

Package ‘LLMing’

October 21, 2025

Title Large Language Model (LLM) Tools for Psychological Text Analysis

Version 1.0.0

Maintainer Lindley Slipetz <ddj6tu@virginia.edu>

Description A collection of large language model (LLM) text analysis methods designed with psychological data in mind. Currently, LLMing (aka ``lemming'') includes a text anomaly detection method based on the angle-based subspace approach described by Zhang, Lin, and Karim (2015)
<[doi:10.1016/j.ress.2015.05.025](https://doi.org/10.1016/j.ress.2015.05.025)>.

License MIT + file LICENSE

Encoding UTF-8

RoxygenNote 7.3.3

Imports Rdpack, quanteda, stopwords, stringi, reticulate, text, dbscan, pracma, stats

RdMacros Rdpack

URL <https://github.com/sliplr19/LLMing>

BugReports <https://github.com/sliplr19/LLMing/issues>

NeedsCompilation no

Author Lindley Slipetz [aut, cre],
Teague Henry [aut],
Siqi Sun [ctb]

Depends R (>= 4.1.0)

Repository CRAN

Date/Publication 2025-10-21 17:50:06 UTC

Contents

| | |
|---------------------|---|
| embed | 2 |
| G_thres | 3 |
| normahalo | 3 |
| pCOS | 4 |

| | |
|-----------------------|---|
| pCOS_row | 4 |
| rep_set | 5 |
| sim_SNN | 5 |
| textanomaly | 6 |
| vector_SNN | 6 |
| z_score | 7 |

Index

8

embed*Embed texts with a Transformer model***Description**

Cleans a text column and converts it to a data frame of numeric vectors via BERT embeddings. For the input data frame, each row is one text entry.

Usage

```
embed(dat, layers, keep_tokens = TRUE, tokens_method = NULL)
```

Arguments

| | |
|----------------------|--|
| dat | A data frame with text data, one text per row |
| layers | Integer vector specifying which model layers to aggregate from. |
| keep_tokens | Logical, keep token-level embeddings in the returned object or discard them to save memory |
| tokens_method | Character scalar controlling how token-level embeddings are aggregated to word types |

Value

A data frame where each row corresponds to one input text and each column is an embedding dimension

```
@examples df <- data.frame( text = c( "I slept well and feel great today!", "I saw from friends and it went well.", "I think I failed that exam. I'm such a disappointment." ) )
```

```
emb_dat <- embed( dat = df, layers = 1, keep_tokens = FALSE, tokens_method = "mean" )
```

| | |
|----------------------|---------------------------------------|
| <code>G_thres</code> | <i>Thresholding of pCOS dataframe</i> |
|----------------------|---------------------------------------|

Description

Converts each column of a pCOS score matrix into binary indicators

Usage

```
G_thres(pCOS_mat, theta)
```

Arguments

| | |
|-----------------------|--------------------------|
| <code>pCOS_mat</code> | Dataframe of pCOS values |
| <code>theta</code> | Numeric threshold |

Value

A matrix of 0s and 1s of which cells meet the threshold

Examples

```
z_dat <- data.frame("A" = rnorm(500,0,1), "B" = rnorm(500,0,1), "C" = rnorm(500,0,1))
snn <- sim_SNN(z_dat, 10, 5)
vec_snn <- vector_SNN(z_dat, snn)
pCOSdat <- pCOS(z_dat, vec_snn)
G <- G_thres(pCOSdat, theta = 0.1)
```

| | |
|------------------------|----------------------------|
| <code>normahalo</code> | <i>Local outlier score</i> |
|------------------------|----------------------------|

Description

Computes a normalized Mahalanobis distance score. Only features with nonzero scores in S receive nonzero Mahalanobis scores.

Usage

```
normahalo(z, rs, S)
```

Arguments

| | |
|-----------------|-----------------------------|
| <code>z</code> | Dataframe of z scores |
| <code>rs</code> | List of reference sets |
| <code>S</code> | Dataframe of numeric values |

Value

A dataframe of local outlier scores

pCOS

pCOS scores for every row of dataframe

Description

Applies pCOS_row() to corresponding rows of two data frames, returning one pCOS value per row.

Usage

`pCOS(z_dat, vec_SNN)`

Arguments

| | |
|----------------------|---|
| <code>z_dat</code> | Numeric dataframe, typically z-scores |
| <code>vec_SNN</code> | Numeric dataframe, typically the output of vector_SNN |

Value

A dataframe with same dimensions as `z_dat`

pCOS_row

Pairwise cosine-style row score

Description

Given two numeric vectors, computes an average cosine-based similarity.

Usage

`pCOS_row(z, v_SNN)`

Arguments

| | |
|--------------------|---|
| <code>z</code> | Numeric vector |
| <code>v_SNN</code> | Numeric vector, same size as <code>z</code> |

Value

A numeric vector

| | |
|---------|--|
| rep_set | <i>The vectors of the shared nearest neighbors</i> |
|---------|--|

Description

Creates a list of the vectors of the top shared nearest neighbors for each row of the z dataframe

Usage

```
rep_set(z, snn)
```

Arguments

- | | |
|-----|---|
| z | Dataframe of values of reference set |
| snn | Dataframe of shared nearest neighbors indices |

Value

A list of dataframes where each row of the dataframe is the vector representation of a given shared nearest neighbor

| | |
|---------|---|
| sim_SNN | <i>Compute shared nearest neighbors</i> |
|---------|---|

Description

Builds a shared nearest neighbors matrix and, for each row (observation), returns the indices of the top neighbors with the largest SNN overlap counts

Usage

```
sim_SNN(z_dat, k, tops)
```

Arguments

- | | |
|-------|--|
| z_dat | A dataframe with numeric columns |
| k | An integer representing number of nearest neighbors |
| tops | An integer representing how many of shared nearest neighbors to return |

Value

A dataframe of top rows with shared nearest neighbors

| | |
|-------------|---------------------------|
| textanomaly | <i>Text anomaly score</i> |
|-------------|---------------------------|

Description

Text anomaly detection method adapted from (Zhang et al. 2015).

Usage

```
textanomaly(dat, k, tops, theta)
```

Arguments

| | |
|-------|--|
| dat | A dataframe with text data, one text per row |
| k | An integer representing number of nearest neighbors |
| tops | An integer representing how many of shared nearest neighbors to return |
| theta | Numeric threshold |

Value

Dataframe of local outlier score

References

Zhang L, Lin J, Karim R (2015). “An angle-based subspace anomaly detection approach to high-dimensional data: With an application to industrial fault detection.” *Reliability Engineering & System Safety*, **142**, 482–497. ISSN 0951-8320, doi:[10.1016/j.ress.2015.05.025](https://doi.org/10.1016/j.ress.2015.05.025).

| | |
|------------|--|
| vector_SNN | <i>Aggregate dataframe into mean feature vectors</i> |
|------------|--|

Description

For each row of the SNN index matrix, this function takes the rows of reference dataframe, z, and computes their column means, yielding one mean vector per observation.

Usage

```
vector_SNN(z, snn)
```

Arguments

| | |
|-----|---|
| z | Numeric dataframe |
| snn | Dataframe of shared nearest neighbors indices |

Value

Dataframe of same dimensions as z

z_score

Z-score on columns

Description

Z-score on columns

Usage

`z_score(dat)`

Arguments

`dat` A dataframe with numeric cells

Value

A dataframe with numeric cells with the same dimensions as dat

Index

embed, 2

G_thres, 3

normahalo, 3

pCOS, 4

pCOS_row, 4

rep_set, 5

sim_SNN, 5

textanomaly, 6

vector_SNN, 6

z_score, 7