>k</code>	The number of clusters to form.
<code>max_iter</code>	The maximum number of iterations for the K-Means algorithm.
<code>nstart</code>	The number of times to randomly initialize the centroids.
<code>distance_metric</code>	The distance metric to use ('euclidean' or 'manhattan').
<code>scale_data</code>	A boolean indicating whether the data should be normalized.
<code>impute_data</code>	The imputation method for missing values ('mean', 'median', 'mode').

Value

A list containing the following elements: - clusters: A vector indicating the cluster of each point. - centers: The coordinates of the centroids of each cluster. - additional_info: Additional information such as total distance and number of iterations.

Examples

```
data(iris)
data_iris <- iris[, -5] # Exclude the species column
results <- KMEANS_FUNCTION(data_iris, k = 3)
print(results$clusters)
```

knn_Function

knn_Function

Description

This function implements a custom K-Nearest Neighbors (KNN) algorithm with data preprocessing options. It predicts the class of a new point based on the k closest neighbors in the feature space.

Usage

```
knn_Function(
  new_points,
  dataset,
  k = 5,
  distance_metric = "gower",
  target_variable,
  scale_data = TRUE,
  impute_data = "mean",
  weight_votes = TRUE
)
```

Arguments

new_points	A dataframe of new points to be classified.
dataset	A dataframe of training data.
k	The number of nearest neighbors to consider.
distance_metric	The distance metric for calculating neighbors ('gower', 'euclidean', 'manhattan').
target_variable	The name of the target variable in 'dataset'.
scale_data	A boolean to indicate whether the data should be normalized.
impute_data	The imputation method for missing values ('mean', 'median', 'mode').
weight_votes	A boolean to indicate whether votes should be weighted by the inverse of the distance.

Value

A list containing 'Predictions' with the predicted class for each new point, 'Data' with the 'new_points' dataframe and an additional column for predictions, 'Distances' with the distances of the k nearest neighbors, and 'Imputed_Values' with the imputed values for missing variables.

Examples

```
# Loading training data (e.g., iris)
data(iris)

# Preparing new points for prediction (e.g., two new observations)
new_points <- data.frame(Sepal.Length = c(5.1, 7.7, 1.3, 0.2, 5.1),
  Sepal.Width = c(3.5, 2.6, 5, 3.7, 3.5),
  Petal.Length = c(1.4, 6.9, 4.5, 6, 3.4),
  Petal.Width = c(10.1, 7.6, 5.6, 8.4, 5.2))

# Calling the custom KNN function
results <- knn_Function(new_points, dataset = iris, k = 3, target_variable = "Species")

# Displaying predictions
print(results$Predictions)
```

Index

`find_Kmeans_best_k`, 2

`find_Knn_best_k`, 2

`KMEANS_FUNCTION`, 3

`knn_Function`, 4