

Package ‘text’

October 14, 2022

Type Package

Title Analyses of Text using Transformers Models from HuggingFace,
Natural Language Processing and Machine Learning

Version 0.9.99.2

Description Link R with Transformers from Hugging Face to transform text variables to word embeddings; where the word embeddings are used to statistically test the mean difference between set of texts, compute semantic similarity scores between texts, predict numerical variables, and visual statistically significant words according to various dimensions etc. For more information see <<https://www.r-text.org>>.

License GPL-3

URL <https://r-text.org/>, <https://github.com/OscarKjell/text/>

BugReports <https://github.com/OscarKjell/text/issues/>

Encoding UTF-8

Archs x64

SystemRequirements Python (>= 3.6.0)

LazyData true

BuildVignettes true

Imports dplyr, tibble, stringi, tidyr, ggplot2, ggrepel, cowplot,
rlang, purrr, magrittr, parsnip, recipes, rsample, reticulate,
tune, workflows, yardstick, future, furr, overlapping

RoxygenNote 7.2.0

Suggests knitr, rmarkdown, testthat, rio, glmnet, randomForest, covr,
xml2, ranger

VignetteBuilder knitr

Depends R (>= 4.00)

NeedsCompilation no

Author Oscar Kjell [aut, cre] (<<https://orcid.org/0000-0002-2728-6278>>),
Salvatore Giorgi [aut] (<<https://orcid.org/0000-0001-7381-6295>>),
Andrew Schwartz [aut] (<<https://orcid.org/0000-0002-6383-3339>>)

Maintainer Oscar Kjell <oscar.kjell@psy.lu.se>

Repository CRAN

Date/Publication 2022-09-20 22:00:02 UTC

R topics documented:

centrality_data_harmony	3
DP_projections_HILS_SWLS_100	4
Language_based_assessment_data_3_100	5
Language_based_assessment_data_8	5
PC_projections_satisfactionwords_40	6
raw_embeddings_1	7
textCentrality	7
textCentralityPlot	8
textClassify	11
textDescriptives	13
textDimName	14
textDistance	15
textDistanceMatrix	16
textDistanceNorm	17
textEmbed	18
textEmbedLayerAggregation	20
textEmbedRawLayers	22
textEmbedStatic	23
textGeneration	24
textModelLayers	26
textModels	27
textModelsRemove	27
textNER	28
textPCA	29
textPCAPlot	30
textPlot	33
textPredict	37
textPredictAll	38
textPredictTest	39
textProjection	40
textProjectionPlot	42
textQA	47
textrpp_initialize	49
textrpp_install	50
textrpp_uninstall	51
textSimilarity	52
textSimilarityMatrix	53
textSimilarityNorm	54
textSimilarityTest	55
textSum	56
textTokenize	58
textTrain	59

<code>centrality_data_harmony</code>	3
<code>textTrainLists</code>	60
<code>textTrainRandomForest</code>	61
<code>textTrainRegression</code>	64
<code>textTranslate</code>	67
<code>textWordPrediction</code>	68
<code>textZeroShot</code>	69
<code>word_embeddings_4</code>	71
Index	72

`centrality_data_harmony`
Example data for plotting a Semantic Centrality Plot.

Description

The dataset is a shortened version of the data sets of Study 1 from Kjell, et al., 2016.

Usage

`centrality_data_harmony`

Format

A data frame with 2,146 and 4 variables:

words unique words

n overall word frequency

central_semantic_similarity cosine semantic similarity to the aggregated word embedding

n_percent frequency in percent

Source

<https://link.springer.com/article/10.1007/s11205-015-0903-z>

DP_projections_HILS_SWLS_100

Data for plotting a Dot Product Projection Plot.

Description

Tibble is the output from textProjection. The dataset is a shortened version of the data sets of Study 3-5 from Kjell, Kjell, Garcia and Sikström 2018.

Usage

DP_projections_HILS_SWLS_100

Format

A data frame with 583 rows and 12 variables:

words unique words

dot.x dot product projection on the x-axes

p_values_dot.x p-value for the word in relation to the x-axes

n_g1.x frequency of the word in group 1 on the x-axes variable

n_g2.x frequency of the word in group 2 on the x-axes variable

dot.y dot product projection on the y-axes

p_values_dot.y p-value for the word in relation to the y-axes

n_g1.y frequency of the word in group 1 on the y-axes variable

n_g2.y frequency of the word in group 2 on the x-axes variable

n overall word frequency

n.percent frequency in percent

N_participant_responses number of participants (as this is needed in the analyses)

Source

<https://psyarxiv.com/er6t7/>

Language_based_assessment_data_3_100

Example text and numeric data.

Description

The dataset is a shortened version of the data sets of Study 3-5 from Kjell, Kjell, Garcia and Sikström 2018.

Usage

Language_based_assessment_data_3_100

Format

A data frame with 100 rows and 4 variables:

harmonywords Word responses from the harmony in life word question

hilstotal total score of the Harmony In Life Scale

swlstotal total score of the Satisfaction With Life Scale

Source

<https://psyarxiv.com/er6t7/>

Language_based_assessment_data_8

Text and numeric data for 10 participants.

Description

The dataset is a shortened version of the data sets of Study 3-5 from Kjell et al., (2018; <https://psyarxiv.com/er6t7/>).

Usage

Language_based_assessment_data_8

Format

A data frame with 40 participants and 8 variables:

harmonywords descriptive words where respondents describe their harmony in life

satisfactionwords descriptive words where respondents describe their satisfaction with life

harmonytexts text where respondents describe their harmony in life

satisfactiontexts text where respondents describe their satisfaction with life

hilstotal total score of the Harmony In Life Scale
swlstotal total score of the Satisfaction With Life Scale
age respondents age in years
gender respondents gender 1=male, 2=female

Source

<https://psyarxiv.com/er6t7/>

PC_projections_satisfactionwords_40

Example data for plotting a Principle Component Projection Plot.

Description

The dataset is a shortened version of the data sets of Study 1 from Kjell, et al., 2016.

Usage

PC_projections_satisfactionwords_40

Format

A data frame.

words unique words

n overall word frequency

Dim_PC1 Principle component value for dimension 1

Dim_PC2 Principle component value for dimension 2

Source

<https://link.springer.com/article/10.1007/s11205-015-0903-z>

raw_embeddings_1	<i>Word embeddings from textEmbedRawLayers function</i>
------------------	---

Description

The dataset is a shortened version of the data sets of Study 3-5 from Kjell, Kjell, Garcia and Sikström 2018.

Usage

```
raw_embeddings_1
```

Format

A list with token-level word embeddings for harmony words.

tokens words

layer_number layer of the transformer model

Dim1:Dim8 Word embeddings dimensions

Source

<https://psyarxiv.com/er6t7/>

textCentrality	<i>Compute semantic similarity score between single words' word embeddings and the aggregated word embedding of all words.</i>
----------------	--

Description

Compute semantic similarity score between single words' word embeddings and the aggregated word embedding of all words.

Usage

```
textCentrality(  
  words,  
  word_embeddings,  
  word_types_embeddings = word_types_embeddings_df,  
  method = "cosine",  
  aggregation = "mean",  
  min_freq_words_test = 0  
)
```

Arguments

words	Word or text variable to be plotted.
word_embeddings	Word embeddings from textEmbed for the words to be plotted (i.e., the aggregated word embeddings for the "words" variable).
word_types_embeddings	Word embeddings from textEmbed for individual words (i.e., the decontextualized word embeddings).
method	Character string describing type of measure to be computed. Default is "cosine" (see also "spearmen", "pearson" as well as measures from textDistance() (which here is computed as 1 - textDistance) including "euclidean", "maximum", "manhattan", "canberra", "binary" and "minkowski").
aggregation	Method to aggregate the word embeddings (default = "mean"; see also "min", "max" or "[CLS]").
min_freq_words_test	Option to select words that have at least occurred a specified number of times (default = 0); when creating the semantic similarity scores.

Value

A dataframe with variables (e.g., including semantic similarity, frequencies) for the individual words that are used for the plotting in the textCentralityPlot function.

See Also

see [textCentralityPlot](#) [textProjection](#)

Examples

```
## Not run:
df_for_plotting <- textCentrality(
  words = Language_based_assessment_data_8$harmonywords,
  word_embeddings = word_embeddings_4$texts$harmonywords,
  word_types_embeddings = word_embeddings_4$word_types
)
df_for_plotting

## End(Not run)
```

textCentralityPlot	<i>Plot words according to semantic similarity to the aggregated word embedding.</i>
--------------------	--

Description

Plot words according to semantic similarity to the aggregated word embedding.

Usage

```

textCentralityPlot(
  word_data,
  min_freq_words_test = 1,
  plot_n_word_extreme = 10,
  plot_n_word_frequency = 10,
  plot_n_words_middle = 10,
  titles_color = "#61605e",
  x_axes = "central_semantic_similarity",
  title_top = "Semantic Centrality Plot",
  x_axes_label = "Semantic Centrality",
  scale_x_axes_lim = NULL,
  scale_y_axes_lim = NULL,
  word_font = NULL,
  centrality_color_codes = c("#EAEAEA", "#85DB8E", "#398CF9", "#9e9d9d"),
  word_size_range = c(3, 8),
  position_jitter_hight = 0,
  position_jitter_width = 0.03,
  point_size = 0.5,
  arrow_transparency = 0.1,
  points_without_words_size = 0.5,
  points_without_words_alpha = 0.5,
  legend_title = "SC",
  legend_x_axes_label = "x",
  legend_x_position = 0.02,
  legend_y_position = 0.02,
  legend_h_size = 0.2,
  legend_w_size = 0.2,
  legend_title_size = 7,
  legend_number_size = 2,
  seed = 1007
)

```

Arguments

word_data Tibble from textPlotData.

min_freq_words_test
 Select words to significance test that have occurred at least `min_freq_words_test` (default = 1).

plot_n_word_extreme
 Number of words per dimension to plot with extreme Supervised Dimension Projection value. (i.e., even if not significant; duplicates are removed).

plot_n_word_frequency
 Number of words to plot according to their frequency. (i.e., even if not significant).

plot_n_words_middle
 Number of words to plot that are in the middle in Supervised Dimension Projection score (i.e., even if not significant; duplicates are removed).

titles_color	Color for all the titles (default: "#61605e").
x_axes	Variable to be plotted on the x-axes (default is "central_semantic_similarity", could also select "n", "n_percent").
title_top	Title (default " ").
x_axes_label	Label on the x-axes.
scale_x_axes_lim	Length of the x-axes (default: NULL, which uses $c(\min(\text{word_data}\$central_semantic_similarity)-0.05, \max(\text{word_data}\$central_semantic_similarity)+0.05)$; change this by e.g., try $c(-5, 5)$).
scale_y_axes_lim	Length of the y-axes (default: NULL, which uses $c(-1, 1)$; change e.g., by trying $c(-5, 5)$).
word_font	Type of font (default: NULL).
centrality_color_codes	Colors of the words selected as plot_n_word_extreme (minimum values), plot_n_words_middle, plot_n_word_extreme (maximum values) and plot_n_word_frequency; the default is $c("#EAEAEA", "#85DB8E", "#398CF9", "#9e9d9d", \text{respectively})$.
word_size_range	Vector with minimum and maximum font size (default: $c(3, 8)$).
position_jitter_high	Jitter height (default: .0).
position_jitter_width	Jitter width (default: .03).
point_size	Size of the points indicating the words' position (default: 0.5).
arrow_transparency	Transparency of the lines between each word and point (default: 0.1).
points_without_words_size	Size of the points not linked to a word (default is to not show the point; , i.e., 0).
points_without_words_alpha	Transparency of the points that are not linked to a word (default is to not show it; i.e., 0).
legend_title	Title of the color legend (default: "(SCP)").
legend_x_axes_label	Label on the color legend (default: "(x)").
legend_x_position	Position on the x coordinates of the color legend (default: 0.02).
legend_y_position	Position on the y coordinates of the color legend (default: 0.05).
legend_h_size	Height of the color legend (default 0.15).
legend_w_size	Width of the color legend (default 0.15).
legend_title_size	Font size of the title (default = 7).
legend_number_size	Font size of the values in the legend (default = 2).
seed	Set different seed.

Value

A 1-dimensional word plot based on similarity to the aggregated word embedding, as well as tibble with processed data used to plot.

See Also

see [textCentrality](#) and [textProjection](#)

Examples

```
# The test-data included in the package is called: centrality_data_harmony
names(centrality_data_harmony)
# Plot
# centrality_plot <- textCentralityPlot(
#   word_data = centrality_data_harmony,
#   min_freq_words_test = 10,
#   plot_n_word_extreme = 10,
#   plot_n_word_frequency = 10,
#   plot_n_words_middle = 10,
#   titles_color = "#61605e",
#   x_axes = "central_semantic_similarity",
#
#   title_top = "Semantic Centrality Plot",
#   x_axes_label = "Semantic Centrality",
#
#   word_font = NULL,
#   centrality_color_codes = c("#EAEAEA", "#85DB8E", "#398CF9", "#9e9d9d"),
#   word_size_range = c(3, 8),
#   point_size = 0.5,
#   arrow_transparency = 0.1,
#   points_without_words_size = 0.5,
#   points_without_words_alpha = 0.5,
# )
# centrality_plot
```

textClassify

Predict label and probability of a text using a pretrained classifier language model. (experimental)

Description

Predict label and probability of a text using a pretrained classifier language model. (experimental)

Usage

```
textClassify(
  x,
  model = "distilbert-base-uncased-finetuned-sst-2-english",
  device = "cpu",
```

```

tokenizer_parallelism = FALSE,
logging_level = "error",
return_incorrect_results = FALSE,
return_all_scores = FALSE,
function_to_apply = "none",
set_seed = 202208
)

```

Arguments

x	(string) A character variable or a tibble/dataframe with at least one character variable.
model	(string) Specification of a pre-trained classifier language model. For full list of options see pretrained classifier models at HuggingFace . For example use "cardiffnlp/twitter-roberta-base-sentiment", "distilbert-base-uncased-finetuned-sst-2-english".
device	(string) Device to use: 'cpu', 'gpu', or 'gpu:k' where k is a specific device number.
tokenizer_parallelism	(boolean) If TRUE this will turn on tokenizer parallelism.
logging_level	(string) Set the logging level. Options (ordered from less logging to more logging): critical, error, warning, info, debug
return_incorrect_results	(boolean) Stop returning some incorrectly formatted/structured results. This setting does CANNOT evaluate the actual results (whether or not they make sense, exist, etc.). All it does is to ensure the returned results are formatted correctly (e.g., does the question-answering dictionary contain the key "answer", is sentiments from textClassify containing the labels "positive" and "negative").
return_all_scores	(boolean) Whether to return all prediction scores or just the one of the predicted class.
function_to_apply	(string) The function to apply to the model outputs to retrieve the scores.
set_seed	(Integer) Set seed. There are four different values: "default": if the model has a single label, will apply the sigmoid function on the output. If the model has several labels, the softmax function will be applied on the output. "sigmoid": Applies the sigmoid function on the output. "softmax": Applies the softmax function on the output. "none": Does not apply any function on the output.

Value

A tibble with predicted labels and scores for each text variable. The comment of the object show the model-name and computation time.

See Also

see [textGeneration](#), [textNER](#), [textSum](#), [textQA](#), [textTranslate](#)

Examples

```
# classifications <- textClassify(x = Language_based_assessment_data_8[1:2, 1:2])
# classifications
# comment(classifications)
```

textDescriptives	<i>Compute descriptive statistics of character variables.</i>
------------------	---

Description

Compute descriptive statistics of character variables.

Usage

```
textDescriptives(
  words,
  compute_total = TRUE,
  entropy_unit = "log2",
  na.rm = TRUE
)
```

Arguments

words	One or several character variables; if its a tibble or dataframe, all the character variables will be selected.
compute_total	Boolean. If the input (words) is a tibble/dataframe with several character variables, a total variable is computed.
entropy_unit	The unit entropy is measured in. The default is to used bits (i.e., log2; see also, "log", "log10"). If a total score for several variables is computed, the text columns are combined using the dplyr unite function. For more information about the entropy see the entropy package and specifically its entropy.plugin function.
na.rm	Option to remove NAs when computing mean, median etc (see under return).

Value

A tibble with descriptive statistics, including variable = the variable names of input "words"; w_total = total number of words in the variable; w_mean = mean number of words in each row of the variable; w_median = median number of words in each row of the variable; w_range_min = smallest number of words of all rows; w_range_max = largest number of words of all rows; w_sd = the standard deviation of the number of words of all rows; unique_tokens = the unique number of tokens (using the word_tokenize function from python package nltk) n_token = number of tokens in the variable (using the word_tokenize function from python package nltk) entropy = the entropy of the variable. It is computed as the Shannon entropy H of a discrete random variable from the specified bin frequencies. (see library entropy and specifically the entropy.plugin function)

See Also

see [textEmbed](#)

Examples

```
## Not run:
textDescriptives(Language_based_assessment_data_8[1:2])

## End(Not run)
```

textDimName	<i>Change the names of the dimensions in the word embeddings.</i>
-------------	---

Description

Change the names of the dimensions in the word embeddings.

Usage

```
textDimName(word_embeddings, dim_names = TRUE)
```

Arguments

word_embeddings	List of word embeddings
dim_names	(boolean) If TRUE the word embedding name will be attached to the name of each dimension; is FALSE, the attached part of the name will be removed.

Value

Word embeddings with changed names.

See Also

see [textEmbed](#)

Examples

```
# Note that dimensions are called Dim1_harmonytexts etc.
word_embeddings_4$texts$harmonytexts
# Here they are changed to just Dim
w_e_T <- textDimName(word_embeddings_4$texts["harmonytexts"],
  dim_names = FALSE
)
# Here they are changed back
w_e_F <- textDimName(w_e_T, dim_names = TRUE)
```

textDistance	<i>Compute the semantic distance between two text variables.</i>
--------------	--

Description

Compute the semantic distance between two text variables.

Usage

```
textDistance(x, y, method = "euclidean", center = FALSE, scale = FALSE)
```

Arguments

x	Word embeddings (from textEmbed).
y	Word embeddings (from textEmbed).
method	Character string describing type of measure to be computed; default is "euclidean" (see also measures from stats::dist() including "maximum", "manhattan", "canberra", "binary" and "minkowski". It is also possible to use "cosine", which computes the cosine distance (i.e., $1 - \cosine(x, y)$).
center	(boolean; from base::scale) If center is TRUE then centering is done by subtracting the column means (omitting NAs) of x from their corresponding columns, and if center is FALSE, no centering is done.
scale	(boolean; from base::scale) If scale is TRUE then scaling is done by dividing the (centered) columns of x by their standard deviations if center is TRUE, and the root mean square otherwise.

Value

A vector comprising semantic distance scores.

See Also

see [textSimilarity](#), [textSimilarityNorm](#) and [textSimilarityTest](#)

Examples

```
library(dplyr)
distance_scores <- textDistance(
  x = word_embeddings_4$texts$harmonytext,
  y = word_embeddings_4$texts$satisfactiontext
)
comment(distance_scores)
```

textDistanceMatrix	<i>Compute semantic distance scores between all combinations in a word embedding</i>
--------------------	--

Description

Compute semantic distance scores between all combinations in a word embedding

Usage

```
textDistanceMatrix(x, method = "euclidean", center = FALSE, scale = FALSE)
```

Arguments

x	Word embeddings (from textEmbed).
method	Character string describing type of measure to be computed; default is "euclidean" (see also measures from stats:dist() including "maximum", "manhattan", "canberra", "binary" and "minkowski". It is also possible to use "cosine", which computes the cosine distance (i.e., $1 - \text{cosine}(x, y)$).
center	(boolean; from base::scale) If center is TRUE then centering is done by subtracting the column means (omitting NAs) of x from their corresponding columns, and if center is FALSE, no centering is done.
scale	(boolean; from base::scale) If scale is TRUE then scaling is done by dividing the (centered) columns of x by their standard deviations if center is TRUE, and the root mean square otherwise.

Value

A matrix of semantic distance scores

See Also

see [textDistanceNorm](#) and [textSimilarityTest](#)

Examples

```
distance_scores <- textDistanceMatrix(word_embeddings_4$texts$harmonytext[1:3, ])  
round(distance_scores, 3)
```

textDistanceNorm	<i>Compute the semantic distance between a text variable and a word norm (i.e., a text represented by one word embedding that represent a construct/concept).</i>
------------------	---

Description

Compute the semantic distance between a text variable and a word norm (i.e., a text represented by one word embedding that represent a construct/concept).

Usage

```
textDistanceNorm(x, y, method = "euclidean", center = FALSE, scale = FALSE)
```

Arguments

x	Word embeddings (from textEmbed).
y	Word embedding from textEmbed (from only one text).
method	Character string describing type of measure to be computed; default is "euclidean" (see also measures from stats:dist() including "maximum", "manhattan", "canberra", "binary" and "minkowski". It is also possible to use "cosine", which computes the cosine distance (i.e., 1 - cosine(x, y)).
center	(boolean; from base::scale) If center is TRUE then centering is done by subtracting the column means (omitting NAs) of x from their corresponding columns, and if center is FALSE, no centering is done.
scale	(boolean; from base::scale) If scale is TRUE then scaling is done by dividing the (centered) columns of x by their standard deviations if center is TRUE, and the root mean square otherwise.

Value

A vector comprising semantic distance scores.

See Also

see [textDistance](#) and [textSimilarityTest](#)

Examples

```
## Not run:
library(dplyr)
library(tibble)
harmonynorm <- c("harmony peace ")
satisfactionnorm <- c("satisfaction achievement")

norms <- tibble::tibble(harmonynorm, satisfactionnorm)
word_embeddings <- word_embeddings_4$texts
```

```

word_embeddings_wordnorm <- textEmbed(norms)
similarity_scores <- textDistanceNorm(
  word_embeddings$harmonytext,
  word_embeddings_wordnorm$harmonynorm
)

## End(Not run)

```

textEmbed	<i>Extract layers and aggregate them to word embeddings, for all character variables in a given dataframe.</i>
-----------	--

Description

Extract layers and aggregate them to word embeddings, for all character variables in a given dataframe.

Usage

```

textEmbed(
  texts,
  model = "bert-base-uncased",
  layers = -2,
  dim_name = TRUE,
  aggregation_from_layers_to_tokens = "concatenate",
  aggregation_from_tokens_to_texts = "mean",
  aggregation_from_tokens_to_word_types = NULL,
  keep_token_embeddings = TRUE,
  tokens_select = NULL,
  tokens_deselect = NULL,
  decontextualize = FALSE,
  model_max_length = NULL,
  max_token_to_sentence = 4,
  tokenizer_parallelism = FALSE,
  device = "gpu",
  logging_level = "error"
)

```

Arguments

texts	A character variable or a tibble/dataframe with at least one character variable.
model	Character string specifying pre-trained language model (default 'bert-base-uncased'). For full list of options see pretrained models at HuggingFace . For example use "bert-base-multilingual-cased", "openai-gpt", "gpt2", "ctrl", "transfo-xl-wt103", "xlnet-base-cased", "xlm-mlm-enfr-1024", "distilbert-base-cased", "roberta-base", or "xlm-roberta-base".

layers	(string or numeric) Specify the layers that should be extracted (default -2 which give the second to last layer). It is more efficient to only extract the layers that you need (e.g., 11). You can also extract several (e.g., 11:12), or all by setting this parameter to "all". Layer 0 is the decontextualized input layer (i.e., not comprising hidden states) and thus should normally not be used. These layers can then be aggregated in the textEmbedLayerAggregation function.
dim_name	Boolean, if TRUE append the variable name after all variable-names in the output. (This differentiates between word embedding dimension names; e.g., Dim1_text_variable_name). see textDimName to change names back and forth.
aggregation_from_layers_to_tokens	(string) Aggregated layers of each token. Method to aggregate the contextualized layers (e.g., "mean", "min" or "max, which takes the minimum, maximum or mean, respectively, across each column; or "concatenate", which links together each word embedding layer to one long row.
aggregation_from_tokens_to_texts	(string) Aggregates to the individual text (i.e., the aggregation of all tokens/words given to the transformer).
aggregation_from_tokens_to_word_types	(string) Aggregates to the word type (i.e., the individual words) rather than texts.
keep_token_embeddings	(boolean) Whether to also keep token embeddings when using texts or word types aggregation.
tokens_select	Option to select word embeddings linked to specific tokens such as [CLS] and [SEP] for the context embeddings.
tokens_deselect	Option to deselect embeddings linked to specific tokens such as [CLS] and [SEP] for the context embeddings.
decontextualize	(boolean) Provide word embeddings of single words as input to the model (these embeddings are, e.g., used for plotting; default is to use). If using this, then set single_context_embeddings to FALSE.
model_max_length	The maximum length (in number of tokens) for the inputs to the transformer model (default the value stored for the associated model).
max_token_to_sentence	(numeric) Maximum number of tokens in a string to handle before switching to embedding text sentence by sentence.
tokenizer_parallelism	(boolean) If TRUE this will turn on tokenizer parallelism. Default FALSE.
device	Name of device to use: 'cpu', 'gpu', or 'gpu:k' where k is a specific device number
logging_level	Set the logging level. Default: "warning". Options (ordered from less logging to more logging): critical, error, warning, info, debug

Value

A tibble with tokens, a column for layer identifier and word embeddings. Note that layer 0 is the input embedding to the transformer

See Also

see [textEmbedLayerAggregation](#), [textEmbedRawLayers](#) and [textDimName](#)

Examples

```
# word_embeddings <- textEmbed(Language_based_assessment_data_8[1:2, 1:2],
#                               layers = 10:11,
#                               aggregation_from_layers_to_tokens = "concatenate",
#                               aggregation_from_tokens_to_texts = "mean",
#                               aggregation_from_tokens_to_word_types = "mean")
## Show information about how the embeddings were constructed
# comment(word_embeddings$texts$satisfactiontexts)
# comment(word_embeddings$word_types)
# comment(word_embeddings$tokens$satisfactiontexts)
```

textEmbedLayerAggregation

Select and aggregate layers of hidden states to form a word embeddings.

Description

Select and aggregate layers of hidden states to form a word embeddings.

Usage

```
textEmbedLayerAggregation(
  word_embeddings_layers,
  layers = "all",
  aggregation_from_layers_to_tokens = "concatenate",
  aggregation_from_tokens_to_texts = "mean",
  return_tokens = FALSE,
  tokens_select = NULL,
  tokens_deselect = NULL
)
```

Arguments

word_embeddings_layers	Layers outputted from textEmbedRawLayers.
layers	The numbers of the layers to be aggregated (e.g., c(11:12) to aggregate the eleventh and twelfth). Note that layer 0 is the input embedding to the transformer, and should normally not be used. Selecting 'all' thus removes layer 0.
aggregation_from_layers_to_tokens	Method to carry out the aggregation among the layers for each word/token, including "min", "max" and "mean" which takes the minimum, maximum or mean across each column; or "concatenate", which links together each layer of the word embedding to one long row. Default is "concatenate"
aggregation_from_tokens_to_texts	Method to carry out the aggregation among the word embeddings for the words/tokens, including "min", "max" and "mean" which takes the minimum, maximum or mean across each column; or "concatenate", which links together each layer of the word embedding to one long row.
return_tokens	If TRUE, provide the tokens used in the specified transformer model.
tokens_select	Option to only select embeddings linked to specific tokens such as "[CLS]" and "[SEP]" (default NULL).
tokens_deselect	Option to deselect embeddings linked to specific tokens such as "[CLS]" and "[SEP]" (default NULL).

Value

A tibble with word embeddings. Note that layer 0 is the input embedding to the transformer, which is normally not used.

See Also

see [textEmbedRawLayers](#) and [textEmbed](#)

Examples

```
# word_embeddings_layers <- textEmbedRawLayers(Language_based_assessment_data_8$harmonywords[1],  
# layers = 11:12)  
# word_embeddings <- textEmbedLayerAggregation(word_embeddings_layers$context, layers = 11)
```

textEmbedRawLayers	<i>Extract layers of hidden states (word embeddings) for all character variables in a given dataframe.</i>
--------------------	--

Description

Extract layers of hidden states (word embeddings) for all character variables in a given dataframe.

Usage

```
textEmbedRawLayers(
  texts,
  model = "bert-base-uncased",
  layers = -2,
  return_tokens = TRUE,
  word_type_embeddings = FALSE,
  decontextualize = FALSE,
  keep_token_embeddings = TRUE,
  device = "cpu",
  tokenizer_parallelism = FALSE,
  model_max_length = NULL,
  max_token_to_sentence = 4,
  logging_level = "error"
)
```

Arguments

texts	A character variable or a tibble/dataframe with at least one character variable.
model	Character string specifying pre-trained language model (default 'bert-base-uncased'). For full list of options see pretrained models at HuggingFace . For example use "bert-base-multilingual-cased", "openai-gpt", "gpt2", "ctrl", "transfo-xl-wt103", "xlnet-base-cased", "xlm-mlm-enfr-1024", "distilbert-base-cased", "roberta-base", or "xlm-roberta-base".
layers	(string or numeric) Specify the layers that should be extracted (default -2, which give the second to last layer). It is more efficient to only extract the layers that you need (e.g., 11). You can also extract several (e.g., 11:12), or all by setting this parameter to "all". Layer 0 is the decontextualized input layer (i.e., not comprising hidden states) and thus should normally not be used. These layers can then be aggregated in the textEmbedLayerAggregation function.
return_tokens	If TRUE, provide the tokens used in the specified transformer model.
word_type_embeddings	(boolean) Wether to provide embeddings for each word/token type.
decontextualize	(boolean) Wether to dectonextualise embeddings (i.e., embedding one word at a time).

keep_token_embeddings	(boolean) Whether to keep token level embeddings in the output (when using word_types aggregation)
device	Name of device to use: 'cpu', 'gpu', or 'gpu:k' where k is a specific device number
tokenizer_parallelism	If TRUE this will turn on tokenizer parallelism. Default FALSE.
model_max_length	The maximum length (in number of tokens) for the inputs to the transformer model (default the value stored for the associated model).
max_token_to_sentence	(numeric) Maximum number of tokens in a string to handle before switching to embedding text sentence by sentence.
logging_level	Set the logging level. Default: "warning". Options (ordered from less logging to more logging): critical, error, warning, info, debug

Value

Returns hiddenstates/layers that can be 1. Can return three different outputA tibble with tokens, column specifying layer and word embeddings. Note that layer 0 is the input embedding to the transformer, and should normally not be used.

See Also

see [textEmbedLayerAggregation](#) and [textEmbed](#)

Examples

```
# texts <- Language_based_assessment_data_8[1:2, 1:2]
# word_embeddings_with_layers <- textEmbedRawLayers(texts, layers = 11:12)
```

textEmbedStatic	<i>Applies word embeddings from a given decontextualized static space (such as from Latent Semantic Analyses) to all character variables</i>
-----------------	--

Description

Applies word embeddings from a given decontextualized static space (such as from Latent Semantic Analyses) to all character variables

Usage

```

textEmbedStatic(
  df,
  space,
  tk_df = "null",
  aggregation_from_tokens_to_texts = "mean",
  dim_name = FALSE,
  tolower = FALSE
)

```

Arguments

df	dataframe that at least contains one character column.
space	decontextualized/static space with a column called "words" and the semantic representations are in columns called Dim1, Dim2 (or V1, V2, ...) and so on (from textSpace, which is not included in the current text package).
tk_df	default "null"; option to use either the "tk" of "df" space (if using textSpace, which has not been implemented yet).
aggregation_from_tokens_to_texts	method to aggregate semantic representation when there are more than a single word. (default is "mean"; see also "min" and "max", "concatenate" and "normalize")
dim_name	Boolean, if TRUE append the variable name after all variable-names in the output. (This differentiates between word embedding dimension names; e.g., Dim1_text_variable_name)
tolower	(boolean) Lower case input.

Value

A list with tibbles for each character variable. Each tibble comprises a column with the text, followed by columns representing the semantic representations of the text. The tibbles are called the same as the original variable.

See Also

see [textEmbed](#)

textGeneration	<i>Predicts the words that will follow a specified text prompt. (experimental)</i>
----------------	--

Description

Predicts the words that will follow a specified text prompt. (experimental)

Usage

```

textGeneration(
  x,
  model = "gpt2",
  device = "cpu",
  tokenizer_parallelism = FALSE,
  logging_level = "warning",
  return_incorrect_results = FALSE,
  return_tensors = FALSE,
  return_text = TRUE,
  return_full_text = TRUE,
  clean_up_tokenization_spaces = FALSE,
  prefix = "",
  handle_long_generation = NULL,
  set_seed = 202208L
)

```

Arguments

x	(string) A variable or a tibble/dataframe with at least one character variable.
model	(string) Specification of a pre-trained language model that have been trained with an autoregressive language modeling objective, which includes the uni-directional models (e.g., gpt2).
device	(string) Device to use: 'cpu', 'gpu', or 'gpu:k' where k is a specific device number
tokenizer_parallelism	(boolean) If TRUE this will turn on tokenizer parallelism.
logging_level	(string) Set the logging level. Options (ordered from less logging to more logging): critical, error, warning, info, debug
return_incorrect_results	(boolean) Stop returning some incorrectly formatted/structured results. This setting does CANNOT evaluate the actual results (whether or not they make sense, exist, etc.). All it does is to ensure the returned results are formatted correctly (e.g., does the question-answering dictionary contain the key "answer", is sentiments from textClassify containing the labels "positive" and "negative").
return_tensors	(boolean) Whether or not the output should include the prediction tensors (as token indices).
return_text	(boolean) Whether or not the outputs should include the decoded text.
return_full_text	(boolean) If FALSE only the added text is returned, otherwise the full text is returned. (This setting is only meaningful if return_text is set to TRUE)
clean_up_tokenization_spaces	(boolean) Option to clean up the potential extra spaces in the returned text.
prefix	(string) Option to add a prefix to prompt.

handle_long_generation

By default, this function does not handle long generation (those that exceed the model maximum length).

set_seed

(Integer) Set seed. (more info :<https://github.com/huggingface/transformers/issues/14033#issuecomment-948385227>). This setting provides some ways to work around the problem: None: default way, where no particular strategy is applied. "hole": Truncates left of input, and leaves a gap that is wide enough to let generation happen. (this might truncate a lot of the prompt and not suitable when generation exceed the model capacity)

Value

A tibble with generated text.

See Also

see [textClassify](#), [textNER](#), [textSum](#), [textQA](#), [textTranslate](#)

Examples

```
# generated_text <- textGeneration("The meaning of life is")
# generated_text
```

textModelLayers

Get the number of layers in a given model.

Description

Get the number of layers in a given model.

Usage

```
textModelLayers(target_model)
```

Arguments

target_model (string) The name of the model to know the number of layers of.

Value

Number of layers.

See Also

see [textModels](#)

Examples

```
## Not run:  
textModelLayers(target_model = "bert-base-uncased")  
  
## End(Not run)
```

textModels	<i>Check downloaded, available models.</i>
------------	--

Description

Check downloaded, available models.

Usage

```
textModels()
```

Value

List of names of models and tokenizers

See Also

see [textModelsRemove](#)

Examples

```
## Not run:  
textModels()  
  
## End(Not run)
```

textModelsRemove	<i>Delete a specified model and model associated files.</i>
------------------	---

Description

Delete a specified model and model associated files.

Usage

```
textModelsRemove(target_model)
```

Arguments

target_model (string) The name of the model to be deleted.

Value

Confirmation whether the model has been deleted.

See Also

see [textModels](#)

Examples

```
## Not run:
textModelsRemove("name-of-model-to-delete")

## End(Not run)
```

textNER

Named Entity Recognition. (experimental)

Description

Named Entity Recognition. (experimental)

Usage

```
textNER(
  x,
  model = "dslim/bert-base-NER",
  device = "cpu",
  tokenizer_parallelism = FALSE,
  logging_level = "error",
  return_incorrect_results = FALSE,
  set_seed = 202208L
)
```

Arguments

x	(string) A variable or a tibble/dataframe with at least one character variable.
model	(string) Specification of a pre-trained language model for token classification that have been fine-tuned on a NER task (e.g., see "dslim/bert-base-NER"). Use for predicting the classes of tokens in a sequence: person, organisation, location or miscellaneous).
device	(string) Device to use: 'cpu', 'gpu', or 'gpu:k' where k is a specific device number
tokenizer_parallelism	(boolean) If TRUE this will turn on tokenizer parallelism.
logging_level	(string) Set the logging level. Options (ordered from less logging to more logging): critical, error, warning, info, debug

`return_incorrect_results` (boolean) Stop returning some incorrectly formatted/structured results. This setting does **CANNOT** evaluate the actual results (whether or not they make sense, exist, etc.). All it does is to ensure the returned results are formatted correctly (e.g., does the question-answering dictionary contain the key "answer", is sentiments from `textClassify` containing the labels "positive" and "negative").

`set_seed` (Integer) Set seed.

Value

A list with tibble(s) with NER classifications for each column.

See Also

see [textClassify](#), [textGeneration](#), [textNER](#), [textSum](#), [textQA](#), [textTranslate](#)

Examples

```
# ner_example <- textNER("Arnes plays football with Daniel")
# ner_example
```

textPCA	<i>Compute 2 PCA dimensions of the word embeddings for individual words.</i>
---------	--

Description

Compute 2 PCA dimensions of the word embeddings for individual words.

Usage

```
textPCA(words, word_types_embeddings = word_types_embeddings_df, seed = 1010)
```

Arguments

`words` Word or text variable to be plotted.

`word_types_embeddings` Word embeddings from `textEmbed` for individual words (i.e., decontextualized embeddings).

`seed` Set different seed.

Value

A dataframe with words, their frequency and two PCA dimensions from the `word_embeddings` for the individual words that is used for the plotting in the `textPCAPlot` function.

See Also

see [textPCAPlot](#)

Examples

```
## Not run:
# Data
df_for_plotting2d <- textPCA(
  words = Language_based_assessment_data_8$harmonywords,
  word_types_embeddings = word_embeddings_4$word_types
)
df_for_plotting2d

## End(Not run)
```

textPCAPlot

Plot words according to 2-D plot from 2 PCA components.

Description

Plot words according to 2-D plot from 2 PCA components.

Usage

```
textPCAPlot(
  word_data,
  min_freq_words_test = 1,
  plot_n_word_extreme = 5,
  plot_n_word_frequency = 5,
  plot_n_words_middle = 5,
  titles_color = "#61605e",
  title_top = "Principal Component (PC) Plot",
  x_axes_label = "PC1",
  y_axes_label = "PC2",
  scale_x_axes_lim = NULL,
  scale_y_axes_lim = NULL,
  word_font = NULL,
  bivariate_color_codes = c("#398CF9", "#60A1F7", "#5dc688", "#e07f6a", "#EAEAEA",
    "#40DD52", "#FF0000", "#EA7467", "#85DB8E"),
  word_size_range = c(3, 8),
  position_jitter_high = 0,
  position_jitter_width = 0.03,
  point_size = 0.5,
  arrow_transparency = 0.1,
  points_without_words_size = 0.2,
  points_without_words_alpha = 0.2,
  legend_title = "PC",
```

```

    legend_x_axes_label = "PC1",
    legend_y_axes_label = "PC2",
    legend_x_position = 0.02,
    legend_y_position = 0.02,
    legend_h_size = 0.2,
    legend_w_size = 0.2,
    legend_title_size = 7,
    legend_number_size = 2,
    seed = 1002
  )

```

Arguments

word_data Dataframe from textPCA

min_freq_words_test Select words to significance test that have occurred at least `min_freq_words_test` (default = 1).

plot_n_word_extreme Number of words that are extreme on Supervised Dimension Projection per dimension. (i.e., even if not significant; per dimensions, where duplicates are removed).

plot_n_word_frequency Number of words based on being most frequent. (i.e., even if not significant).

plot_n_words_middle Number of words plotted that are in the middle in Supervised Dimension Projection score (i.e., even if not significant; per dimensions, where duplicates are removed).

titles_color Color for all the titles (default: "#61605e")

title_top Title (default " ")

x_axes_label Label on the x-axes.

y_axes_label Label on the y-axes.

scale_x_axes_lim Manually set the length of the x-axes (default = NULL, which uses `ggplot2::scale_x_continuous(limits = scale_x_axes_lim)`; change e.g., by trying `c(-5, 5)`).

scale_y_axes_lim Manually set the length of the y-axes (default = NULL; which uses `ggplot2::scale_y_continuous(limits = scale_y_axes_lim)`; change e.g., by trying `c(-5, 5)`).

word_font Font type (default: NULL).

bivariate_color_codes The different colors of the words (default: `c("#398CF9", "#60A1F7", "#5dc688", "#e07f6a", "#EAEAEA", "#40DD52", "#FF0000", "#EA7467", "#85DB8E")`).

word_size_range Vector with minimum and maximum font size (default: `c(3, 8)`).

position_jitter_hight Jitter height (default: `.0`).

`position_jitter_width` Jitter width (default: .03).
`point_size` Size of the points indicating the words' position (default: 0.5).
`arrow_transparency` Transparency of the lines between each word and point (default: 0.1).
`points_without_words_size` Size of the points not linked with a words (default is to not show it, i.e., 0).
`points_without_words_alpha` Transparency of the points not linked with a words (default is to not show it, i.e., 0).
`legend_title` Title on the color legend (default: "(PCA)").
`legend_x_axes_label` Label on the color legend (default: "(x)").
`legend_y_axes_label` Label on the color legend (default: "(y)").
`legend_x_position` Position on the x coordinates of the color legend (default: 0.02).
`legend_y_position` Position on the y coordinates of the color legend (default: 0.05).
`legend_h_size` Height of the color legend (default 0.15).
`legend_w_size` Width of the color legend (default 0.15).
`legend_title_size` Font size (default: 7).
`legend_number_size` Font size of the values in the legend (default: 2).
`seed` Set different seed.

Value

A 1- or 2-dimensional word plot, as well as tibble with processed data used to plot..

See Also

see [textPCA](#)

Examples

```

# The test-data included in the package is called: DP_projections_HILS_SWLS_100

# Supervised Dimension Projection Plot
principle_component_plot_projection <- textPCAPlot(PC_projections_satisfactionwords_40)
principle_component_plot_projection

names(DP_projections_HILS_SWLS_100)

```

textPlot	<i>Plot words from textProjection() or textWordPrediction().</i>
----------	--

Description

Plot words from textProjection() or textWordPrediction().

Usage

```
textPlot(
  word_data,
  k_n_words_to_test = FALSE,
  min_freq_words_test = 1,
  min_freq_words_plot = 1,
  plot_n_words_square = 3,
  plot_n_words_p = 5,
  plot_n_word_extreme = 5,
  plot_n_word_frequency = 5,
  plot_n_words_middle = 5,
  titles_color = "#61605e",
  y_axes = FALSE,
  p_alpha = 0.05,
  p_adjust_method = "none",
  title_top = "Supervised Dimension Projection",
  x_axes_label = "Supervised Dimension Projection (SDP)",
  y_axes_label = "Supervised Dimension Projection (SDP)",
  scale_x_axes_lim = NULL,
  scale_y_axes_lim = NULL,
  word_font = NULL,
  bivariate_color_codes = c("#398CF9", "#60A1F7", "#5dc688", "#e07f6a", "#EAEAEA",
    "#40DD52", "#FF0000", "#EA7467", "#85DB8E"),
  word_size_range = c(3, 8),
  position_jitter_hight = 0,
  position_jitter_width = 0.03,
  point_size = 0.5,
  arrow_transparency = 0.1,
  points_without_words_size = 0.2,
  points_without_words_alpha = 0.2,
  legend_title = "SDP",
  legend_x_axes_label = "x",
  legend_y_axes_label = "y",
  legend_x_position = 0.02,
  legend_y_position = 0.02,
  legend_h_size = 0.2,
  legend_w_size = 0.2,
  legend_title_size = 7,
  legend_number_size = 2,
```

```

group_embeddings1 = FALSE,
group_embeddings2 = FALSE,
projection_embedding = FALSE,
aggregated_point_size = 0.8,
aggregated_shape = 8,
aggregated_color_G1 = "black",
aggregated_color_G2 = "black",
projection_color = "blue",
seed = 1005,
explore_words = NULL,
explore_words_color = "#ad42f5",
explore_words_point = "ALL_1",
explore_words_aggregation = "mean",
remove_words = NULL,
n_contrast_group_color = NULL,
n_contrast_group_remove = FALSE,
space = NULL,
scaling = FALSE
)

```

Arguments

word_data Dataframe from textProjection

k_n_words_to_test
Select the k most frequent words to significance test ($k = \sqrt{100 \cdot N}$; $N =$ number of participant responses). Default = TRUE.

min_freq_words_test
Select words to significance test that have occurred at least `min_freq_words_test` (default = 1).

min_freq_words_plot
Select words to plot that has occurred at least `min_freq_words_plot` times.

plot_n_words_square
Select number of significant words in each square of the figure to plot. The significant words, in each square is selected according to most frequent words.

plot_n_words_p Number of significant words to plot on each(positive and negative) side of the x-axes and y-axes, (where duplicates are removed); selects first according to lowest p-value and then according to frequency. Hence, on a two dimensional plot it is possible that `plot_n_words_p = 1` yield 4 words.

plot_n_word_extreme
Number of words that are extreme on Supervised Dimension Projection per dimension. (i.e., even if not significant; per dimensions, where duplicates are removed).

plot_n_word_frequency
Number of words based on being most frequent. (i.e., even if not significant).

plot_n_words_middle
Number of words plotted that are in the middle in Supervised Dimension Projection score (i.e., even if not significant; per dimensions, where duplicates are removed).

titles_color	Color for all the titles (default: "#61605e")
y_axes	If TRUE, also plotting on the y-axes (default is FALSE). Also plotting on y-axes produces a two dimension 2-dimensional plot, but the textProjection function has to have had a variable on the y-axes.
p_alpha	Alpha (default = .05).
p_adjust_method	Method to adjust/correct p-values for multiple comparisons (default = "holm"; see also "none", "hochberg", "hommel", "bonferroni", "BH", "BY", "fdr").
title_top	Title (default " ")
x_axes_label	Label on the x-axes.
y_axes_label	Label on the y-axes.
scale_x_axes_lim	Manually set the length of the x-axes (default = NULL, which uses ggplot2::scale_x_continuous(limits = scale_x_axes_lim); change e.g., by trying c(-5, 5)).
scale_y_axes_lim	Manually set the length of the y-axes (default = NULL; which uses ggplot2::scale_y_continuous(limits = scale_y_axes_lim); change e.g., by trying c(-5, 5)).
word_font	Font type (default: NULL).
bivariate_color_codes	The different colors of the words. Note that, at the moment, two squares should not have the exact same colour-code because the numbers within the squares of the legend will then be aggregated (and show the same, incorrect value). (default: c("#398CF9", "#60A1F7", "#5dc688", "#e07f6a", "#EAEAEA", "#40DD52", "#FF0000", "#EA7467", "#85DB8E")).
word_size_range	Vector with minimum and maximum font size (default: c(3, 8)).
position_jitter_hight	Jitter height (default: .0).
position_jitter_width	Jitter width (default: .03).
point_size	Size of the points indicating the words' position (default: 0.5).
arrow_transparency	Transparency of the lines between each word and point (default: 0.1).
points_without_words_size	Size of the points not linked with a words (default is to not show it, i.e., 0).
points_without_words_alpha	Transparency of the points not linked with a words (default is to not show it, i.e., 0).
legend_title	Title on the color legend (default: "(SDP)").
legend_x_axes_label	Label on the color legend (default: "(x)").
legend_y_axes_label	Label on the color legend (default: "(y)").

`legend_x_position` Position on the x coordinates of the color legend (default: 0.02).
`legend_y_position` Position on the y coordinates of the color legend (default: 0.05).
`legend_h_size` Height of the color legend (default 0.15).
`legend_w_size` Width of the color legend (default 0.15).
`legend_title_size` Font size (default: 7).
`legend_number_size` Font size of the values in the legend (default: 2).
`group_embeddings1` Shows a point representing the aggregated word embedding for group 1 (default = FALSE).
`group_embeddings2` Shows a point representing the aggregated word embedding for group 2 (default = FALSE).
`projection_embedding` Shows a point representing the aggregated direction embedding (default = FALSE).
`aggregated_point_size` Size of the points representing the `group_embeddings1`, `group_embeddings2` and `projection_embedding`.
`aggregated_shape` Shape type of the points representing the `group_embeddings1`, `group_embeddings2` and `projection_embedding`.
`aggregated_color_G1` Color
`aggregated_color_G2` Color
`projection_color` Color
`seed` Set different seed.
`explore_words` Explore where specific words are positioned in the embedding space. For example, `c("happy content", "sad down")`.
`explore_words_color` Specify the color(s) of the words being explored. For example `c("#ad42f5", "green")`
`explore_words_point` Specify the names of the point for the aggregated word embeddings of all the explored words.
`explore_words_aggregation` Specify how to aggregate the word embeddings of the explored words.
`remove_words` manually remove words from the plot (which is done just before the words are plotted so that the `remove_words` are part of previous counts/analyses).
`n_contrast_group_color` Set color to words that have higher frequency (N) on the other opposite side of its dot product projection (default = NULL).

n_contrast_group_remove	Remove words that have higher frequency (N) on the other opposite side of its dot product projection (default = FALSE).
space	Provide a semantic space if using static embeddings and wanting to explore words.
scaling	Scaling word embeddings before aggregation.

Value

A 1- or 2-dimensional word plot, as well as tibble with processed data used to plot.

See Also

see [textProjection](#)

Examples

```
# The test-data included in the package is called: DP_projections_HILS_SWLS_100

# Supervised Dimension Projection Plot
plot_projection <- textPlot(
  word_data = DP_projections_HILS_SWLS_100,
  k_n_words_to_test = FALSE,
  min_freq_words_test = 1,
  plot_n_words_square = 3,
  plot_n_words_p = 3,
  plot_n_word_extreme = 1,
  plot_n_word_frequency = 1,
  plot_n_words_middle = 1,
  y_axes = FALSE,
  p_alpha = 0.05,
  title_top = "Supervised Dimension Projection (SDP)",
  x_axes_label = "Low vs. High HILS score",
  y_axes_label = "Low vs. High SWLS score",
  p_adjust_method = "bonferroni",
  scale_y_axes_lim = NULL
)
plot_projection

names(DP_projections_HILS_SWLS_100)
```

textPredict

Predict scores or classification from, e.g., textTrain.

Description

Predict scores or classification from, e.g., textTrain.

Usage

```
textPredict(
  model_info,
  word_embeddings,
  x_append = NULL,
  type = NULL,
  dim_names = TRUE,
  ...
)
```

Arguments

model_info	(model object) Model info (e.g., saved output from <code>textTrain</code> , <code>textTrainRegression</code> or <code>textRandomForest</code>).
word_embeddings	(tibble) Word embeddings
x_append	(tibble) Variables to be appended after the word embeddings (x).
type	(string) Type of prediction; e.g., "prob", "class".
dim_names	(boolean) Account for specific dimension names from <code>textEmbed()</code> (rather than generic names including Dim1, Dim2 etc.). If FALSE the models need to have been trained on word embeddings created with <code>dim_names</code> FALSE, so that embeddings were only called Dim1, Dim2 etc.
...	Setting from <code>stats::predict</code> can be called.

Value

Predicted scores from word embeddings.

See Also

see [textTrain](#) [textTrainLists](#) [textTrainRandomForest](#) [textSimilarityTest](#)

Examples

```
word_embeddings <- word_embeddings_4
ratings_data <- Language_based_assessment_data_8
```

`textPredictAll` *Predict from several models, selecting the correct input*

Description

Predict from several models, selecting the correct input

Usage

```
textPredictAll(models, word_embeddings, x_append = NULL, ...)
```

Arguments

models	Object containing several models.
word_embeddings	List of word embeddings (if using word embeddings from more than one text-variable use dim_names = TRUE throughout the pipeline).
x_append	A tibble/dataframe with additional variables used in the training of the models (optional).
...	Settings from textPredict.

Value

A tibble with predictions.

See Also

see [textPredict](#) and [textTrain](#)

Examples

```
# x <- Language_based_assessment_data_8[1:2, 1:2]
# word_embeddings_with_layers <- textEmbedLayersOutput(x, layers = 11:12)
```

textPredictTest	<i>Significance testing correlations If only y1 is provided a t-test is computed, between the absolute error from yhat1-y1 and yhat2-y1.</i>
-----------------	--

Description

If y2 is provided a bootstrapped procedure is used to compare the correlations between y1 and yhat1 versus y2 and yhat2. This is achieved by creating two distributions of correlations using bootstrapping; and then finally compute the distributions overlap.

Usage

```
textPredictTest(
  y1,
  y2 = NULL,
  yhat1,
  yhat2,
  paired = TRUE,
  bootstraps_times = 10000,
  seed = 6134,
  ...
)
```

Arguments

y1	The observed scores (i.e., what was used to predict when training a model).
y2	The second observed scores (default = NULL; i.e., for when comparing models that are predicting different outcomes. In this case a bootstrap procedure is used to create two distributions of correlations that are compared (see description above).
yhat1	The predicted scores from model 1.
yhat2	The predicted scores from model 2 that will be compared with model 1.
paired	Paired test or not in stats::t.test (default TRUE).
bootstraps_times	Number of bootstraps (when providing y2).
seed	Set different seed.
...	Settings from stats::t.test or overlapping::overlap (e.g., plot = TRUE).

Value

Comparison of correlations either a t-test or the overlap of a bootstrapped procedure (see \$OV).

See Also

see [textTrain](#) [textPredict](#)

Examples

```
# Example random data
y1 <- runif(10)
yhat1 <- runif(10)
y2 <- runif(10)
yhat2 <- runif(10)

boot_test <- textPredictTest(y1, yhat1, y2, yhat2, bootstraps_times = 10)
```

textProjection	<i>Compute Supervised Dimension Projection and related variables for plotting words.</i>
----------------	--

Description

Compute Supervised Dimension Projection and related variables for plotting words.

Usage

```

textProjection(
  words,
  word_embeddings,
  word_types_embeddings,
  x,
  y = NULL,
  pca = NULL,
  aggregation = "mean",
  split = "quartile",
  word_weight_power = 1,
  min_freq_words_test = 0,
  mean_centering = FALSE,
  mean_centering2 = FALSE,
  Npermutations = 10000,
  n_per_split = 50000,
  seed = 1003
)

```

Arguments

words	Word or text variable to be plotted.
word_embeddings	Word embeddings from textEmbed for the words to be plotted (i.e., the aggregated word embeddings for the "words" parameter).
word_types_embeddings	Word embeddings from textEmbed for individual words (i.e., decontextualized embeddings).
x	Numeric variable that the words should be plotted according to on the x-axes.
y	Numeric variable that the words should be plotted according to on the y-axes (y=NULL).
pca	Number of PCA dimensions applied to the word embeddings in the beginning of the function. A number below 1 takes out % of variance; An integer specify number of components to extract. (default is NULL as this setting has not yet been evaluated).
aggregation	Method to aggregate the word embeddings (default = "mean"; see also "min", "max", and "[CLS]").
split	Method to split the axes (default = "quartile" involving selecting lower and upper quartile; see also "mean"). However, if the variable is only containing two different values (i.e., being dichotomous) mean split is used.
word_weight_power	Compute the power of the frequency of the words and multiply the word embeddings with this in the computation of aggregated word embeddings for group low (1) and group high (2). This increases the weight of more frequent words.

min_freq_words_test	Option to select words that have occurred a specified number of times (default = 0); when creating the Supervised Dimension Projection line (i.e., single words receive Supervised Dimension Projection and p-value).
mean_centering	Boolean; separately mean centering the Group 1 split aggregation embedding, and the Group 2 split aggregation embedding
mean_centering2	Boolean; separately mean centering the G1 and G2 split aggregation embeddings
Npermutations	Number of permutations in the creation of the null distribution.
n_per_split	Setting to split Npermutations to avoid reaching computer memory limits; set it lower than Npermutations <- and the higher it is set the faster the computation completes, but too high may lead to abortion.
seed	Set different seed.

Value

A dataframe with variables (e.g., including Supervised Dimension Projection, frequencies, p-values) for the individual words that is used for the plotting in the textProjectionPlot function.

Examples

```
# Data
# Pre-processing data for plotting
## Not run:
df_for_plotting <- textProjection(
  words = Language_based_assessment_data_8$harmonywords,
  word_embeddings = word_embeddings_4$texts$harmonywords,
  word_types_embeddings = word_embeddings_4$word_types,
  x = Language_based_assessment_data_8$hilstotal,
  split = "mean",
  Npermutations = 10,
  n_per_split = 1
)
df_for_plotting

## End(Not run)
#' @seealso see \code{\link{textProjectionPlot}}
```

textProjectionPlot *Plot words according to Supervised Dimension Projection.*

Description

Plot words according to Supervised Dimension Projection.

Usage

```
textProjectionPlot(  
  word_data,  
  k_n_words_to_test = FALSE,  
  min_freq_words_test = 1,  
  min_freq_words_plot = 1,  
  plot_n_words_square = 3,  
  plot_n_words_p = 5,  
  plot_n_word_extreme = 5,  
  plot_n_word_frequency = 5,  
  plot_n_words_middle = 5,  
  titles_color = "#61605e",  
  y_axes = FALSE,  
  p_alpha = 0.05,  
  p_adjust_method = "none",  
  title_top = "Supervised Dimension Projection",  
  x_axes_label = "Supervised Dimension Projection (SDP)",  
  y_axes_label = "Supervised Dimension Projection (SDP)",  
  scale_x_axes_lim = NULL,  
  scale_y_axes_lim = NULL,  
  word_font = NULL,  
  bivariate_color_codes = c("#398CF9", "#60A1F7", "#5dc688", "#e07f6a", "#EAEAEA",  
    "#40DD52", "#FF0000", "#EA7467", "#85DB8E"),  
  word_size_range = c(3, 8),  
  position_jitter_hight = 0,  
  position_jitter_width = 0.03,  
  point_size = 0.5,  
  arrow_transparency = 0.1,  
  points_without_words_size = 0.2,  
  points_without_words_alpha = 0.2,  
  legend_title = "SDP",  
  legend_x_axes_label = "x",  
  legend_y_axes_label = "y",  
  legend_x_position = 0.02,  
  legend_y_position = 0.02,  
  legend_h_size = 0.2,  
  legend_w_size = 0.2,  
  legend_title_size = 7,  
  legend_number_size = 2,  
  group_embeddings1 = FALSE,  
  group_embeddings2 = FALSE,  
  projection_embedding = FALSE,  
  aggregated_point_size = 0.8,  
  aggregated_shape = 8,  
  aggregated_color_G1 = "black",  
  aggregated_color_G2 = "black",  
  projection_color = "blue",  
  seed = 1005,
```

```

explore_words = NULL,
explore_words_color = "#ad42f5",
explore_words_point = "ALL_1",
explore_words_aggregation = "mean",
remove_words = NULL,
n_contrast_group_color = NULL,
n_contrast_group_remove = FALSE,
space = NULL,
scaling = FALSE
)

```

Arguments

<code>word_data</code>	Dataframe from textProjection
<code>k_n_words_to_test</code>	Select the k most frequent words to significance test ($k = \sqrt{100 \cdot N}$; N = number of participant responses). Default = TRUE.
<code>min_freq_words_test</code>	Select words to significance test that have occurred at least <code>min_freq_words_test</code> (default = 1).
<code>min_freq_words_plot</code>	Select words to plot that has occurred at least <code>min_freq_words_plot</code> times.
<code>plot_n_words_square</code>	Select number of significant words in each square of the figure to plot. The significant words, in each square is selected according to most frequent words.
<code>plot_n_words_p</code>	Number of significant words to plot on each(positive and negative) side of the x-axes and y-axes, (where duplicates are removed); selects first according to lowest p-value and then according to frequency. Hence, on a two dimensional plot it is possible that <code>plot_n_words_p = 1</code> yield 4 words.
<code>plot_n_word_extreme</code>	Number of words that are extreme on Supervised Dimension Projection per dimension. (i.e., even if not significant; per dimensions, where duplicates are removed).
<code>plot_n_word_frequency</code>	Number of words based on being most frequent. (i.e., even if not significant).
<code>plot_n_words_middle</code>	Number of words plotted that are in the middle in Supervised Dimension Projection score (i.e., even if not significant; per dimensions, where duplicates are removed).
<code>titles_color</code>	Color for all the titles (default: "#61605e")
<code>y_axes</code>	If TRUE, also plotting on the y-axes (default is FALSE). Also plotting on y-axes produces a two dimension 2-dimensional plot, but the textProjection function has to have had a variable on the y-axes.
<code>p_alpha</code>	Alpha (default = .05).
<code>p_adjust_method</code>	Method to adjust/correct p-values for multiple comparisons (default = "holm"; see also "none", "hochberg", "hommel", "bonferroni", "BH", "BY", "fdr").

title_top	Title (default " ")
x_axes_label	Label on the x-axes.
y_axes_label	Label on the y-axes.
scale_x_axes_lim	Manually set the length of the x-axes (default = NULL, which uses <code>ggplot2::scale_x_continuous(limits = scale_x_axes_lim)</code> ; change e.g., by trying <code>c(-5, 5)</code>).
scale_y_axes_lim	Manually set the length of the y-axes (default = NULL; which uses <code>ggplot2::scale_y_continuous(limits = scale_y_axes_lim)</code> ; change e.g., by trying <code>c(-5, 5)</code>).
word_font	Font type (default: NULL).
bivariate_color_codes	The different colors of the words. Note that, at the moment, two squares should not have the exact same colour-code because the numbers within the squares of the legend will then be aggregated (and show the same, incorrect value). (default: <code>c("#398CF9", "#60A1F7", "#5dc688", "#e07f6a", "#EAEAEA", "#40DD52", "#FF0000", "#EA7467", "#85DB8E")</code>).
word_size_range	Vector with minimum and maximum font size (default: <code>c(3, 8)</code>).
position_jitter_hight	Jitter height (default: <code>.0</code>).
position_jitter_width	Jitter width (default: <code>.03</code>).
point_size	Size of the points indicating the words' position (default: <code>0.5</code>).
arrow_transparency	Transparency of the lines between each word and point (default: <code>0.1</code>).
points_without_words_size	Size of the points not linked with a words (default is to not show it, i.e., <code>0</code>).
points_without_words_alpha	Transparency of the points not linked with a words (default is to not show it, i.e., <code>0</code>).
legend_title	Title on the color legend (default: <code>"(SDP)"</code>).
legend_x_axes_label	Label on the color legend (default: <code>"(x)"</code>).
legend_y_axes_label	Label on the color legend (default: <code>"(y)"</code>).
legend_x_position	Position on the x coordinates of the color legend (default: <code>0.02</code>).
legend_y_position	Position on the y coordinates of the color legend (default: <code>0.05</code>).
legend_h_size	Height of the color legend (default <code>0.15</code>).
legend_w_size	Width of the color legend (default <code>0.15</code>).
legend_title_size	Font size (default: <code>7</code>).

<code>legend_number_size</code>	Font size of the values in the legend (default: 2).
<code>group_embeddings1</code>	Shows a point representing the aggregated word embedding for group 1 (default = FALSE).
<code>group_embeddings2</code>	Shows a point representing the aggregated word embedding for group 2 (default = FALSE).
<code>projection_embedding</code>	Shows a point representing the aggregated direction embedding (default = FALSE).
<code>aggregated_point_size</code>	Size of the points representing the <code>group_embeddings1</code> , <code>group_embeddings2</code> and <code>projection_embedding</code>
<code>aggregated_shape</code>	Shape type of the points representing the <code>group_embeddings1</code> , <code>group_embeddings2</code> and <code>projection_embedding</code>
<code>aggregated_color_G1</code>	Color
<code>aggregated_color_G2</code>	Color
<code>projection_color</code>	Color
<code>seed</code>	Set different seed.
<code>explore_words</code>	Explore where specific words are positioned in the embedding space. For example, <code>c("happy content", "sad down")</code> .
<code>explore_words_color</code>	Specify the color(s) of the words being explored. For example <code>c("#ad42f5", "green")</code>
<code>explore_words_point</code>	Specify the names of the point for the aggregated word embeddings of all the explored words.
<code>explore_words_aggregation</code>	Specify how to aggregate the word embeddings of the explored words.
<code>remove_words</code>	manually remove words from the plot (which is done just before the words are plotted so that the <code>remove_words</code> are part of previous counts/analyses).
<code>n_contrast_group_color</code>	Set color to words that have higher frequency (N) on the other opposite side of its dot product projection (default = NULL).
<code>n_contrast_group_remove</code>	Remove words that have higher frequency (N) on the other opposite side of its dot product projection (default = FALSE).
<code>space</code>	Provide a semantic space if using static embeddings and wanting to explore words.
<code>scaling</code>	Scaling word embeddings before aggregation.

Value

A 1- or 2-dimensional word plot, as well as tibble with processed data used to plot.

See Also

see [textProjection](#)

Examples

```
# The test-data included in the package is called: DP_projections_HILS_SWLS_100

# Supervised Dimension Projection Plot
plot_projection <- textProjectionPlot(
  word_data = DP_projections_HILS_SWLS_100,
  k_n_words_to_test = FALSE,
  min_freq_words_test = 1,
  plot_n_words_square = 3,
  plot_n_words_p = 3,
  plot_n_word_extreme = 1,
  plot_n_word_frequency = 1,
  plot_n_words_middle = 1,
  y_axes = FALSE,
  p_alpha = 0.05,
  title_top = "Supervised Dimension Projection (SDP)",
  x_axes_label = "Low vs. High HILS score",
  y_axes_label = "Low vs. High SWLS score",
  p_adjust_method = "bonferroni",
  scale_y_axes_lim = NULL
)
plot_projection

names(DP_projections_HILS_SWLS_100)
```

textQA

Question Answering. (experimental)

Description

Question Answering. (experimental)

Usage

```
textQA(
  question,
  context,
  model = "",
  device = "cpu",
  tokenizer_parallelism = FALSE,
  logging_level = "warning",
```

```

return_incorrect_results = FALSE,
top_k = 1L,
doc_stride = 128L,
max_answer_len = 15L,
max_seq_len = 384L,
max_question_len = 64L,
handle_impossible_answer = FALSE,
set_seed = 202208L
)

```

Arguments

question	(string) A question
context	(string) The context(s) where the model will look for the answer.
model	(string) HuggingFace name of a pre-trained language model that have been fine-tuned on a question answering task.
device	(string) Device to use: 'cpu', 'gpu', or 'gpu:k' where k is a specific device number
tokenizer_parallelism	(boolean) If TRUE this will turn on tokenizer parallelism.
logging_level	(string) Set the logging level. Options (ordered from less logging to more logging): critical, error, warning, info, debug
return_incorrect_results	(boolean) Stop returning some incorrectly formatted/structured results. This setting does CANNOT evaluate the actual results (whether or not they make sense, exist, etc.). All it does is to ensure the returned results are formatted correctly (e.g., does the question-answering dictionary contain the key "answer", is sentiments from textClassify containing the labels "positive" and "negative").
top_k	(integer) (int) Indicates number of possible answer span(s) to get from the model output.
doc_stride	(integer) If the context is too long to fit with the question for the model, it will be split into overlapping chunks. This setting controls the overlap size.
max_answer_len	(integer) Max answer size to be extracted from the model's output.
max_seq_len	(integer) The max total sentence length (context + question) in tokens of each chunk passed to the model. If needed, the context is split in chunks (using doc_stride as overlap).
max_question_len	(integer) The max question length after tokenization. It will be truncated if needed.
handle_impossible_answer	(boolean) Whether or not impossible is accepted as an answer.
set_seed	(Integer) Set seed.

Value

Answers.

See Also

see [textClassify](#), [textGeneration](#), [textNER](#), [textSum](#), [textQA](#), [textTranslate](#)

Examples

```
# qa_examples <- textQA(question = "Which colour have trees?",
#   context = "Trees typically have leaves, are mostly green and like water.")
```

texttrpp_initialize *Initialize text required python packages*

Description

Initialize text required python packages to call from R.

Usage

```
texttrpp_initialize(
  python_executable = NULL,
  virtualenv = NULL,
  condaenv = "texttrpp_condaenv",
  ask = FALSE,
  refresh_settings = FALSE,
  save_profile = FALSE,
  check_env = TRUE,
  textEmbed_test = FALSE,
  prompt = TRUE
)
```

Arguments

python_executable	the full path to the Python executable, for which text required python packages is installed.
virtualenv	set a path to the Python virtual environment with text required python packages installed Example: virtualenv = "~/myenv"
condaenv	set a path to the anaconda virtual environment with text required python packages installed Example: condaenv = "myenv"
ask	logical; if FALSE, use the first text required python packages installation found; if TRUE, list available text required python packages installations and prompt the user for which to use. If another (e.g. python_executable) is set, then this value will always be treated as FALSE.
refresh_settings	logical; if TRUE, text will ignore the saved settings in the profile and initiate a search of new settings.

save_profile	logical; if TRUE, the current text required python packages setting will be saved for the future use.
check_env	logical; check whether conda/virtual environment generated by textrpp_install() exists
textEmbed_test	logical; Test whether function (textEmbed) that requires python packages works.
prompt	logical; asking whether user wants to set the environment as default.

textrpp_install	<i>Install text required python packages in conda or virtualenv environment</i>
-----------------	---

Description

Install text required python packages (rpp) in a self-contained environment. For macOS and Linux-based systems, this will also install Python itself via a "miniconda" environment, for textrpp_install. Alternatively, an existing conda installation may be used, by specifying its path. The default setting of "auto" will locate and use an existing installation automatically, or download and install one if none exists.

For Windows, automatic installation of miniconda installation is not currently available, so the user will need to **miniconda (or Anaconda) manually**.

If you wish to install Python in a "virtualenv", use the textrpp_install_virtualenv function. It requires that you have a python version and path to it (such as "/usr/local/bin/python3.9" for Mac and Linux.).

Usage

```
textrpp_install(
  conda = "auto",
  update_conda = FALSE,
  force_conda = FALSE,
  rpp_version = "rpp_version_system_specific_defaults",
  python_version = "python_version_system_specific_defaults",
  envname = "textrpp_condaenv",
  pip = TRUE,
  python_path = NULL,
  prompt = TRUE
)

textrpp_install_virtualenv(
  rpp_version = c("torch==1.11.0", "transformers==4.19.2", "numpy", "nltk"),
  python_path = "/usr/local/bin/python3.9",
  pip_version = NULL,
  envname = "textrpp_virtualenv",
  prompt = TRUE
)
```

Arguments

conda	character; path to conda executable. Default "auto" which automatically find the path
update_conda	Boolean; update to the latest version of Miniconda after install? (should be combined with force_conda = TRUE)
force_conda	Boolean; force re-installation if Miniconda is already installed at the requested path?
rpp_version	character; default is "rpp_version_system_specific_defaults", because different systems require different combinations of python version and packages. It is also possible to specify your own, such as c('torch==0.4.1', 'transformers==3.3.1').
python_version	character; default is "python_version_system_specific_defaults". You can specify your Python version for the condaenv yourself. installation.
envname	character; name of the conda-environment to install text required python packages. Default is "texttrpp_condaenv".
pip	TRUE to use pip for installing rpp If FALSE, conda package manager with conda-forge channel will be used for installing rpp.
python_path	character; path to Python in virtualenv installation
prompt	logical; ask whether to proceed during the installation
pip_version	character;

Examples

```
## Not run:
# install text required python packages in a miniconda environment (macOS and Linux)
texttrpp_install(prompt = FALSE)

# install text required python packages to an existing conda environment
texttrpp_install(conda = "~/anaconda/bin/")

## End(Not run)
## Not run:
# install text required python packages in a virtual environment
texttrpp_install_virtualenv()

## End(Not run)
```

texttrpp_uninstall *Uninstall texttrpp conda environment*

Description

Removes the conda environment created by texttrpp_install()

Usage

```
texttrpp_uninstall(conda = "auto", prompt = TRUE, envname = "texttrpp_condaenv")
```

Arguments

conda	path to conda executable, default to "auto" which automatically finds the path
prompt	logical; ask whether to proceed during the installation
envname	character; name of conda environment to remove

textSimilarity	<i>Compute the semantic similarity between two text variables.</i>
----------------	--

Description

Compute the semantic similarity between two text variables.

Usage

```
textSimilarity(x, y, method = "cosine", center = TRUE, scale = FALSE)
```

Arguments

x	Word embeddings from textEmbed.
y	Word embeddings from textEmbed.
method	Character string describing type of measure to be computed. Default is "cosine" (see also "spearmen", "pearson" as well as measures from textDistance() (which here is computed as 1 - textDistance) including "euclidean", "maximum", "manhattan", "canberra", "binary" and "minkowski").
center	(boolean; from base::scale) If center is TRUE then centering is done by subtracting the column means (omitting NAs) of x from their corresponding columns, and if center is FALSE, no centering is done.
scale	(boolean; from base::scale) If scale is TRUE then scaling is done by dividing the (centered) columns of x by their standard deviations if center is TRUE, and the root mean square otherwise.

Value

A vector comprising semantic similarity scores.

See Also

see [textDistance](#), [textSimilarityNorm](#) and [textSimilarityTest](#)

Examples

```
library(dplyr)
similarity_scores <- textSimilarity(
  x = word_embeddings_4$texts$harmonytext,
  y = word_embeddings_4$texts$satisfactiontext
)
comment(similarity_scores)
```

textSimilarityMatrix *Compute semantic similarity scores between all combinations in a word embedding*

Description

Compute semantic similarity scores between all combinations in a word embedding

Usage

```
textSimilarityMatrix(x, method = "cosine", center = TRUE, scale = FALSE)
```

Arguments

x	Word embeddings from textEmbed.
method	Character string describing type of measure to be computed. Default is "cosine" (see also "spearmen", "pearson" as well as measures from textDistance() (which here is computed as 1 - textDistance) including "euclidean", "maximum", "manhattan", "canberra", "binary" and "minkowski").
center	(boolean; from base::scale) If center is TRUE then centering is done by subtracting the column means (omitting NAs) of x from their corresponding columns, and if center is FALSE, no centering is done.
scale	(boolean; from base::scale) If scale is TRUE then scaling is done by dividing the (centered) columns of x by their standard deviations if center is TRUE, and the root mean square otherwise.

Value

A matrix of semantic similarity scores

See Also

see [textSimilarityNorm](#) and [textSimilarityTest](#)

Examples

```
similarity_scores <- textSimilarityMatrix(word_embeddings_4$texts$harmonytext[1:3, ])  
round(similarity_scores, 3)
```

textSimilarityNorm	<i>Compute the semantic similarity between a text variable and a word norm (i.e., a text represented by one word embedding that represent a construct).</i>
--------------------	---

Description

Compute the semantic similarity between a text variable and a word norm (i.e., a text represented by one word embedding that represent a construct).

Usage

```
textSimilarityNorm(x, y, method = "cosine", center = TRUE, scale = FALSE)
```

Arguments

x	Word embeddings from textEmbed.
y	Word embedding from textEmbed (from only one text).
method	Character string describing type of measure to be computed. Default is "cosine" (see also "spearmen", "pearson" as well as measures from textDistance() (which here is computed as 1 - textDistance) including "euclidean", "maximum", "manhattan", "canberra", "binary" and "minkowski").
center	(boolean; from base::scale) If center is TRUE then centering is done by subtracting the column means (omitting NAs) of x from their corresponding columns, and if center is FALSE, no centering is done.
scale	(boolean; from base::scale) If scale is TRUE then scaling is done by dividing the (centered) columns of x by their standard deviations if center = TRUE, and the root mean square otherwise.

Value

A vector comprising semantic similarity scores.

See Also

see [textSimilarity](#) and [textSimilarityTest](#)

Examples

```
## Not run:
library(dplyr)
library(tibble)
harmonynorm <- c("harmony peace ")
satisfactionnorm <- c("satisfaction achievement")

norms <- tibble::tibble(harmonynorm, satisfactionnorm)
word_embeddings <- word_embeddings_4$texts
```

```

word_embeddings_wordnorm <- textEmbed(norms)
similarity_scores <- textSimilarityNorm(
  word_embeddings$harmonytext,
  word_embeddings_wordnorm$harmonynorm
)

## End(Not run)

```

textSimilarityTest *EXPERIMENTAL: Test whether there is a significant difference in meaning between two sets of texts (i.e., between their word embeddings).*

Description

EXPERIMENTAL: Test whether there is a significant difference in meaning between two sets of texts (i.e., between their word embeddings).

Usage

```

textSimilarityTest(
  x,
  y,
  similarity_method = "cosine",
  Npermutations = 10000,
  method = "paired",
  center = FALSE,
  scale = FALSE,
  alternative = "greater",
  output.permutations = TRUE,
  N_cluster_nodes = 1,
  seed = 1001
)

```

Arguments

x	Set of word embeddings from textEmbed.
y	Set of word embeddings from textEmbed.
similarity_method	Character string describing type of measure to be computed; default is "cosine" (see also measures from textDistance (here computed as 1 - textDistance()) including "euclidean", "maximum", "manhattan", "canberra", "binary" and "minkowski").
Npermutations	Number of permutations (default 10000).
method	Compute a "paired" or an "unpaired" test.
center	(boolean; from base::scale) If center is TRUE then centering is done by subtracting the column means (omitting NAs) of x from their corresponding columns, and if center is FALSE, no centering is done.

scale	(boolean; from base::scale) If scale is TRUE then scaling is done by dividing the (centered) columns of x by their standard deviations if center is TRUE, and the root mean square otherwise.
alternative	Use a two or one-sided test (select one of: "two_sided", "less", "greater").
output.permutations	If TRUE, returns permuted values in output.
N_cluster_nodes	Number of cluster nodes to use (more makes computation faster; see parallel package).
seed	Set different seed.

Value

A list with a p-value, similarity score estimate and permuted values if output.permutations=TRUE.

Examples

```
x <- word_embeddings_4$texts$harmonywords
y <- word_embeddings_4$texts$satisfactionwords
textSimilarityTest(x,
  y,
  method = "paired",
  Npermutations = 100,
  N_cluster_nodes = 1,
  alternative = "two_sided"
)
```

textSum	<i>Summarize texts. (experimental)</i>
---------	--

Description

Summarize texts. (experimental)

Usage

```
textSum(
  x,
  min_length = 10L,
  max_length = 20L,
  model = "t5-small",
  device = "cpu",
  tokenizer_parallelism = FALSE,
  logging_level = "warning",
  return_incorrect_results = FALSE,
  return_text = TRUE,
  return_tensors = FALSE,
```



```

    clean_up_tokenization_spaces = FALSE,
    set_seed = 202208L
  )

```

Arguments

x (string) A variable or a tibble/dataframe with at least one character variable.

min_length (explicit integer; e.g., 10L) The minimum number of tokens in the summed output.

max_length (explicit integer higher than min_length; e.g., 20L) The maximum number of tokens in the summed output.

model (string) Specification of a pre-trained language model that have been fine-tuned on a summarization task, such as 'bart-large-cnn', 't5-small', 't5-base', 't5-large', 't5-3b', 't5-11b'.

device (string) Device to use: 'cpu', 'gpu', or 'gpu:k' where k is a specific device number.

tokenizer_parallelism (boolean) If TRUE this will turn on tokenizer parallelism.

logging_level (string) Set the logging level. Options (ordered from less logging to more logging): critical, error, warning, info, debug

return_incorrect_results (boolean) Stop returning some incorrectly formatted/structured results. This setting does **CANNOT** evaluate the actual results (whether or not they make sense, exist, etc.). All it does is to ensure the returned results are formatted correctly (e.g., does the question-answering dictionary contain the key "answer", is sentiments from textClassify containing the labels "positive" and "negative").

return_text (boolean) Whether or not the outputs should include the decoded text.

return_tensors (boolean) Whether or not the output should include the prediction tensors (as token indices).

clean_up_tokenization_spaces (boolean) Option to clean up the potential extra spaces in the returned text.

set_seed (Integer) Set seed.

Value

A tibble with summed text(s).

See Also

see [textClassify](#), [textGeneration](#), [textNER](#), [textSum](#), [textQA](#), [textTranslate](#)

Examples

```

# sum_examples <- textSum(Language_based_assessment_data_8[1:2,1:2],
# min_length = 5L,
# max_length = 10L)

```

textTokenize	<i>Tokenize according to different huggingface transformers</i>
--------------	---

Description

Tokenize according to different huggingface transformers

Usage

```
textTokenize(
  texts,
  model = "bert-base-uncased",
  max_token_to_sentence = 4,
  device = "cpu",
  tokenizer_parallelism = FALSE,
  model_max_length = NULL,
  logging_level = "error"
)
```

Arguments

texts	A character variable or a tibble/dataframe with at least one character variable.
model	Character string specifying pre-trained language model (default 'bert-base-uncased'). For full list of options see pretrained models at HuggingFace . For example use "bert-base-multilingual-cased", "openai-gpt", "gpt2", "ctrl", "transfo-xl-wt103", "xlnet-base-cased", "xlm-mlm-enfr-1024", "distilbert-base-cased", "roberta-base", or "xlm-roberta-base".
max_token_to_sentence	(numeric) Maximum number of tokens in a string to handle before switching to embedding text sentence by sentence.
device	Name of device to use: 'cpu', 'gpu', or 'gpu:k' where k is a specific device number
tokenizer_parallelism	If TRUE this will turn on tokenizer parallelism. Default FALSE.
model_max_length	The maximum length (in number of tokens) for the inputs to the transformer model (default the value stored for the associated model).
logging_level	Set the logging level. Default: "warning". Options (ordered from less logging to more logging): critical, error, warning, info, debug

Value

Returns tokens according to specified huggingface transformer.

See Also

see [textEmbed](#)

Examples

```
# tokens <- textTokenize("hello are you?")
```

textTrain	<i>Train word embeddings to a numeric (ridge regression) or categorical (random forest) variable.</i>
-----------	---

Description

Train word embeddings to a numeric (ridge regression) or categorical (random forest) variable.

Usage

```
textTrain(x, y, force_train_method = "automatic", ...)
```

Arguments

x	Word embeddings from textEmbed (or textEmbedLayerAggregation). Can analyze several variables at the same time; but if training to several outcomes at the same time use a tibble within the list as input rather than just a tibble input (i.e., keep the name of the wordembedding).
y	Numeric variable to predict. Can be several; although then make sure to have them within a tibble (this is required even if it is only one outcome but several word embeddings variables).
force_train_method	default is "automatic", so if y is a factor random_forest is used, and if y is numeric ridge regression is used. This can be overridden using "regression" or "random_forest".
...	Arguments from textTrainRegression or textTrainRandomForest the textTrain function.

Value

A correlation between predicted and observed values; as well as a tibble of predicted values.

See Also

[textTrainRegression](#) [textTrainRandomForest](#) [textTrainLists](#) [textSimilarityTest](#)

Examples

```
## Not run:
results <- textTrain(
  x = word_embeddings_4$texts$harmonytext,
  y = Language_based_assessment_data_8$hilstotal
)

## End(Not run)
```

textTrainLists	<i>Individually trains word embeddings from several text variables to several numeric or categorical variables. It is possible to have word embeddings from one text variable and several numeric/categorical variables; or vice versa, word embeddings from several text variables to one numeric/categorical variable. It is not possible to mix numeric and categorical variables.</i>
----------------	---

Description

Individually trains word embeddings from several text variables to several numeric or categorical variables. It is possible to have word embeddings from one text variable and several numeric/categorical variables; or vice versa, word embeddings from several text variables to one numeric/categorical variable. It is not possible to mix numeric and categorical variables.

Usage

```
textTrainLists(
  x,
  y,
  force_train_method = "automatic",
  save_output = "all",
  method_cor = "pearson",
  eval_measure = "rmse",
  p_adjust_method = "holm",
  ...
)
```

Arguments

x	Word embeddings from textEmbed (or textEmbedLayerAggregation).
y	Tibble with several numeric or categorical variables to predict. Please note that you cannot mix numeric and categorical variables.
force_train_method	Default is "automatic"; see also "regression" and "random_forest".
save_output	Option not to save all output; default "all". see also "only_results" and "only_results_predictions".
method_cor	A character string describing type of correlation (default "Pearson").

eval_measure Type of evaluative measure to assess models on.
 p_adjust_method Method to adjust/correct p-values for multiple comparisons (default = "holm"; see also "none", "hochberg", "hommel", "bonferroni", "BH", "BY", "fdr").
 ... Arguments from textTrainRegression or textTrainRandomForest the textTrain function.

Value

Correlations between predicted and observed values.

See Also

see [textTrain](#) [textTrainRegression](#) [textTrainRandomForest](#)

Examples

```
## Not run:
word_embeddings <- word_embeddings_4$texts[1:2]
ratings_data <- Language_based_assessment_data_8[5:6]
results <- textTrainLists(
  x = word_embeddings,
  y = ratings_data
)
results
comment(results)

## End(Not run)
```

textTrainRandomForest *Train word embeddings to a categorical variable using random forest.*

Description

Train word embeddings to a categorical variable using random forrest.

Usage

```
textTrainRandomForest(
  x,
  y,
  x_append = NULL,
  cv_method = "validation_split",
  outside_folds = 10,
  outside_strata_y = "y",
  outside_breaks = 4,
  inside_folds = 3/4,
```

```

inside_strata_y = "y",
inside_breaks = 4,
mode_rf = "classification",
preprocess_step_center = FALSE,
preprocess_scale_center = FALSE,
preprocess_PCA = NA,
extremely_randomised_splitrule = "extratrees",
mtry = c(1, 10, 20, 40),
min_n = c(1, 10, 20, 40),
trees = c(1000),
eval_measure = "bal_accuracy",
model_description = "Consider writing a description of your model here",
multi_cores = "multi_cores_sys_default",
save_output = "all",
seed = 2020,
...
)

```

Arguments

x	Word embeddings from textEmbed.
y	Categorical variable to predict.
x_append	Variables to be appended after the word embeddings (x); if wanting to prepend them before the word embeddings use the option first = TRUE. If not wanting to train with word embeddings, set x = NULL.
cv_method	Cross-validation method to use within a pipeline of nested outer and inner loops of folds (see nested_cv in rsample). Default is using cv_folds in the outside folds and "validation_split" using rsample::validation_split in the inner loop to achieve a development and assessment set (note that for validation_split the inside_folds should be a proportion, e.g., inside_folds = 3/4); whereas "cv_folds" uses rsample::vfold_cv to achieve n-folds in both the outer and inner loops.
outside_folds	Number of folds for the outer folds (default = 10).
outside_strata_y	Variable to stratify according (default "y"; can also set to NULL).
outside_breaks	The number of bins wanted to stratify a numeric stratification variable in the outer cross-validation loop.
inside_folds	Number of folds for the inner folds (default = 3/4).
inside_strata_y	Variable to stratify according (default "y"; can also set to NULL).
inside_breaks	The number of bins wanted to stratify a numeric stratification variable in the inner cross-validation loop.
mode_rf	Default is "classification" ("regression" is not supported yet).
preprocess_step_center	normalizes dimensions to have a mean of zero; default is set to TRUE. For more info see (step_center in recipes).

preprocess_scale_center	normalize dimensions to have a standard deviation of one. For more info see (step_scale in recipes).
preprocess_PCA	Pre-processing threshold for PCA. Can select amount of variance to retain (e.g., .90 or as a grid c(0.80, 0.90)); or number of components to select (e.g., 10). Default is "min_halving", which is a function that selects the number of PCA components based on number of participants and feature (word embedding dimensions) in the data. The formula is: preprocess_PCA = round(max(min(number_features/2), number_participants/2), min(50, number_features))).
extremely_randomised_splitrule	default: "extratrees", which thus implement a random forest; can also select: NULL, "gini" or "hellinger"; if these are selected your mtry settings will be overridden (see Geurts et al. (2006) Extremely randomized trees for details; and see the ranger r-package for details on implementations).
mtry	hyper parameter that may be tuned; default:c(1, 20, 40),
min_n	hyper parameter that may be tuned; default: c(1, 20, 40)
trees	Number of trees to use (default 1000).
eval_measure	Measure to evaluate the models in order to select the best hyperparameters default "roc_auc"; see also "accuracy", "bal_accuracy", "sens", "spec", "precision", "kappa", "f_measure".
model_description	Text to describe your model (optional; good when sharing the model with others).
multi_cores	If TRUE it enables the use of multiple cores if the computer system allows for it (i.e., only on unix, not windows). Hence it makes the analyses considerably faster to run. Default is "multi_cores_sys_default", where it automatically uses TRUE for Mac and Linux and FALSE for Windows.
save_output	Option not to save all output; default "all". see also "only_results" and "only_results_predictions".
seed	Set different seed.
...	For example settings in yardstick::accuracy to set event_level (e.g., event_level = "second").

Value

A list with roc_curve_data, roc_curve_plot, truth and predictions, preprocessing_recipe, final_model, model_description chisq and fishers test as well as evaluation measures, e.g., including accuracy, f_meas and roc_auc (for details on these measures see the yardstick r-package documentation).

See Also

see [textTrainLists](#) [textSimilarityTest](#)

Examples

```
results <- textTrainRandomForest(
```

```

x = word_embeddings_4$texts$harmonywords,
y = as.factor(Language_based_assessment_data_8$gender),
trees = c(1000, 1500),
mtry = c(1), # this is short because of testing
min_n = c(1), # this is short because of testing
multi_cores = FALSE # This is FALSE due to CRAN testing and Windows machines.
)

```

textTrainRegression *Train word embeddings to a numeric variable.*

Description

Train word embeddings to a numeric variable.

Usage

```

textTrainRegression(
  x,
  y,
  x_append = NULL,
  cv_method = "validation_split",
  outside_folds = 10,
  outside_strata_y = "y",
  outside_breaks = 4,
  inside_folds = 3/4,
  inside_strata_y = "y",
  inside_breaks = 4,
  model = "regression",
  eval_measure = "default",
  preprocess_step_center = TRUE,
  preprocess_step_scale = TRUE,
  preprocess_PCA = NA,
  penalty = 10^seq(-16, 16),
  mixture = c(0),
  first_n_predictors = NA,
  impute_missing = FALSE,
  method_cor = "pearson",
  model_description = "Consider writing a description of your model here",
  multi_cores = "multi_cores_sys_default",
  save_output = "all",
  seed = 2020,
  ...
)

```


Arguments

x	Word embeddings from textEmbed (or textEmbedLayerAggregation). If several word embeddings are provided in a list they will be concatenated.
y	Numeric variable to predict.
x_append	Variables to be appended after the word embeddings (x); if wanting to prepend them before the word embeddings use the option first = TRUE. If not wanting to train with word embeddings, set x = NULL.
cv_method	Cross-validation method to use within a pipeline of nested outer and inner loops of folds (see nested_cv in rsample). Default is using cv_folds in the outside folds and "validation_split" using rsample::validation_split in the inner loop to achieve a development and assessment set (note that for validation_split the inside_folds should be a proportion, e.g., inside_folds = 3/4); whereas "cv_folds" uses rsample::vfold_cv to achieve n-folds in both the outer and inner loops.
outside_folds	Number of folds for the outer folds (default = 10).
outside_strata_y	Variable to stratify according (default y; can set to NULL).
outside_breaks	The number of bins wanted to stratify a numeric stratification variable in the outer cross-validation loop.
inside_folds	The proportion of data to be used for modeling/analysis; (default proportion = 3/4). For more information see validation_split in rsample.
inside_strata_y	Variable to stratify according (default y; can set to NULL).
inside_breaks	The number of bins wanted to stratify a numeric stratification variable in the inner cross-validation loop.
model	Type of model. Default is "regression"; see also "logistic" for classification.
eval_measure	Type of evaluative measure to select models from. Default = "rmse" for regression and "bal_accuracy" for logistic. For regression use "rsq" or "rmse"; and for classification use "accuracy", "bal_accuracy", "sens", "spec", "precision", "kappa", "f_measure", or "roc_auc", (for more details see the yardstick package).
preprocess_step_center	normalizes dimensions to have a mean of zero; default is set to TRUE. For more info see (step_center in recipes).
preprocess_step_scale	normalize dimensions to have a standard deviation of one. For more info see (step_scale in recipes).
preprocess_PCA	Pre-processing threshold for PCA (to skip this step set it to NA). Can select amount of variance to retain (e.g., .90 or as a grid c(0.80, 0.90)); or number of components to select (e.g., 10). Default is "min_halving", which is a function that selects the number of PCA components based on number of participants and feature (word embedding dimensions) in the data. The formula is: preprocess_PCA = round(max(min(number_features/2), number_participants/2), min(50, number_features))).
penalty	hyper parameter that is tuned

mixture	A number between 0 and 1 (inclusive) that reflects the proportion of L1 regularization (i.e. lasso) in the model (for more information see the linear_reg-function in the parsnip-package). When mixture = 1, it is a pure lasso model while mixture = 0 indicates that ridge regression is being used (specific engines only).
first_n_predictors	by default this setting is turned off (i.e., NA). To use this method, set it to the highest number of predictors you want to test. Then the X first dimensions are used in training, using a sequence from Kjell et al., 2019 paper in Psychological Methods. Adding 1, then multiplying by 1.3 and finally rounding to the nearest integer (e.g., 1, 3, 5, 8). This option is currently only possible for one embedding at the time.
impute_missing	default FALSE (can be set to TRUE if something else than word_embeddings are trained).
method_cor	Type of correlation used in evaluation (default "pearson"; can set to "spearman" or "kendall").
model_description	Text to describe your model (optional; good when sharing the model with others).
multi_cores	If TRUE it enables the use of multiple cores if the computer system allows for it (i.e., only on unix, not windows). Hence it makes the analyses considerably faster to run. Default is "multi_cores_sys_default", where it automatically uses TRUE for Mac and Linux and FALSE for Windows.
save_output	Option not to save all output; default "all". see also "only_results" and "only_results_predictions".
seed	Set different seed.
...	For example settings in yardstick::accuracy to set event_level (e.g., event_level = "second").

Value

A (one-sided) correlation test between predicted and observed values; tibble of predicted values, as well as information about the model (preprocessing_recipe, final_model and model_description).

See Also

see [textEmbedLayerAggregation](#) [textTrainLists](#) [textTrainRandomForest](#) [textSimilarityTest](#)

Examples

```
results <- textTrainRegression(
  x = word_embeddings_4$texts$harmonytext,
  y = Language_based_assessment_data_8$hilstotal,
  multi_cores = FALSE # This is FALSE due to CRAN testing and Windows machines.
)
```

textTranslate	<i>Translation. (experimental)</i>
---------------	------------------------------------

Description

Translation. (experimental)

Usage

```
textTranslate(
  x,
  source_lang = "",
  target_lang = "",
  model = "xlm-roberta-base",
  device = "cpu",
  tokenizer_parallelism = FALSE,
  logging_level = "warning",
  return_incorrect_results = FALSE,
  return_tensors = FALSE,
  return_text = TRUE,
  clean_up_tokenization_spaces = FALSE,
  set_seed = 202208L
)
```

Arguments

x	(string) The text to be translated.
source_lang	(string) The input language. Might be needed for multilingual models (it will not have any effect for single pair translation models). using ISO 639-1 Code, such as: "en", "zh", "es", "fr", "de", "it", "sv", "da", "nn".
target_lang	(string) The desired language output. Might be required for multilingual models (will not have any effect for single pair translation models).
model	(string) Specify a pre-trained language model that have been fine-tuned on a translation task.
device	(string) Name of device to use: 'cpu', 'gpu', or 'gpu:k' where k is a specific device number
tokenizer_parallelism	(boolean) If TRUE this will turn on tokenizer parallelism.
logging_level	(string) Set the logging level. Options (ordered from less logging to more logging): critical, error, warning, info, debug
return_incorrect_results	(boolean) Stop returning some incorrectly formatted/structured results. This setting does CANNOT evaluate the actual results (whether or not they make sense, exist, etc.). All it does is to ensure the returned results are formatted correctly (e.g., does the question-answering dictionary contain the key "answer", is sentiments from textClassify containing the labels "positive" and "negative").

return_tensors (boolean) Whether or not to include the predictions' tensors as token indices in the outputs.
 return_text (boolean) Whether or not to also output the decoded texts.
 clean_up_tokenization_spaces (boolean) Whether or not to clean the output from potential extra spaces.
 set_seed (Integer) Set seed.

Value

A tibble with translated text.

See Also

see [textClassify](#), [textGeneration](#), [textNER](#), [textSum](#), and [textQA](#)

Examples

```
# translation_example <- text::textTranslate(
#   Language_based_assessment_data_8[1,1:2],
#   source_lang = "en",
#   target_lang = "fr",
#   model = "t5-base")
```

textWordPrediction	<i>Compute predictions based on single words for plotting words. The word embeddings of single words are trained to predict the mean value associated with that word. P-values does NOT work yet.</i>
--------------------	---

Description

Compute predictions based on single words for plotting words. The word embeddings of single words are trained to predict the mean value associated with that word. P-values does NOT work yet.

Usage

```
textWordPrediction(
  words,
  word_types_embeddings = word_types_embeddings_df,
  x,
  y = NULL,
  seed = 1003,
  case_insensitive = TRUE,
  text_remove = "[()]",
  ...
)
```

Arguments

words	Word or text variable to be plotted.
word_types_embeddings	Word embeddings from textEmbed for individual words (i.e., decontextualized embeddings).
x	Numeric variable that the words should be plotted according to on the x-axes.
y	Numeric variable that the words should be plotted according to on the y-axes (y=NULL).
seed	Set different seed.
case_insensitive	When TRUE all words are made lower case.
text_remove	Remove special characters
...	Training options from textTrainRegression().

Value

A dataframe with variables (e.g., including trained (out of sample) predictions, frequencies, p-values) for the individual words that is used for the plotting in the textProjectionPlot function.

Examples

```
# Data
# Pre-processing data for plotting
## Not run:
df_for_plotting <- textWordPrediction(
  words = Language_based_assessment_data_8$harmonywords,
  word_types_embeddings = word_embeddings_4$word_types,
  x = Language_based_assessment_data_8$hilstotal
)
df_for_plotting

## End(Not run)
#' @seealso see \link{textProjection}
```

textZeroShot

Zero Shot Classification (Experimental)

Description

Zero Shot Classification (Experimental)

Usage

```

textZeroShot(
    sequences,
    candidate_labels,
    hypothesis_template = "This example is {}. ",
    multi_label = FALSE,
    model = "",
    device = "cpu",
    tokenizer_parallelism = FALSE,
    logging_level = "error",
    return_incorrect_results = FALSE,
    set_seed = 202208L
)

```

Arguments

sequences	(string) The sequence(s) to classify (not that they will be truncated if the model input is too large).
candidate_labels	(string) The set of class labels that is possible in the to classification of each sequence. It may be a single label, a string of comma-separated labels, or a list of labels.
hypothesis_template	(string; optional) The template that is used for turning each of the label into an NLI-style hypothesis. This template must include a "" or similar syntax so that the candidate label can be inserted into the template. For example, the default template is "This example is ." With the candidate label "sports", this would be fed into the model like "<cls> sequence to classify <sep> This example is sports . <sep>". The default template works well in many cases, but it may be worthwhile to experiment with different templates depending on the task setting (see https://huggingface.co/docs/transformers/).
multi_label	(boolean; optional) It indicates whether multiple candidate labels can be true. If FALSE, the scores are normalized such that the sum of the label likelihoods for each sequence is 1. If TRUE, the labels are considered independent and probabilities are normalized for each candidate by doing a softmax of the entailment score vs. the contradiction score.
model	(string) Specify a pre-trained language model that have been fine-tuned on a translation task.
device	(string) Name of device to use: 'cpu', 'gpu', or 'gpu:k' where k is a specific device number
tokenizer_parallelism	(boolean) If TRUE this will turn on tokenizer parallelism.
logging_level	(string) Set the logging level. Options (ordered from less logging to more logging): critical, error, warning, info, debug
return_incorrect_results	(boolean) Stop returning some incorrectly formatted/structured results. This setting does CANOT evaluate the actual results (whether or not they make sense,

exist, etc.). All it does is to ensure the returned results are formatted correctly (e.g., does the question-answering dictionary contain the key "answer", is sentiments from textClassify containing the labels "positive" and "negative").

`set_seed` (Integer) Set seed.

Value

A tibble with the result with the following keys: `sequence` (string) The imputed sequence. `labels` (string) The labels sorted in the order of likelihood. `scores` (numeric) The probabilities for each of the labels.

See Also

see [textClassify](#), [textGeneration](#), [textNER](#), [textSum](#), [textQA](#), [textTranslate](#)

Examples

```
# ZeroShot_example <- text::textZeroShot(sequences = c("I play football",
# "The forest is wonderful"),
# candidate_labels = c("sport", "nature", "research"),
# model = "facebook/bart-large-mnli")
```

word_embeddings_4 *Word embeddings for 4 text variables for 40 participants*

Description

The dataset is a shortened version of the data sets of Study 3-5 from Kjell, Kjell, Garcia and Sikström 2018.

Usage

```
word_embeddings_4
```

Format

A list with word embeddings for harmony words, satisfaction words, harmony text, satisfaction text and decontextualized word embeddings. BERT-base embeddings based on mean aggregation of layer 11 and 12.

words words

n word frequency

Dim1:Dim768 Word embeddings dimensions

Source

<https://psyarxiv.com/er6t7/>

Index

* datasets

- centrality_data_harmony, 3
 - DP_projections_HILS_SWLS_100, 4
 - Language_based_assessment_data_3_100, 5
 - Language_based_assessment_data_8, 5
 - PC_projections_satisfactionwords_40, 6
 - raw_embeddings_1, 7
 - word_embeddings_4, 71
- centrality_data_harmony, 3
- DP_projections_HILS_SWLS_100, 4
- Language_based_assessment_data_3_100, 5
- Language_based_assessment_data_8, 5
- PC_projections_satisfactionwords_40, 6
- raw_embeddings_1, 7
- textCentrality, 7, 11
- textCentralityPlot, 8, 8
- textClassify, 11, 26, 29, 49, 57, 68, 71
- textDescriptives, 13
- textDimName, 14, 19, 20
- textDistance, 15, 17, 52
- textDistanceMatrix, 16
- textDistanceNorm, 16, 17
- textEmbed, 14, 18, 21, 23, 24, 58
- textEmbedLayerAggregation, 20, 20, 23, 66
- textEmbedRawLayers, 20, 21, 22
- textEmbedStatic, 23
- textGeneration, 12, 24, 29, 49, 57, 68, 71
- textModelLayers, 26
- textModels, 26, 27, 28
- textModelsRemove, 27, 27
- textNER, 12, 26, 28, 29, 49, 57, 68, 71
- textPCA, 29, 32
- textPCAPlot, 30, 30
- textPlot, 33
- textPredict, 37, 39, 40
- textPredictAll, 38
- textPredictTest, 39
- textProjection, 8, 11, 37, 40, 47
- textProjectionPlot, 42
- textQA, 12, 26, 29, 47, 49, 57, 68, 71
- textrpp_initialize, 49
- textrpp_install, 50
- textrpp_install_virtualenv (textrpp_install), 50
- textrpp_uninstall, 51
- textSimilarity, 15, 52, 54
- textSimilarityMatrix, 53
- textSimilarityNorm, 15, 52, 53, 54
- textSimilarityTest, 15–17, 38, 52–54, 55, 59, 63, 66
- textSum, 12, 26, 29, 49, 56, 57, 68, 71
- textTokenize, 58
- textTrain, 38–40, 59, 61
- textTrainLists, 38, 59, 60, 63, 66
- textTrainRandomForest, 38, 59, 61, 61, 66
- textTrainRegression, 59, 61, 64
- textTranslate, 12, 26, 29, 49, 57, 67, 71
- textWordPrediction, 68
- textZeroShot, 69
- word_embeddings_4, 71