

Package ‘rcttext’

August 21, 2024

Title Tools for impact analysis in randomized trials with text outcomes

Version 0.0.0.9000

Description A flexible and user-friendly toolkit for performing impact analysis in randomized trials with outcomes generated through human, machine, and/or hybrid scoring of text data. Provides functionality for feature extraction and aggregation, applying supervised and unsupervised machine learning models for semi-automated text scoring, estimating model-assisted treatment impacts with respect to text outcomes under various randomized designs, visually representing found impacts on text outcomes, and additional functionality for performing text analysis using existing frameworks, especially quanteda.

Encoding UTF-8

Imports quanteda,
textreg,
sampling,
tm,
dplyr,
caret,
quanteda.textstats,
plotrix,
randomizr,
stringi

RdMacros Rdpack

RoxygenNote 7.3.1

Depends R (>= 2.10)

LazyData true

Suggests knitr,
rmarkdown,
testthat (>= 3.0.0)

Config/testthat/edition 3

VignetteBuilder knitr

R topics documented:

apply_hunspell	2
clean_features	4
clean_text	5
estimate_impacts	6

extract_liwc	7
extract_taaco	8
extract_w2v	9
generate_features	10
mini_glove	12
ML_iteration	12
ML_plot	14
plot_ccs	17
plot_textfx	17
predict_scores	18
prep_external	18
repair_spelling	19
results.tab	20
run_ccs	20
textfx	21
textfx_terms	22
textML	23
textsamp	24
toy_reads	26
train_ensemble	26

Index	28
--------------	-----------

apply_hunspell	<i>Force spell-correct</i>
----------------	----------------------------

Description

Try to spell correct unrecognized words as best as able using the hunspell package. This will take the first suggestion of hunspell and just replace the text with the suggestion, word by word.

Usage

```
apply_hunspell(
  text,
  additional_words = NULL,
  skip_prefix = NULL,
  threshold = 1,
  to_lower = FALSE,
  verbose = FALSE
)
```

Arguments

text	The text to spell check (list of character vectors)
additional_words	List of words that should be considered as spelled correctly.
skip_prefix	List of prefixes that if a word has we should skip over
threshold	Don't try to correct anything that only occurs threshold or less times in corpus.
to_lower	If TRUE will keep everything lowercase, regardless of spelling suggestions.

Details

All words of 2 or 1 character are skipped.

Value

Vector of text, spell-corrected we hope.

Examples

```
## Example 1: Single string correction

# Text to be checked for spelling errors
txt = "This function seplaces the textt withh the suggestion"

# Check for spelling errors and apply replacements
txt_rep = apply_hunspell( txt, verbose=TRUE, threshold = 0 )
txt_rep

## Example 2: Multiple texts correction

# Load example dataframe with multiple texts
data("toy_reads")

# Extract texts from the dataframe
txts = toy_reads$text
length(txts)

# Check for spelling errors and apply replacements
txts_rep = apply_hunspell( txts, verbose=TRUE, threshold = 0 )
which( txts_rep != txts )

# Display original and corrected texts
txts[7]
txts_rep[7]

txts[20]
txts_rep[20]

# Check for spelling errors and apply replacements with additional accepted words
txts_rep2 = apply_hunspell( txts,
                           verbose=TRUE,
                           additional_words = c( "dont", "rainforest" ),
                           threshold = 0 )
                           which( txts_rep2 != txts )

# Display original and corrected texts with additional accepted words
txts[7]
txts_rep2[7]

txts[20]
txts_rep2[20]
```

clean_features	<i>Clean and simplify feature set</i>
----------------	---------------------------------------

Description

Given set of generated features, simplify set of features using carat package.

Usage

```
clean_features(
  meta,
  ignore = NULL,
  remove.lc = TRUE,
  uniqueCut = 1,
  freqCut = 99,
  cor = 0.95,
  verbose = FALSE
)
```

Arguments

meta	A matrix or dataframe containing the set of only numeric features to be simplified.
ignore	List of column names to ignore when simplifying (e.g., ID column and other columns that should be preserved).
remove.lc	TRUE means remove collinear combinations of features.
uniqueCut	Param for carat's nearZeroVar
freqCut	Param for carat's nearZeroVar
cor	Cutoff of how correlated features should be before dropping one.
verbose	Print out progress to console.

Value

Updated 'meta' in matrix form with fewer columns of the preserved features.

Examples

```
## Example 1
# Create a 10x20 matrix with random numbers between 0 and 1
random_matrix <- matrix(runif(200), nrow = 10, ncol = 20)

# Create a 10x20 dataframe with random numbers between 0 and 1
random_df <- data.frame(matrix(runif(200), nrow = 10, ncol = 20))

# Apply the clean_features function to the matrix
cleaned_matrix <- clean_features(random_matrix)

# Apply the clean_features function to the dataframe and convert back to dataframe
cleaned_df <- data.frame(clean_features(random_df))

## Example 2
```

```
# Apply the clean_features function with a correlation cutoff (0.10)
cleaned_matrix <- clean_features(random_matrix, cor = 0.10)

# Apply the clean_features function to the dataframe with a correlation cutoff (0.90)
cleaned_df <- data.frame(clean_features(random_df, cor = 0.50))
```

clean_text

*Generic text cleaning function***Description**

Quick and easy text cleaning. Take given corpus and remove punctuation, remove whitespace, convert everything to lowercase. This function is used for pre-processing text within the generate_features function. The main use for clean_text is to check which elements of your text are converted to empty strings. Empty strings may result in clean_features dropping desired columns.

Usage

```
clean_text(x, split_hyphens = TRUE)
```

Arguments

x	Character vector or corpus object.
split_hyphens	A logical indicating whether hyphenated words should be treated as two tokens (split at the hyphen).

Examples

```
## Typo correction: given corpus x -> given corpus

# Texts to be cleaned
txts = c( "THIS FUNCTION CONVERTS EVERYTHING TO LOWERCASE.",
          "This function removes punctuation.....",
          "This function removes      whitespace.",
          "This-function-splits-hyphens",
          "The main use for clean_text is to check which elements of
          your text are converted to empty strings.",
          " " )

## Example 1: Convert to Lowercase
txts[1]
clean_text( txts[1] )

## Example 2: Remove Punctuation
txts[2]
clean_text( txts[2] )

## Example 3: Remove White spaces
txts[3]
clean_text( txts[3] )

## Example 4: Split Hyphens
txts[4]
clean_text( txts[4] )
```

```
clean_text( txts[4], split_hyphens = FALSE )

## Example 5: Clean Entire Vector of Texts)
clean_text( txts )
```

estimate_impacts

*Estimate treatment impacts for hybrid-scored text outcomes***Description**

Given text from a randomized trial with a binary treatment, where a subset of the documents have been human-scored, this function computes model-assisted estimates for the average treatment effect with respect to the human-coded outcome.

Usage

```
estimate_impacts(
  y.obs,
  yhat,
  Z,
  wts = NULL,
  design = c("crd", "multi", "cluster", "rcbd"),
  siteID = NULL,
  clusterID = NULL,
  data,
  adjust = NULL
)
```

Arguments

y.obs	A vector of human-coded scores (with NAs for unscored documents).
yhat	A vector of predicted scores estimated via <code>predict_scores</code> .
Z	Indicator for treatment assignment.
wts	Sampling weights for which documents were human scored. Assumed uniform if null.
design	Type of design used for random assignment (complete randomization, multisite randomized, cluster randomized, and blocked and cluster randomized).
siteID	Vector of IDs for site, for multi-site randomized experiments.
clusterID	Vector of IDs for cluster, for cluster-randomized experiments.
data	A <code>data.frame</code> of subject-level identifiers, demographic variables, group membership, and/or other pre-treatment covariates.
adjust	(optional) character vector or named list of variables in the data matrix to adjust for when estimating treatment impacts.

Value

A model object for estimating treatment impact across an array of features.

extract_liwc	<i>Functions for processing and appending output from Linguistic Inquiry Word Count (LIWC) software</i>
--------------	---

Description

Functions for processing and appending output from Linguistic Inquiry Word Count (LIWC) software

Usage

```
extract_liwc(file, meta = NULL, ID.liwc = 1, ID.meta = NULL, clean = TRUE)
```

Arguments

file	character path to LIWC output file
meta	optional data frame to attach LIWC output to
ID.liwc	ID column (either name or column number) of document IDs in the liwc file.
ID.meta	If meta is specified, character vector of the document ID column to use for merging, corresponding with IDs in ID.liwc.
clean	should LIWC output be cleaned prior to appending (e.g., remove linear combinations, repeat variables, etc.)

Examples

```
# This function extracts features from a LIWC-generated-generated data
# (CSV or dataframe).
# The txt_data should be a dataframe
# and have the same column name(s) as "reads_liwc".

data("reads_liwc")
data("example_meta")

reads_liwc = reads_liwc
txt_data = meta

all.feats = extract_liwc( file = "reads_liwc",
                          meta = txt_data,
                          ID.liwc = c( "ID" ),
                          ID.meta = c( "ID" ),
                          clean = FALSE )
```

extract_taaco

Import and merge text features generated using TAACO

Description

This method helps support feature extraction using TAACO. You will have to use the external TACCO program to generate these features; these methods just help move back and forth from R to TACCO.

Usage

```
extract_taaco(
  file,
  meta = NULL,
  ID.meta = NULL,
  drop_para = FALSE,
  drop_sent = TRUE
)
```

Arguments

file	Filename where TAACO results are stored
meta	Optional data.frame with additional document-level variables to include in output.
ID.meta	If meta is specified, character vector with name of variables used for merging.
drop_para	Drop paragraph level measures of cohesion from features (TRUE/FALSE).
drop_sent	Drop adjacent overlap between sentences (TRUE/FALSE).

Details

See prep_external() for generating text files that would be ready for TAACO analysis. Call this on the corpus to make files that can be read in and processed easily.

Once external processing is complete, extract_taaco() reads output and log files produced by the TAACO program and returns a data.frame that can be merged with other feature sets.

The "Filename" column in the read file should be the IDs (with a '.txt' suffix that will be dropped). The results can then be merged with 'meta', if passed,

Value

Returns a data.frame of text features.

Examples

```
# This function extracts features from a TAACO-generated data
# (CSV or dataframe).
# The txt_data should be a dataframe
# and have the same column name(s) as "reads_tacco"

data("reads_taaco")
data("example_meta")
```



```
reads_tacco = reads_taaco
txt_data = meta

all.feats = extract_taaco( "reads_taaco",
                           meta = txt_data,
                           ID.meta = "ID" )
```

extract_w2v	<i>Compute document-level feature vectors from a pre-trained embedding model.</i>
-------------	---

Description

This function generates a vector embedding for each word in a string using a set of pre-trained word vectors such as GloVe (Pennington et al. 2014) and returns the mean vector projection across all words in a document.

Usage

```
extract_w2v(x, meta = NULL, model = NULL)
```

Arguments

x	A corpus object or character vector of text documents.
meta	Dataframe corresponding to the corpus x. If passed, and non-NULL, all features will be added to this dataframe.
model	User-specified model object pointing to a custom pre-trained embedding model, represented as a matrix or data frame where the first column is the word/token and the following columns are numeric vectors. If NULL, use default "mini_glove" embeddings on 1000 common words (not recommended).

Value

A list of data frames containing the Word2Vec projections of the corpus

References

Mikolov T, Chen K, Corrado G, Dean J (2013). "Efficient estimation of word representations in vector space." *arXiv preprint arXiv:1301.3781*. Pennington J, Socher R, Manning C (2014). "Glove: Global vectors for word representation." In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 1532–1543.

Examples

```
# The txt_data should be a dataframe

library(textdata)
data("example_meta")

txt_data = meta
```

```
glove.50d = embedding_glove6b(dimensions = 50)

all.feats = extract_w2v( clean_text( text$text ),
                        meta = txt_data,
                        model = glove.50d )
```

generate_features	<i>Generate an array of text features</i>
-------------------	---

Description

Generates a rich feature representation for documents provided as a character vector or [quanteda::corpus\(\)](#) object by applying an array of linguistic and syntactic indices, available text analysis dictionaries, and pre-trained embedding models to all documents.

Usage

```
generate_features(
  x,
  meta = NULL,
  lex = TRUE,
  sent = TRUE,
  ld = "all",
  clean_features = TRUE,
  read = c("ARI", "Coleman", "DRP", "ELF", "Flesch", "Flesch.Kincaid",
    "meanWordSyllables"),
  terms = NULL,
  preProc = list(uniqueCut = 1, freqCut = 99, cor = 0.95, remove.lc = TRUE),
  verbose = FALSE,
  ignore = NULL
)
```

Arguments

x	A corpus object or character vector of text documents.
meta	Dataframe corresponding to the corpus x. If passed, and non-NULL, all features will be added to this dataframe.
lex	Logical, indicating whether to compute lexical indices including measures of lexical diversity, readability, and entropy
sent	Logical, indicating whether to compute sentiment analysis features from available dictionaries
ld	character vector defining lexical diversity measures to compute; see quanteda.textstats::textstat_lexdiv
clean_features	TRUE means implement cleaning step where we drop features with no variation and collinear features. (This happens before any term generation features are added.)
read	character vector defining readability measures to compute; see quanteda.textstats::textstat_readability
terms	character vector of terms to evaluate as standalone features based on document-level frequency (case-insensitive). Not cleaned by clean_features.

preProc	Named list of arguments passed to <code>caret::preProcess()</code> for applying pre-processing transformations across the set of text features (e.g., removing collinear features).
ignore	List of column names to ignore when simplifying (e.g., ID column and other columns that should be preserved).
...	(optional) additional arguments passed to <code>quanteda::tokens()</code> for text pre-processing.

Value

A data.frame of available text features, one row per document, one column per feature.

Examples

```
# Note: This function does not work with the dataframe with 1 object.
# The 'meta' argument requires a dataframe with at least two rows

## Example 1: Basic Feature Generation

# Create a small dataframe with 2 texts (objects).
df = data.frame(
  text = c( "This function generates an array of text features",
            "This function generates a rich feature representation as a character
            vector or quanteda::corpus() object by applying an array of linguistic
            and syntactic indices" ))

# Generated text features without simplifying the set of features
feats1 = generate_features( df$text,
                           meta = df,
                           clean_features = FALSE )

## Example 2: Feature Generation with Example Data and Customization

# Load example dataframe with multiple texts
data( "toy_reads" )

# Generate text features without simplifying the set of features
feats2 = generate_features( toy_reads$text,
                           meta=toy_reads,
                           clean_features = FALSE,
                           ignore = "ID" )

# Generate preliminary text features,
# simplifying the set of features and specifying sent, read, ld
feats3 = generate_features( toy_reads$text, meta=toy_reads,
                           sent = TRUE,
                           clean_features = TRUE,
                           read = c("Flesch", "Flesch.Kincaid", "ARI"),
                           ld=c("TTR", "R", "K"),
                           ignore=c("ID"),
                           verbose = TRUE )
```

mini_glove	<i>Mini glove dataset</i>
------------	---------------------------

Description

glove embeddings (50 dimensional) for 1000 common words beyond those words listed in several stopword lists provided by quanteda.

Usage

```
mini_glove
```

Format

A [matrix](#) with 1000 rows and 50 columns

Details

Original glove embeddings downloaded from Stanford CITE.

1000 by 50 matrix, each row is an embedding. Rownames are words.

ML_iteration	<i>Iterative Machine Learning with Performance Tracking</i>
--------------	---

Description

Trains a machine learning model multiple times, splitting the data into training and testing sets for each iteration. Evaluates and stores performance metrics for each run, allowing for assessment of model stability and generalization.

Usage

```
ML_iteration(
  x,
  y,
  n_iteration = 1,
  training_portion = 0.8,
  trCon = NULL,
  preProc = NULL,
  n.tune = NULL,
  model = NULL,
  grid = NULL,
  outcome = NULL
)
```

Arguments

<code>x</code>	Matrix or data frame of predictor variables.
<code>y</code>	Vector of outcome variable (continuous or categorical).
<code>n_iteration</code>	Number of iterations (defaults to 1).
<code>training_portion</code>	Proportion of data used for training (defaults to 0.8).
<code>trCon</code>	Train control object for caret (optional).
<code>preProc</code>	Preprocessing methods for caret (optional).
<code>n.tune</code>	Number of tuning parameter combinations (optional).
<code>model</code>	Machine learning model to use (e.g., "rf", "xgbTree").
<code>grid</code>	Tuning grid for hyperparameter search (optional).
<code>outcome</code>	Type of outcome variable: "continuous" or "categorical".

Details

This function provides a standardized way to train and evaluate models within a larger analysis, facilitating comparison across different models, hyperparameters, or repeated runs. It handles both regression (continuous outcome) and classification (categorical outcome) problems, providing appropriate evaluation metrics for each.

Value

A list containing results for each iteration, including: * The trained model ('ML_model') * Predicted values on the test set ('predicted_values') * Values of predictor variables on the test set ('test_values') * Performance metrics (RMSE, MAE, correlation, R-squared for continuous; confusion matrix, accuracy, kappa values for categorical) * Execution time for each iteration ('execution_time')

See Also

* 'ML_iterations': A wrapper function for performing multiple iterations of 'ML_iteration' with different training proportions. * 'caret::train': The underlying function used for model training.

Examples

```
# Load texts
data("toy_reads")

# Generate text features
feats = generate_features( toy_reads$text, meta=toy_reads,
                           sent = TRUE,
                           clean_features = TRUE,
                           read = c("Flesch", "Flesch.Kincaid", "ARI"),
                           ld=c("TTR", "R", "K"),
                           ignore=c("ID"),
                           verbose = TRUE )

# Preprocess the feature space to remove collinear features
# and features with near-zero variance
X_all = dplyr::select( feats,
                       -ID, -Q1, -Q2, -text, -more )
X_all = predict(caret::preProcess( X_all, method = c("nzv", "corr"),
```

```

                                uniqueCut=2, cutoff=0.95), X_all )
caret::findLinearCombos(X_all) # sanity check to make sure no redundant features

# Transform all variables as numeric variables
X_all[] <- lapply(X_all,
                  function(x) if(is.character(x)) as.numeric(as.factor(x)) else x)

# Extract outcome variables
all_Scores <- toy_reads$Q1

## Set parameters
X <- X_all
Y <- all_Scores
n_iter <- 2
n_tune <- 2
control <- caret::trainControl(method = 'cv')
preProc <- 'zv'
outcome <- 'continuous'
portion <- 0.7
best_mod <- 'rf'

## Run ML_iteration
random_forest_Score = ML_iteration (x = X, y = Y,
                                    n_iteration = n_iter,
                                    training_portion = portion,
                                    trCon = control,
                                    preProc = preProc,
                                    n.tune = n_tune,
                                    model = best_mod,
                                    outcome = outcome)

```

ML_plot

Create a Line Plot with Points using ggplot2

Description

This function creates a line plot with points using the ggplot2 package. It allows for extensive customization of aesthetics and layout.

Usage

```

ML_plot(
  data = NULL,
  x = NULL,
  y = NULL,
  color = NULL,
  group = NULL,
  ylim_min = NULL,
  ylim_max = NULL,
  x_lable = NULL,
  y_lable = NULL,
  fig_title = NULL,

```

```

    legend_title = "",
    legend.position = "bottom"
)

```

Arguments

<code>data</code>	A dataframe containing the data to be plotted.
<code>x</code>	A string representing the column name for the x-axis.
<code>y</code>	A string representing the column name for the y-axis.
<code>color</code>	A string representing the column name to be used for coloring the lines and points.
<code>group</code>	A string representing the column name to be used for grouping the lines.
<code>ylim_min</code>	A numeric value representing the minimum limit for the y-axis.
<code>ylim_max</code>	A numeric value representing the maximum limit for the y-axis.
<code>x_lable</code>	A string for the label of the x-axis.
<code>y_lable</code>	A string for the label of the y-axis.
<code>fig_title</code>	A string for the title of the plot.
<code>legend_title</code>	A string for the title of the legend.
<code>legend.position</code>	A string representing the position of the legend in the plot (default is "bottom").

Value

A ggplot object representing the created plot.

Examples

```

# Load texts
data("toy_reads")

# Generate text features
feats = generate_features( toy_reads$text, meta=toy_reads,
                           sent = TRUE,
                           clean_features = TRUE,
                           read = c("Flesch", "Flesch.Kincaid", "ARI"),
                           ld=c("TTR", "R", "K"),
                           ignore=c("ID"),
                           verbose = TRUE )

# Preprocess the feature space to remove collinear features
# and features with near-zero variance
X_all = dplyr::select( feats,
                       -ID, -Q1, -Q2, -text, -more )
X_all = predict(caret::preProcess( X_all, method = c("nzv", "corr"),
                                   uniqueCut=2, cutoff=0.95), X_all )
caret::findLinearCombos(X_all) # sanity check to make sure no redundant features

# Transform all variables as numeric variables
X_all[] <- lapply(X_all,
                  function(x) if(is.character(x)) as.numeric(as.factor(x)) else x)

# Extract outcome variables

```

```

all_Scores <- toy_reads$Q1

## Set parameters
X <- X_all
Y <- all_Scores
n_iter <- 2
n_tune <- 2
control <- caret::trainControl(method = 'cv')
preProc <- 'zv'
outcome <- 'continuous'

## Define the percentages for training portions
percentages <- c(.20, 0.40, 0.60, 0.80)

## Loop through each percentage
best_mod <- 'rf'
random_forest_Scores = ML_iterations ( x = X, y = Y,
                                       n_iteration = n_iter,
                                       training_portions = percentages,
                                       trCon = control,
                                       preProc = preProc,
                                       n.tune = n_tune,
                                       model = best_mod,
                                       outcome = outcome )

best_mod <- 'rrf'
regularized_rf_Scores = ML_iterations ( x = X, y = Y,
                                       n_iteration = n_iter,
                                       training_portions = percentages,
                                       trCon = control,
                                       preProc = preProc,
                                       n.tune = n_tune,
                                       model = best_mod,
                                       outcome = outcome )

## Evaluate the metrics
eval_random_forest_Scores = eval_metrics( random_forest_Scores,
                                           outcome = outcome )

eval_regularized_rf_Scores = eval_metrics( regularized_rf_Scores,
                                           outcome = outcome )

## Create data frames for data visualization
all_Scores = rbind( eval_random_forest_Scores,
                    eval_regularized_rf_Scores )

## Create a line plot (Continuous outcomes)
ML_plot( data = all_Scores,
         x = "training", y = "rmse", color = "model", group = "model",
         ylim_min= 0.5, ylim_max= 0.7,
         x_label = "Training Set (%)", y_label = "RMSE",
         fig_title = "RMSE by Training set % for Different Models",
         legend_title = "", legend.position = "bottom" )

```

plot_ccs	<i>Plot the results from a CCS run</i>
----------	--

Description

This function provides a visualization of the set of words and phrases found to differ systematically between treatment and control groups

Usage

```
plot_ccs(out, xlim = NULL, xadj = c(-0.025, 0.025), ...)
```

Arguments

out	a <code>textreg.result()</code> object
xlim	limits for x-axis
xadj	adjustments to the lower and upper limits on the x-axis of the plot
...	additional arguments passed to plot

plot_textfx	<i>Plot the results from an impact analysis with text outcomes</i>
-------------	--

Description

This function provides a visualization of the set of textual features found to differ systematically between treatment and control groups.

Usage

```
plot_textfx(out, alpha = 0.05, cols = F, group = T, xlim = NULL, ...)
```

Arguments

out	a model object output from <code>estimate_impacts()</code>
alpha	the threshold for determining statistical significance
cols	should effects be colored by direction (red for negative impacts, blue for positive impacts)
group	(optional) should effects be grouped by category (e.g., higher-level summary measures, linguistic features, etc.)
xlim	(optional) bounds for x-axis
...	additional arguments passed to plot

predict_scores	<i>Extract predictions from a fitted text scoring model.</i>
----------------	--

Description

This function computes the predicted scores for a collection of documents based on the results of a trained ensemble learner.

Usage

```
predict_scores(fit, X, Z)
```

Arguments

fit	a model or list of models to use for prediction
newdata	an optional data frame or matrix of predictors
na.action	the method for handling missing data
...	additional arguments to pass to predict.train

Value

A vector of predictions

prep_external	<i>Prepare text documents for analysis using external programs</i>
---------------	--

Description

Text pre-processing and corpus management functions to provide compatibility with external text analysis programs and standalone software packages such as Linguistic Inquiry Word Count (LIWC), the Tool for Automated Analysis of Cohesion (TAACO) and the Sentiment Analysis and Social Cognition Engine (SEANCE).

Usage

```
prep_external(x, dir, docnames = NULL, preProc = NULL)
```

Arguments

x	A corpus object or character vector of text documents.
dir	Name of directory where the generated intermediate text files should be stored.
docnames	Optional character string specifying file names for each document in x.
preProc	Optional text pre-processing function(s) (e.g., stemming) to apply prior to writing text files for analysis in external programs.

References

Pennebaker JW, Booth RJ, Boyd RL, Francis ME (2015). "Linguistic Inquiry and Word Count: LIWC 2015." www.liwc.net. Crossley SA, Kyle K, McNamara DS (2016). "The tool for the automatic analysis of text cohesion (TAACO): Automatic assessment of local, global, and text cohesion." *Behavior research methods*, **48**(4), 1227–1237. Crossley SA, Kyle K, McNamara DS (2017). "Sentiment Analysis and Social Cognition Engine (SEANCE): An automatic tool for sentiment, social cognition, and social-order analysis." *Behavior research methods*, **49**(3), 803–821.

repair_spelling	<i>Replace all words in dictionary with alternates</i>
-----------------	--

Description

Given text as a list of character strings, and a dictionary as a two-column dataframe with the first column being misspelled words and the second being correct spelling, swap all misspelled words with the correct spellings.

Usage

```
repair_spelling(text, dictionary, to_words = NULL)
```

Arguments

text	Character vector
dictionary	Dataframe with two columns of text, or a list of words.
to_words	If dictionary is a list of words, this is list of corresponding words.
Character	vector, revised version of text.

Examples

```
## Example 1: Basic Word Replacement
# Texts to be repaired
txt = "This function replaces alll wordss in dictionary with alternates"

# Repair spelling errors by replacing specified words with their correct forms
txt_rep = repair_spelling( txt,
                           c( "alll", "wordss" ),
                           c( "all", "word" ))

txt_rep

## Example 2: Spelling Correction on a Data Frame Column
data( "toy_reads" )

# Repair spelling errors in the text column of the dataframe

toy_reads$text_rep = repair_spelling( toy_reads$text,
                                       c( "No", "the", "what", "My" ),
                                       c( "NO.", "THE", "WHAT", "MY" ) )

# View the original and repaired texts
view(data.frame(toy_reads$text, toy_reads$text_rep))
```

results.tab	<i>Make results table for grid CCS run</i>
-------------	--

Description

Make results table for grid CCS run

Usage

```
results.tab(result, corp, Z)
```

Arguments

result	a textreg.result() object
corp	a corpus or character vector to calculate term frequencies across
Z	an indicator for treatment assignment
clusterID	optional vector of cluster ID's
...	additional arguments passed to textreg() .

Value

a [textreg.result\(\)](#) object.

run_ccs	<i>Perform Concise Comparative Summarization across a grid of tuning parameters</i>
---------	---

Description

Wrapper for [textreg::textreg\(\)](#).

Determine the penalty C that will zero out the textreg model for a series of randomly permuted labelings with random assignment dictated by a blocked and cluster-randomized experiment.

Usage

```
run_ccs(x, Z, clusterID = NULL)

## S3 method for class 'threshold.C'
cluster(
  x,
  Z,
  design = c("crd", "multi", "cluster", "rcbd"),
  clusterID = NULL,
  siteID = NULL,
  R,
  ...
)
```

Arguments

x	a corpus, character vector of text documents, or set of text features.
Z	an indicator for treatment assignment
clusterID	vector of cluster ID's
design	Type of design used for random assignment (complete randomization, multisite randomized, cluster randomized, and blocked and cluster randomized).
siteID	vector of block ID's
R	Number of times to scramble treatment assignment labels
...	additional arguments passed to <code>textreg()</code> .

Details

Method repeatedly generates +1/-1 vectors within the given blocking structure with blocks of +1/-1 within the clustering vector, and then finds a threshold C for each permutation.

Value

a `textreg.result()` object.

List of numbers. First is the threshold C for the passed labeling. Remainder are the reference distribution based on the permutations.

textfx	<i>Given text from a randomized trial with a binary treatment, this function computes estimates for the average treatment effect with respect to an array of text-based outcomes</i>
--------	--

Description

Given text from a randomized trial with a binary treatment, this function computes estimates for the average treatment effect with respect to an array of text-based outcomes

Usage

```
textfx(
  x,
  Z,
  adj = NULL,
  data,
  wts = NULL,
  design = list(siteID = NULL, clusterID = NULL)
)
```

Arguments

x	A character vector of text documents or a feature matrix
Z	Indicator for treatment assignment.
adj	(optional) character vector or named list of variables in the data matrix to adjust for when estimating treatment impacts.

data	A <code>data.frame</code> of subject-level identifiers, demographic variables, group membership, and/or other pre-treatment covariates.
wt	Sampling weights for documents. Assumed uniform if null.
design	For multi-site and cluster randomized experiments, a named list of vectors containing site IDs and/or cluster IDs.
mcp	character string specifying the correction method to be applied to adjust for multiple comparisons. Defaults to no adjustments. See p.adjust for available adjustment methods.

Value

A model object for estimating treatment impact across an array of features.

textfx_terms	<i>Estimates the average treatment effect on the frequency and prevalence of a list of specified terms and phrases</i>
--------------	--

Description

Estimates the average treatment effect on the frequency and prevalence of a list of specified terms and phrases

Usage

```
textfx_terms(x, Z, terms, ...)
```

Arguments

x	A character vector of text documents or a feature matrix
Z	Indicator for treatment assignment.
terms	Terms and phrases to evaluate
...	optional parameters passed to quanteda::tokens()

Value

A vector showing the frequency and prevalence of the specified terms within each treatment group and results of the hypothesis test comparing prevalence across groups.

textML	<i>Model-assisted impact analysis through hybrid human/machine text scoring</i>
--------	---

Description

A wrapper function for the multiple steps of generating features, training a scoring model on the human-coded data, predicting scores, and comparing human v. machine estimates.

Usage

```
textML(
  x,
  y,
  z = NULL,
  wts = NULL,
  design = c("crd", "multi", "cluster", "rcbd"),
  siteID = NULL,
  clusterID = NULL,
  max.features = NULL,
  ...
)
```

Arguments

x	a corpus or character vector of text documents.
y	a vector of human-coded scores. Set elements to 'NA' for documents not previously scored.
z	optional indicator for treatment assignment. If specified, separate ensembles will be trained for each treatment group;
wts	Sampling weights for which documents were human scored. Assumed uniform if null.
design	Type of design used for random assignment (complete randomization, multisite randomized, cluster randomized, and blocked and cluster randomized).
siteID	Vector of IDs for site, for multi-site randomized experiments.
clusterID	Vector of IDs for cluster, for cluster-randomized experiments.
max.features	maximum number of text features to use for model training. Defaults to 'NULL' (no strict limit)
...	additional arguments passed to train .

Details

This function takes in a corpus of text documents (or a set of computed text features) along with a sample of human-coded outcome values, and trains an ensemble of machine learning models to predict the outcome as a function of the machine measures of text.

Value

a textML model object

textsamp

*Select a random sample of documents***Description**

Functions to select random samples of documents using different sampling schemes and/or along different design criteria.

Proportions will be calculated as $\text{round}(\text{size} * n_k / n)$, where n is the number of documents total.

(In particular, this method will sample $\text{round}(\text{size} / K)$ documents from each cluster, where K is the total number of clusters.) - sounds like equal_size_samples within stratified sampling - n

Usage

```
textsamp(
  x,
  size = NULL,
  prob = NULL,
  wt.fn = NULL,
  scheme = NULL,
  sampling = NULL,
  sampling_control = NULL,
  method = c("srswor", "srswr", "systematic", "poisson"),
  return.data = TRUE,
  ...
)

textsamp_strata(x, size, by = NULL, equal_size_samples = FALSE, ...)

textsamp_cluster(x, by = NULL, n_clusters, cluster_size = NULL, ...)
```

Arguments

<code>x</code>	A corpus object or character vector of text documents or dataframe.
<code>size</code>	The number of documents to sample across all clusters. If <code>NULL</code> , all documents will be sampled from each selected cluster.
<code>prob</code>	a vector of probability weights for each document.
<code>wt.fn</code>	a function for generating probability weights; ignored when <code>prob</code> is used. See Details.
<code>scheme</code>	optional sampling scheme to implement. NOT YET IMPLEMENTED.
<code>sampling</code>	the type of sampling to perform. Options include 'strata' and 'cluster'.
<code>sampling_control</code>	A list containing control parameters for the specified sampling method (see details)
<code>method</code>	the following methods are implemented: simple random sampling without replacement ('srswor'), simple random sampling with replacement ('srswr'), Poisson sampling ('poisson'), systematic sampling ('systematic'); if method is missing, the default method is srswor.

<code>return.data</code>	logical; if TRUE, the function returns the subset of <code>x</code> that are sampled. FALSE returns a vector of row numbers corresponding to the sampled documents.
<code>...</code>	additional arguments passed on to <code>'textsamp'</code> . Cannot include <code>'scheme'</code> .
<code>by</code>	a <code>data.frame</code> with document-level grouping variable(s) or character vector with names of variables in <code>'docvars(x)'</code>
<code>equal_size_samples</code>	TRUE means sample same number of documents from each strata. FALSE means sample proportional to strata size, as described above.
<code>n_clusters</code>	The number of clusters to sample. If not provided, an error will be thrown.
<code>cluster_size</code>	The number of documents to sample from each selected cluster. If NULL, all documents within each cluster will be sampled.
<code>n_cluster</code>	the number of clusters to sample from data

Details

Select a random sample of documents

Functions to select random samples of documents using different sampling schemes and/or along different design criteria.

For stratified sampling (`'sampling = "strata"'`), `'sampling_control'` should be a list with: * `'by'`: A character vector of column names in `'x'` to define the strata. * `'equal_size_samples'`: Logical. If TRUE, sample equal numbers from each stratum.

For cluster sampling (`'sampling = "cluster"'`), `'sampling_control'` should be a list with: * `'by'`: A character vector of column names in `'x'` to define the clusters. * `'n_clusters'`: The number of clusters to sample. * `'cluster_size'`: (Optional) The number of documents to sample from each cluster.

If multiple clustering variables are passed, this will make clusters as all the `_unique_` combinations of these variables. E.g., if A has values of 1, 2, 3 and B has values of 1, 2, then there could be up to six clusters.

Value

Returns either the sampled data or a vector of rownumbers of sampled documents.

Examples

```
# Load example dataframe
data("toy_reads")

## Example 1: Sample 4 documents using the default method
##           (simple random sampling without replacement)
textamp_df = textsamp(toy_reads, size = 4)

## Example 2: Sample 8 documents using Poisson sampling
textamp_df2 = textsamp(toy_reads, size = 8, method = "poisson")

## Example 3: Sample 8 documents using systematic sampling,
##           but only return the row numbers
textamp_df3 = textsamp(toy_reads, size = 8,
                      method = "systematic", return.data = FALSE)

# Load example dataframe
data("toy_reads")
```

```
## Example 1: Stratified Sampling by a Single Variable
stratified_df1 <- textsamp_strata(toy_reads, size = 6, by = "more")

Example 2: Stratified Sampling by Multiple Variables
stratified_df2 <- textsamp_strata(toy_reads, size = 10,
                                by = c("Q1", "more"))

Example 3: Stratified Sampling with Unequal Sample Sizes per Stratum
stratified_df3 <- textsamp_strata(toy_reads, size = 10,
                                by = "more", equal_size_samples = FALSE)
```

toy_reads	<i>Dataset with 20 essays from READS pilot data</i>
-----------	---

Description

Used for testing and illustration of rcttext functions.

Usage

```
toy_reads
```

Format

A [data.frame](#) with 5 columns and 20 rows

Details

5 column data.frame, ID is the id of subject, Q1, Q2, more are meta information on scoring, and text contains character string of the text of the essay.

train_ensemble	<i>Train an ensemble learner for semi-supervised text scoring</i>
----------------	---

Description

This function takes in a corpus of text documents or a set of computed text features, along with a sample of human-coded outcome values and trains an ensemble of machine learning models to predict the outcome as a function of machine measures of text.

Usage

```
train_ensemble(
  X,
  Z,
  Yobs,
  coded,
  n.tune = 3,
  bounds = NULL,
```

```
preProc = "zv",  
model = "rf",  
seeds = NA,  
...  
)
```

Arguments

n.tune	an integer denoting the amount of granularity in the tuning parameter grid. By default, this argument is the number of levels for each tuning parameters that should be generated by train .
bounds	a vector (y1, y2) specifying the lower and upper limits for prediction
...	additional arguments passed to trainControl .
x	a data.frame or matrix of numeric text features.
y	a vector of human-coded scores for the outcome of interest.
z	optional indicator for treatment assignment. If specified, separate ensembles will be trained for each treatment group;
cvf	number of folds for cross validation
return.all	should all component models be returned? If 'FALSE', returns only the fitted ensemble(s).

Value

a fitted model object

Index

*** data**
 mini_glove, [12](#)
 toy_reads, [26](#)

apply_hunspell, [2](#)

clean_features, [4](#)
clean_text, [5](#)
cluster.threshold.C(run_ccs), [20](#)
corpus, [9](#), [10](#), [18](#), [24](#)

data.frame, [26](#)

estimate_impacts, [6](#)
extract_liwc, [7](#)
extract_taaco, [8](#)
extract_w2v, [9](#)

generate_features, [10](#)

matrix, [12](#)
mini_glove, [12](#)
ML_iteration, [12](#)
ML_plot, [14](#)

p.adjust, [22](#)
plot_ccs, [17](#)
plot_textfx, [17](#)
predict_scores, [18](#)
prep_external, [18](#)

quanteda.textstats::textstat_lexdiv(),
 [10](#)
quanteda.textstats::textstat_readability(),
 [10](#)
quanteda::corpus(), [10](#)
quanteda::tokens(), [11](#), [22](#)

repair_spelling, [19](#)
results.tab, [20](#)
run_ccs, [20](#)

textfx, [21](#)
textfx_terms, [22](#)
textML, [23](#)

textreg(), [20](#), [21](#)
textreg.result(), [17](#), [20](#), [21](#)
textreg::textreg(), [20](#)
textsamp, [24](#)
textsamp_cluster(textsamp), [24](#)
textsamp_strata(textsamp), [24](#)
toy_reads, [26](#)
train, [23](#), [27](#)
train_ensemble, [26](#)
trainControl, [27](#)