

Package ‘tada’

May 9, 2022

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,

RdMacros Rdpack

RoxygenNote 7.1.1

R topics documented:

| | |
|----------------------------|----|
| estimate_impacts | 2 |
| extract_taaco | 3 |
| get_dimnames | 3 |
| glove.50d | 4 |
| plot.ccs | 4 |
| plot.textfx | 5 |
| predict_scores | 5 |
| prep_external | 6 |
| results.tab | 6 |
| run_ccs | 7 |
| tada | 8 |
| textfx | 9 |
| textML | 10 |
| textsamp | 11 |
| train_ensemble | 12 |

| | |
|--------------|-----------|
| Index | 13 |
|--------------|-----------|

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 predict_scores. |
| 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 data.frame 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_taaco | <i>Manage and merge text features generated using TAACO</i> |
|---------------|---|

Description

Tools to support feature extraction using TAACO. `prep_taaco()` prepares a corpus for analysis in TAACO. `extract_taaco()` reads output and log files produced by TAACO program and returns a `data.frame` that can be merged with other feature sets.

Usage

```
extract_taaco(file, data = NULL, idvar = NULL)
```

```
prep_taaco(x, dir, docnames = NULL)
```

Arguments

| | |
|-----------------------|---|
| <code>file</code> | Filename where TAACO results are stored |
| <code>data</code> | Optional <code>data.frame</code> with additional document-level variables to include in output. |
| <code>idvar</code> | If <code>data</code> is specified, character vector with name(s) of variables used for merging. |
| <code>x</code> | A <code>[quanteda::corpus()]</code> object or character vector of text documents. |
| <code>dir</code> | Name of directory where TAACO intermediate text files should be stored. |
| <code>docnames</code> | Optional character string specifying file names for each document in <code>x</code> . |

Value

Returns a `data.frame` of text features.

| | |
|--------------|--|
| get_dimnames | <i>Internal functions for processing LIWC output</i> |
|--------------|--|

Description

Internal functions for processing LIWC output

Usage

```
get_dimnames()
```

| | |
|-----------|---|
| glove.50d | <i>Dataset containing word embeddings on 50 dimensions based on GloVe pre-trained embedding model</i> |
|-----------|---|

Description

Dataset containing word embeddings on 50 dimensions based on GloVe pre-trained embedding model

Usage

```
glove.50d
```

Format

A [data.frame](#) with 400,000 terms and 50 dimensions of word embeddings

Details

50-dimensional word embedding vectors for 400,000 terms and phrases based on GloVe pre-trained embedding model

| | |
|----------|--|
| 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

```
## S3 method for class 'ccs'
plot(out, xadj = c(-0.025, 0.025), ...)
```

Arguments

| | |
|------|---|
| out | a textreg.result() object |
| xadj | adjustments to the lower and upper limits on the x-axis of the plot |
| ... | additional arguments passed to plot |

plot.textfx

Plot the results from an impact analysis with text outcomes

Description

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

Usage

```
## S3 method for class 'textfx'
plot(out, alpha = 0.05, cols = F, group = NULL, ...)
```

Arguments

| | |
|-------|---|
| 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.) |
| ... | additional arguments passed to plot |
| x | a model object output from estimate_impacts() |

predict_scores

Extract predictions from a fitted text scoring model.

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, newdata, na.action = na.omit, ...)
```

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 TAACO 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.

| | |
|-------------|--|
| 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 textreg() . |

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.

| | |
|------|---|
| tada | <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

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

Arguments

| | |
|---------|--|
| x | A <code>corpus</code> object or character vector of text documents. |
| 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 <code>quanteda.textstats::textstat_lexdiv</code> |
| read | character vector defining readability measures to compute; see <code>quanteda.textstats::textstat_readability</code> |
| terms | character vector of terms to evaluate as standalone features based on document-level frequency (case-insensitive) |
| 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) |
| ... | (optional) additional arguments passed to <code>quanteda::tokens()</code> for text pre-processing. |

Value

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

References

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.

| | |
|--------|--|
| 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,  
  mcp = "none",  
  wts = NULL,  
  design = list(siteID = NULL, clusterID = NULL)  
)
```

Arguments

| | |
|--------|--|
| x | A character vector of text documents or a feature matrix returned by tada |
| 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. |
| 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. |
| wts | 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. |

Value

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

| | |
|--------|---|
| 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 | <i>Select a random sample of documents</i> |
|----------|--|

Description

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

Usage

```
textsamp(
  x,
  size = length(x),
  prob = NULL,
  wt.fn = NULL,
  scheme = NULL,
  method = c("srswr", "srswor", "systematic", "poisson")
)

textsamp_strata(x, by = NULL, ...)

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

Arguments

| | |
|---------------------|---|
| <code>x</code> | A corpus object or character vector of text documents. |
| <code>size</code> | a non-negative integer giving the number of documents to sample. |
| <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 |
| <code>method</code> | the following methods are implemented: simple random sampling without replacement (<code>'srswor'</code>), simple random sampling with replacement (<code>'srswr'</code>), Poisson sampling (<code>'poisson'</code>), systematic sampling (<code>'systematic'</code>); if <code>method</code> is missing, the default method is <code>srswor</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>...</code> | additional arguments passed on to <code>'textsamp'</code> . Cannot include <code>'scheme'</code> . |

Value

Returns a `data.frame` containing identifiers for the selected documents.

train_ensemble

*Train an ensemble learner for semi-supervised text scoring***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,
  y,
  z = NULL,
  n.tune = 3,
  cvf = 5,
  bounds = NULL,
  ...,
  return.all = TRUE
)
```

Arguments

| | |
|------------|--|
| x | a <code>data.frame</code> 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; |
| 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 . |
| cvf | number of folds for cross validation |
| bounds | a vector (y1, y2) specifying the lower and upper limits for prediction |
| ... | additional arguments passed to trainControl . |
| return.all | should all component models be returned? If 'FALSE', returns only the fitted ensemble(s). |

Value

a fitted model object

Index

- * **data**
 - glove.50d, 4
- cluster.threshold.C(run_ccs), 7
- corpus, 6, 8, 11
- data.frame, 4
- estimate_impacts, 2
- extract_taaco, 3
- get_dimnames, 3
- glove.50d, 4
- p.adjust, 9
- plot.ccs, 4
- plot.textfx, 5
- predict_scores, 5
- prep_external, 6
- prep_taaco(extract_taaco), 3
- quanteda.textstats::textstat_lexdiv(), 8
- quanteda.textstats::textstat_readability(), 8
- quanteda::corpus(), 8
- quanteda::tokens(), 8
- results.tab, 6
- run_ccs, 7
- tada, 8, 9
- textfx, 9
- textML, 10
- textreg(), 6, 7
- textreg.result(), 4, 6, 7
- textreg::textreg(), 7
- textsamp, 11
- textsamp_cluster(textsamp), 11
- textsamp_strata(textsamp), 11
- train, 10, 12
- train_ensemble, 12
- trainControl, 12