

# Package ‘textmatch’

December 6, 2019

**Title** Toolkit for Matching Textual Data and Evaluating Textual Similarity

**Version** 0.0.0.9000

**Description** What the package does (one paragraph).

**Depends** R (>= 3.5.2)

**License** What license is it under?

**Encoding** UTF-8

**LazyData** false

**RoxygenNote** 7.0.2

**Imports** dplyr,  
data.table,  
quantda,  
stats,  
progress,  
tm,  
stm,  
optmatch,  
textir,  
text2vec,  
cem,  
Rfast,  
reshape2,  
Matrix,  
cobalt

**Suggests** knitr,  
rmarkdown,  
softmaxreg,  
glmnet

**VignetteBuilder** knitr

## R topics documented:

corp_contrast . . . . .	2
FoxCNNcorpus . . . . .	2
FoxCNNmeta . . . . .	3
FoxCNNsurvey . . . . .	3
get_bal_optmatch . . . . .	4

get_CEM . . . . .	5
get_matches . . . . .	5
get_similarity_scores . . . . .	6
get_word2vec_glove . . . . .	6
pair_distances . . . . .	7
quality_model . . . . .	7
SE . . . . .	8
select_matches . . . . .	8
textPS . . . . .	9
textPS_dist . . . . .	10
transform_dfm . . . . .	10
transform_stm . . . . .	11
<b>Index</b>	<b>12</b>

---

corp_contrast	<i>Corpus contrasts</i>
---------------	-------------------------

---

**Description**

Corpus contrasts

**Usage**

corp\_contrast(corp, Z)

**Arguments**

- |      |                                       |
|------|---------------------------------------|
| corp | a fitted <a href="#">stm</a> object   |
| Z    | an indicator for treatment assignment |

**Value**

Corpus contrast summaries

---

FoxCNNcorpus	<i>Corpus with 1,565 articles from CNN and 1,796 articles from Fox News.</i>
--------------	--

---

**Description**

Corpus with 1,565 articles from CNN and 1,796 articles from Fox News.

**Usage**

FoxCNNcorpus

**Format**

A [data.frame](#) with 3,361 observations of 5 articles

## Details

Corpus of front-page news articles published online by CNN or Fox News from 12/20/2014 to 05/09/2015 containing 1,565 articles from CNN and 1,796 articles from Fox News. Data contains article identifiers corresponding to data in FoxCNNmeta as well as raw and cleaned text data.

## References

Mozier et al. (2019) [Matching with Text Data: An Experimental Evaluation of Methods for Matching Documents and of Measuring Match Quality](#). *Political Analysis*, Forthcoming.

---

FoxCNNmeta	<i>Dataset containing metadata for FoxCNN corpus with 1,565 articles from CNN and 1,796 articles from Fox News.</i>
------------	---

---

## Description

Dataset containing metadata for FoxCNN corpus with 1,565 articles from CNN and 1,796 articles from Fox News.

## Usage

FoxCNNmeta

## Format

A [data.frame](#) with 3,361 observations of 5 articles

## Details

Metadata for the FoxCNN corpus including article names, original URLs, and dates of publication.

## References

Mozier et al. (2019) [Matching with Text Data: An Experimental Evaluation of Methods for Matching Documents and of Measuring Match Quality](#). *Political Analysis*, Forthcoming.

---

FoxCNNsurvey	<i>Dataset containing distance measurements and average quality scores for sample of 505 pairs of matched Fox and CNN articles evaluated by human coders on Mechanical Turk. Raw quality scores are regressed on distance measurements to fit a predictive model for match quality as a function of the 117 metrics considered.</i>
--------------	---

---

## Description

Dataset containing distance measurements and average quality scores for sample of 505 pairs of matched Fox and CNN articles evaluated by human coders on Mechanical Turk. Raw quality scores are regressed on distance measurements to fit a predictive model for match quality as a function of the 117 metrics considered.

**Usage**

```
FoxCNNsurvey
```

**Format**

A [data.frame](#) with 505 pairs and 108 features.

**Details**

Distance measurements and average quality scores for a sample of 505 matched pairs of documents evaluated by human coders.

**References**

Mozer et al. (2019) \"Matching with Text Data: An Experimental Evaluation of Methods for Matching Documents and of Measuring Match Quality\". *Political Analysis*, Forthcoming.

---

get_bal_optmatch	<i>Generates balance diagnostics for many possible matched datasets which can be used to make comparisons between matching methods.</i>
------------------	---

---

**Description**

Generates balance diagnostics for many possible matched datasets which can be used to make comparisons between matching methods.

**Usage**

```
get_bal_optmatch(matchobj, Z, covs)
```

**Arguments**

matchobj	an <a href="#">optmatch</a> object or a matched dataset
Z	a vector of treatment indicators
covs	a matrix or data frame of covariates to assess balance on

**Value**

A [data.frame](#) MNIR sufficient reduction scores for a corpus

---

get_CEM	<i>Create a data frame of pairs of documents obtained through coarsened exact matching (CEM) within a specified number of bins and return indices for matched sets</i>
---------	--

---

### Description

Create a data frame of pairs of documents obtained through coarsened exact matching (CEM) within a specified number of bins and return indices for matched sets

### Usage

```
get_CEM(x, Z, rep.name, cuts, caliper_fun, verbose = FALSE, SR = NULL)
```

### Arguments

x	a text representation
Z	a vector of treatment indicators
rep.name	a string or character with the name of the matching method
cuts	a function for how the variables will be binned. Defaults to "median"

### Value

A [data.frame](#) of indices for matched pairs of documents

---

get_matches	<i>Similarity and distance computation between documents or features</i>
-------------	--

---

### Description

These functions compute distance matrices from a text representation where each row is a document and each column is a feature to measure distance over based on treatment indicator Z

### Usage

```
get_matches(
  dist,
  Z,
  dist.name,
  caliper_fun,
  within.calip = "quantile",
  tol = 0.001
)
```

**Arguments**

<code>Z</code>	a vector of treatment indicators
<code>dist.name</code>	a string or character with the name of the matching method
<code>caliper_fun</code>	an optional function specifying the caliper to enforce when matching
<code>x</code>	a matrix of pairwise distances for all potential matches of treatment and control units. See <a href="#">pair_distances</a> .

---

`get_similarity_scores` *This function calculates an input character vector's similarity matrix according to the measures contained in the predictive model.*

---

**Description**

This function calculates an input character vector's similarity matrix according to the measures contained in the predictive model.

**Usage**

```
get_similarity_scores(x)
```

**Arguments**

<code>x</code>	A character vector where each element is a document
----------------	---

**Value**

A data frame of rows ( $n * n - 1$ ) and columns 16; each column is one of the constituent similarity measures

---

`get_word2vec_glove` *This function calculates the Word2Vec embeddings*

---

**Description**

This function calculates the Word2Vec embeddings

**Usage**

```
get_word2vec_glove(dir.source, corpus)
```

**Arguments**

<code>dat</code>	FoxCNN corpus to calculate Word2Vec scores for
------------------	--

**Value**

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

---

pair_distances	<i>Similarity and distance computation between documents or features</i>
----------------	--

---

### Description

These functions compute distance matrices from a text representation where each row is a document and each column is a feature to measure distance over based on treatment indicator Z

### Usage

```
pair_distances(
  dat,
  Z,
  include = c("cosine", "euclidean", "mahalanobis"),
  form = "data.frame",
  verbose = FALSE
)
```

### Arguments

include	Which distances to calculate
form	Should the distances be returned as a list of matrices or condensed into a single data frame?
x	a matrix text representation with rows corresponding to each document in a corpus and columns that represent summary measures of the text (e.g., word counts, topic proportions, etc.). Acceptable forms include a valid <b>quanteda</b> dfm object, a <b>tm</b> Document-Term Matrix, a matrix of estimated topic proportions, or a vector of estimated propensity scores.

### Value

A matrix showing pairwise distances for all potential matches of treatment and control units under various distance metrics

---

quality_model	<i>Fitted model for pairwise match quality as a function of 117 distance metrics calculated in Mozer et al. (2019). Trained on "FoxCNNsurvey" dataset.</i>
---------------	--

---

### Description

Fitted model for pairwise match quality as a function of 117 distance metrics calculated in Mozer et al. (2019). Trained on "FoxCNNsurvey" dataset.

### Usage

```
quality_model
```

**Format**

A `glmnet` model object.

**Details**

Fitted model for predicting the match quality score for a given pair of text documents as a function of 117 distance measurements.

**References**

Mozer et al. (2019) \"Matching with Text Data: An Experimental Evaluation of Methods for Matching Documents and of Measuring Match Quality\". *Political Analysis*, Forthcoming.

---

SE	<i>Utility functions for textmatch objects</i>
----	--

---

**Description**

Utility functions for textmatch objects

**Usage**

`SE(v, z)`

**Arguments**

`v` a vector of covariate or outcome data

**Value**

`z` a vector of treatment indicators

---

<code>select_matches</code>	<i>Given a model for match quality and a corpus of documents, calculate the estimated match quality of each potential pairing of treated and control documents and return a matched dataset containing all pairs of documents with estimated quality above a specified threshold</i>
-----------------------------	--

---

**Description**

Given a model for match quality and a corpus of documents, calculate the estimated match quality of each potential pairing of treated and control documents and return a matched dataset containing all pairs of documents with estimated quality above a specified threshold

**Usage**

`select_matches(corpus, Z, mod, threshold)`



**Arguments**

corpus	description
Z	treatment indicator
threshold	for quality scores to return

**Value**

A [data.frame](#) of matched pairs of documents

---

textPS	<i>Calculate a univariate projection and pairwise distances for a corpus of text documents with binary treatment indicator Z.</i>
--------	---

---

**Description**

Calculate a univariate projection and pairwise distances for a corpus of text documents with binary treatment indicator Z.

Given a text representation x and treatment indicator Z, estimate a univariate projection of the text data that can serve as a sufficient statistic for the confounding information captured within the text.

**Usage**

```
textPS(x, Z, verbose = FALSE, g = 1)
```

```
textPS(x, Z, verbose = FALSE, g = 1)
```

**Arguments**

x	a TDM or DFM text representation
Z	a vector of treatment indicators
normalize	should the TDM features be normalized? Defaults to TRUE.
return.df	Should the distances be returned as a data frame? Default is TRUE.

**Value**

A [data.frame](#) MNIR sufficient reduction scores for a corpus

A [data.frame](#) MNIR sufficient reduction scores for a corpus

---

textPS_dist	<i>Given a text representation <math>x</math> and treatment indicator <math>Z</math>, construct a distance matrix for the pairwise distances between treatment and control documents based on the textual propensity score.</i>
-------------	---

---

### Description

Given a text representation  $x$  and treatment indicator  $Z$ , construct a distance matrix for the pairwise distances between treatment and control documents based on the textual propensity score.

### Usage

```
textPS_dist(x, Z, normalize = TRUE, return.df = TRUE, verbose = FALSE)
```

### Arguments

$x$	a TDM or DFM text representation
$Z$	a vector of treatment indicators
normalize	should the TDM features be normalized? Defaults to TRUE.
return.df	Should the distances be returned as a data frame? Default is TRUE.

---

transform_dfm	<i>Applies bounds, weights, and/or coarsening schemes to a dfm or document frequency matrix to reduce the dimension of the data, reduce noise, or apply other design rules (e.g. - to exclude words that occur in too few or too many documents).</i>
---------------	---

---

### Description

Applies bounds, weights, and/or coarsening schemes to a dfm or document frequency matrix to reduce the dimension of the data, reduce noise, or apply other design rules (e.g. - to exclude words that occur in too few or too many documents).

### Usage

```
transform_dfm(x, bounds, tfidf = FALSE, verbose = TRUE)
```

### Arguments

$x$	a matrix text representation with rows corresponding to each document in a corpus and columns that represent summary measures of the text (e.g., word counts, topic proportions, etc.). Acceptable forms include a valid <b>quanteda</b> dfm object, a <b>tm</b> Document-Term Matrix, or a matrix of estimated topic proportions.
bounds	a vector of lower and upper bounds to enforce. Defaults to excluding any terms that appear in only one document and any terms that appear in every document
tfidf	optional scheme to use for weighting the DTM. Defaults to FALSE.
verbose	indicator for verbosity

### Value

A bounded DFM

---

transform_stm	<i>Refits a STM with a content-based covariate so that all document-level topic-proportions are estimated "as-treated". Also allows for calculation of the SR sufficient reduction and optional coarsening to reduce the dimension of the data, reduce noise, or apply other design rules (e.g. - to exclude words that occur in too few or too many documents).</i>
---------------	--

---

### Description

Refits a STM with a content-based covariate so that all document-level topic-proportions are estimated "as-treated". Also allows for calculation of the SR sufficient reduction and optional coarsening to reduce the dimension of the data, reduce noise, or apply other design rules (e.g. - to exclude words that occur in too few or too many documents).

### Usage

```
transform_stm(mod, out, Z, calc.SR = FALSE, coarsen = FALSE, simplex = FALSE)
```

### Arguments

mod	a fitted <a href="#">stm</a> object
out	the original call to the STM
Z	an indicator for treatment assignment
calc.SR	an indicator for returning the sufficient reduction. Default is TRUE.
coarsen	an indicator for returning the coarsened STM

### Value

A bounded DFM

# Index

- \*Topic **datasets**
  - quality\_model, [7](#)
- \*Topic **data**
  - FoxCNNcorpus, [2](#)
  - FoxCNNmeta, [3](#)
  - FoxCNNsurvey, [3](#)
- corp\_contrast, [2](#)
- data.frame, [2–5](#), [9](#)
- FoxCNNcorpus, [2](#)
- FoxCNNmeta, [3](#)
- FoxCNNsurvey, [3](#)
- get\_bal\_optmatch, [4](#)
- get\_CEM, [5](#)
- get\_matches, [5](#)
- get\_similarity\_scores, [6](#)
- get\_word2vec\_glove, [6](#)
- glmnet, [8](#)
- optmatch, [4](#)
- pair\_distances, [6](#), [7](#)
- quality\_model, [7](#)
- SE, [8](#)
- select\_matches, [8](#)
- stm, [2](#), [11](#)
- textPS, [9](#)
- textPS\_dist, [10](#)
- transform\_dfm, [10](#)
- transform\_stm, [11](#)