

# Package ‘RTextTools’

July 2, 2014

**Type** Package

**Title** Automatic Text Classification via Supervised Learning

**Version** 1.4.2

**Date** 2014-01-18

**Author** Timothy P. Jurka, Loren Collingwood, Amber E. Boydston, Emiliano Grossman, Wouter van Atteveldt

**Maintainer** Timothy P. Jurka <tpjurka@ucdavis.edu>

**Depends** R (>= 2.15.0), SparseM

**Imports** methods, randomForest, tree, nnet, tm, e1071, ipred, caTools, maxent, glmnet, tau

**Description** RTextTools is a machine learning package for automatic text classification that makes it simple for novice users to get started with machine learning, while allowing experienced users to easily experiment with different settings and algorithm combinations. The package includes nine algorithms for ensemble classification (svm, slda, boosting, bagging, random forests, glmnet, decision trees, neural networks, maximum entropy), comprehensive analytics, and thorough documentation.

**License** GPL-3

**URL** <http://www.rtexttools.com/>

**NeedsCompilation** yes

**Repository** CRAN

**Date/Publication** 2014-01-19 09:07:18

## R topics documented:

analytics-class . . . . .	2
analytics_virgin-class . . . . .	3
classify_model . . . . .	4
classify_models . . . . .	5
create_analytics . . . . .	6
create_container . . . . .	7
create_ensembleSummary . . . . .	8
create_matrix . . . . .	9
create_precisionRecallSummary . . . . .	10
create_scoreSummary . . . . .	11
cross_validate . . . . .	12
getStemLanguages . . . . .	13
matrix_container-class . . . . .	14
NYTimes . . . . .	15
print_algorithms . . . . .	16
read_data . . . . .	17
recall_accuracy . . . . .	18
summary.analytics . . . . .	19
summary.analytics_virgin . . . . .	19
train_model . . . . .	20
train_models . . . . .	22
USCongress . . . . .	23
wordStem . . . . .	24
<b>Index</b>	<b>26</b>

---

analytics-class	<i>an S4 class containing the analytics for a classified set of documents.</i>
-----------------	--

---

### Description

An S4 class containing the analytics for a classified set of documents. This includes a label summary, document summary, ensemble summary, and algorithm summary. This class is returned if `virgin=FALSE` in `create_container`.

### Objects from the Class

Objects could in principle be created by calls of the form `new("analytics", ...)`. The preferred form is to have them created via a call to `create_analytics`.

**Slots**

`label_summary` Object of class "data.frame": stores the analytics for each label, including the percent coded accurately and how much overcoding occurred

`document_summary` Object of class "data.frame": stores the analytics for each document, including all available raw data associated with the learning process

`algorithm_summary` Object of class "data.frame": stores precision, recall, and F-score statistics for each algorithm, broken down by label

`ensemble_summary` Object of class "matrix": stores the accuracy and coverage for an n-algorithm ensemble scoring

**Author(s)**

Timothy P. Jurka <tpjurka@ucdavis.edu>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100, size=100, replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"], data["Subject"]), language="english",
  removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix, data$Topic.Code, trainSize=1:75, testSize=76:100,
  virgin=FALSE)
models <- train_models(container, algorithms=c("MAXENT", "SVM"))
results <- classify_models(container, models)
analytics <- create_analytics(container, results)

summary(analytics)
```

---

analytics\_virgin-class

*an S4 class containing the analytics for a classified set of documents.*

---

**Description**

An S4 class containing the analytics for a classified set of documents. This includes a label summary and a document summary. This class is returned if `virgin=TRUE` in `create_container`.

**Objects from the Class**

Objects could in principle be created by calls of the form `new("analytics_virgin", ...)`. The preferred form is to have them created via a call to `create_analytics`.

**Slots**

`label_summary` Object of class "data.frame": stores the analytics for each label, including how many documents were classified with each label

`document_summary` Object of class "data.frame": stores the analytics for each document, including all available raw data associated with the learning process

**Author(s)**

Timothy P. Jurka <tpjurka@ucdavis.edu>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100,size=100,replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"],data["Subject"]), language="english",
removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix,data$Topic.Code,trainSize=1:75, testSize=76:100,
virgin=TRUE)
models <- train_models(container, algorithms=c("MAXENT","SVM"))
results <- classify_models(container, models)
analytics <- create_analytics(container, results)

summary(analytics)
```

---

classify_model	<i>makes predictions from a train_model() object.</i>
----------------	---

---

**Description**

Uses a trained model from the [train\\_model](#) function to classify new data.

**Usage**

```
classify_model(container, model, s=0.01, ...)
```

**Arguments**

container	Class of type <a href="#">matrix_container-class</a> generated by the <a href="#">create_container</a> function.
model	Slot for trained SVM, SLDA, boosting, bagging, RandomForests, glmnet, decision tree, neural network, or maximum entropy model generated by <a href="#">train_model</a> .
s	Penalty parameter lambda for <b>glmnet</b> classification.
...	Additional parameters to be passed into the predict function of any algorithm.

**Details**

Only one model may be passed in at a time for classification. See [train\\_models](#) and [classify\\_models](#) to train and classify using multiple algorithms.

**Value**

Returns a `data.frame` of predicted codes and probabilities for the specified algorithm.

**Author(s)**

Loren Collingwood <lorenc2@uw.edu>, Timothy P. Jurka <tpjurka@ucdavis.edu>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100,size=100,replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"],data["Subject"]), language="english",
  removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix,data$Topic.Code,trainSize=1:75, testSize=76:100,
  virgin=FALSE)
maxent_model <- train_model(container,"MAXENT")
maxent_results <- classify_model(container,maxent_model)
```

---

classify\_models      *makes predictions from a train\_models() object.*

---

**Description**

Uses a trained model from the [train\\_models](#) function to classify new data.

**Usage**

```
classify_models(container, models, ...)
```

**Arguments**

container	Class of type <a href="#">matrix_container-class</a> generated by the <a href="#">create_container</a> function.
models	List of models to be used for classification generated by <a href="#">train_models</a> .
...	Other parameters to be passed on to <a href="#">classify_model</a> .

**Details**

Use the list returned by [train\\_models](#) to use multiple models for classification.

**Author(s)**

Wouter Van Atteveldt <wouter@vanatteveldt.com>, Timothy P. Jurka <tpjurka@ucdavis.edu>

**Examples**

```

library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100,size=100,replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"],data["Subject"]), language="english",
removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix,data$Topic.Code,trainSize=1:75, testSize=76:100,
virgin=FALSE)
models <- train_models(container, algorithms=c("MAXENT","SVM"))
results <- classify_models(container, models)

```

---

create_analytics	<i>creates an object of class analytics given classification results.</i>
------------------	---

---

**Description**

Takes the results from functions [classify\\_model](#) or [classify\\_models](#) and computes various statistics to help interpret the data.

**Usage**

```
create_analytics(container, classification_results, b=1)
```

**Arguments**

container	Class of type <a href="#">matrix_container-class</a> generated by the <a href="#">create_container</a> function.
classification_results	A <code>cbind()</code> of result objects returned by <a href="#">classify_model</a> , or the object returned by <a href="#">classify_models</a> .
b	b-value for generating precision, recall, and F-scores statistics.

**Value**

Object of class [analytics\\_virgin-class](#) or [analytics-class](#) has either two or four slots respectively, depending on whether the virgin flag is set to TRUE or FALSE in [create\\_container](#). They can be accessed using the @ operator for S4 classes (e.g. `analytics@document_summary`).

**Author(s)**

Timothy P. Jurka <tpjurka@ucdavis.edu>, Loren Collingwood <lorenc2@uw.edu>

## Examples

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100, size=100, replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"], data["Subject"]), language="english",
  removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix, data$Topic.Code, trainSize=1:75, testSize=76:100,
  virgin=FALSE)
models <- train_models(container, algorithms=c("MAXENT", "SVM"))
results <- classify_models(container, models)
analytics <- create_analytics(container, results)
```

---

create_container	<i>creates a container for training, classifying, and analyzing documents.</i>
------------------	--

---

## Description

Given a DocumentTermMatrix from the **tm** package and corresponding document labels, creates a container of class `matrix_container-class` that can be used for training and classification (i.e. `train_model`, `train_models`, `classify_model`, `classify_models`)

## Usage

```
create_container(matrix, labels, trainSize=NULL, testSize=NULL, virgin)
```

## Arguments

matrix	A document-term matrix of class DocumentTermMatrix or TermDocumentMatrix from the <b>tm</b> package, or generated by <code>create_matrix</code> .
labels	A factor or vector of labels corresponding to each document in the matrix.
trainSize	A range (e.g. 1:1000) specifying the number of documents to use for training the models. Can be left blank for classifying corpora using saved models that don't need to be trained.
testSize	A range (e.g. 1:1000) specifying the number of documents to use for classification. Can be left blank for training on all data in the matrix.
virgin	A logical (TRUE or FALSE) specifying whether to treat the classification data as virgin data or not.

## Value

A container of class `matrix_container-class` that can be passed into other functions such as `train_model`, `train_models`, `classify_model`, `classify_models`, and `create_analytics`.

## Author(s)

Timothy P. Jurka <tpjurka@ucdavis.edu>, Loren Collingwood <lorenc2@uw.edu>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100, size=100, replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"], data["Subject"]), language="english",
  removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix, data$Topic.Code, trainSize=1:75, testSize=76:100,
  virgin=FALSE)
```

---

```
create_ensembleSummary
```

*creates a summary with ensemble coverage and precision.*

---

**Description**

Creates a summary with ensemble coverage and precision values for an ensemble greater than the threshold specified.

**Usage**

```
create_ensembleSummary(document_summary)
```

**Arguments**

```
document_summary
```

The document\_summary slot from the [analytics-class](#) generated by [create\\_analytics](#).

**Details**

This summary is created in the [create\\_analytics](#) function. Note that a threshold value of 3 will return ensemble coverage and precision statistics for topic codes that had 3 or more (i.e.  $\geq 3$ ) algorithms agree on the same topic code.

**Author(s)**

Loren Collingwood <lorenc2@uw.edu>, Timothy P. Jurka <tpjurka@ucdavis.edu>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100, size=100, replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"], data["Subject"]), language="english",
  removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix, data$Topic.Code, trainSize=1:75, testSize=76:100,
  virgin=FALSE)
models <- train_models(container, algorithms=c("MAXENT", "SVM"))
results <- classify_models(container, models)
analytics <- create_analytics(container, results)
```



```
ensemble <- create_ensembleSummary(analytics@document_summary)
ensemble
```

---

create\_matrix                      *creates a document-term matrix to be passed into create\_container().*

---

## Description

Creates an object of class `DocumentTermMatrix` from **tm** that can be used in the `create_container` function.

## Usage

```
create_matrix(textColumns, language="english", minDocFreq=1, maxDocFreq=Inf,
minWordLength=3, maxWordLength=Inf, ngramLength=1, originalMatrix=NULL,
removeNumbers=FALSE, removePunctuation=TRUE, removeSparseTerms=0,
removeStopwords=TRUE, stemWords=FALSE, stripWhitespace=TRUE, toLower=TRUE,
weighting=weightTf)
```

## Arguments

<code>textColumns</code>	Either character vector (e.g. <code>data\$Title</code> ) or a <code>cbind()</code> of columns to use for training the algorithms (e.g. <code>cbind(data\$Title,data\$Subject)</code> ).
<code>language</code>	The language to be used for stemming the text data.
<code>minDocFreq</code>	The minimum number of times a word should appear in a document for it to be included in the matrix. See package <b>tm</b> for more details.
<code>maxDocFreq</code>	The maximum number of times a word should appear in a document for it to be included in the matrix. See package <b>tm</b> for more details.
<code>minWordLength</code>	The minimum number of letters a word or n-gram should contain to be included in the matrix. See package <b>tm</b> for more details.
<code>maxWordLength</code>	The maximum number of letters a word or n-gram should contain to be included in the matrix. See package <b>tm</b> for more details.
<code>ngramLength</code>	The number of words to include per n-gram for the document-term matrix.
<code>originalMatrix</code>	The original <code>DocumentTermMatrix</code> used to train the models. If supplied, will adjust the new matrix to work with saved models.
<code>removeNumbers</code>	A logical parameter to specify whether to remove numbers.
<code>removePunctuation</code>	A logical parameter to specify whether to remove punctuation.
<code>removeSparseTerms</code>	See package <b>tm</b> for more details.
<code>removeStopwords</code>	A logical parameter to specify whether to remove stopwords using the language specified in <code>language</code> .

stemWords	A logical parameter to specify whether to stem words using the language specified in language.
stripWhitespace	A logical parameter to specify whether to strip whitespace.
toLower	A logical parameter to specify whether to make all text lowercase.
weighting	Either weightTf or tm::weightTfIdf. See package <b>tm</b> for more details.

**Author(s)**

Timothy P. Jurka <tpjurka@ucdavis.edu>, Loren Collingwood <lorenc2@uw.edu>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100, size=100, replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"], data["Subject"]), language="english",
removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
```

---

create\_precisionRecallSummary

*creates a summary with precision, recall, and F1 scores.*

---

**Description**

Creates a summary with precision, recall, and F1 scores for each algorithm broken down by unique label.

**Usage**

```
create_precisionRecallSummary(container, classification_results, b_value = 1)
```

**Arguments**

container	Class of type <code>matrix_container-class</code> generated by the <code>create_container</code> function.
classification_results	A <code>cbind()</code> of result objects returned by <code>classify_model</code> , or the object returned by <code>classify_models</code> .
b_value	b-value for generating precision, recall, and F-scores statistics.

**Author(s)**

Loren Collingwood <lorenc2@uw.edu>, Timothy P. Jurka <tpjurka@ucdavis.edu>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100,size=100,replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"],data["Subject"]), language="english",
removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix,data$Topic.Code,trainSize=1:75, testSize=76:100,
virgin=FALSE)
models <- train_models(container, algorithms=c("MAXENT","SVM"))
results <- classify_models(container, models)
precision_recall_f1 <- create_precisionRecallSummary(container, results)
```

---

create\_scoreSummary    *creates a summary with the best label for each document.*

---

**Description**

Creates a summary with the best label for each document, determined by highest algorithm certainty, and highest consensus (i.e. most number of algorithms agreed).

**Usage**

```
create_scoreSummary(container, classification_results)
```

**Arguments**

**container**            Class of type `matrix_container-class` generated by the `create_container` function.

**classification\_results**    A `cbind()` of result objects returned by `classify_model`, or the object returned by `classify_models`.

**Author(s)**

Timothy P. Jurka <tpjurka@ucdavis.edu>, Loren Collingwood <lorenc2@uw.edu>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100,size=100,replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"],data["Subject"]), language="english",
removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix,data$Topic.Code,trainSize=1:75, testSize=76:100,
virgin=FALSE)
models <- train_models(container, algorithms=c("MAXENT","SVM"))
results <- classify_models(container, models)
score_summary <- create_scoreSummary(container, results)
```

---

cross\_validate      *used for cross-validation of various algorithms.*

---

## Description

Performs n-fold cross-validation of specified algorithm.

## Usage

```
cross_validate(container, nfold, algorithm = c("SVM", "SLDA", "BOOSTING",
"BAGGING", "RF", "GLMNET", "TREE", "NNET", "MAXENT"), seed = NA,
method = "C-classification", cross = 0, cost = 100, kernel = "radial",
maxitboost = 100, maxitglm = 10^5, size = 1, maxitnnet = 1000, MaxNWts = 10000,
rang = 0.1, decay = 5e-04, ntree = 200, l1_regularizer = 0, l2_regularizer = 0,
use_sgd = FALSE, set_heldout = 0, verbose = FALSE)
```

## Arguments

container	Class of type <a href="#">matrix_container-class</a> generated by the <a href="#">create_container</a> function.
nfold	Number of folds to perform for cross-validation.
algorithm	A string specifying which algorithm to use. Use <a href="#">print_algorithms</a> to see a list of options.
seed	Random seed number used to replicate cross-validation results.
method	Method parameter for SVM implementation. See <b>e1071</b> documentation for more details.
cross	Cross parameter for SVM implementation. See <b>e1071</b> documentation for more details.
cost	Cost parameter for SVM implementation. See <b>e1071</b> documentation for more details.
kernel	Kernel parameter for SVM implementation. See <b>e1071</b> documentation for more details.
maxitboost	Maximum iterations parameter for boosting implementation. See <b>caTools</b> documentation for more details.
maxitglm	Maximum iterations parameter for glmnet implementation. See <b>glmnet</b> documentation for more details.
size	Size parameter for neural networks implementation. See <b>nnet</b> documentation for more details.
maxitnnet	Maximum iterations for neural networks implementation. See <b>nnet</b> documentation for more details.
MaxNWts	Maximum number of weights parameter for neural networks implementation. See <b>nnet</b> documentation for more details.

rang	Range parameter for neural networks implementation. See <b>nnet</b> documentation for more details.
decay	Decay parameter for neural networks implementation. See <b>nnet</b> documentation for more details.
nree	Number of trees parameter for RandomForests implementation. See <b>randomForest</b> documentation for more details.
l1_regularizer	An numeric turning on L1 regularization and setting the regularization parameter. A value of 0 will disable L1 regularization. See <b>maxent</b> documentation for more details.
l2_regularizer	An numeric turning on L2 regularization and setting the regularization parameter. A value of 0 will disable L2 regularization. See <b>maxent</b> documentation for more details.
use_sgd	A logical indicating that SGD parameter estimation should be used. Defaults to FALSE. See <b>maxent</b> documentation for more details.
set_heldout	An integer specifying the number of documents to hold out. Sets a held-out subset of your data to test against and prevent overfitting. See <b>maxent</b> documentation for more details.
verbose	A logical specifying whether to provide descriptive output about the training process. Defaults to FALSE, or no output. See <b>maxent</b> documentation for more details.

**Author(s)**

Loren Collingwood <lorenc2@uw.edu>, Timothy P. Jurka <tpjurka@ucdavis.edu>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100, size=100, replace=FALSE), ]
matrix <- create_matrix(cbind(data["Title"], data["Subject"]), language="english",
  removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix, data$Topic.Code, trainSize=1:75, testSize=76:100,
  virgin=FALSE)
svm <- cross_validate(container, 2, algorithm="SVM")
maxent <- cross_validate(container, 2, algorithm="MAXENT")
```

---

getStemLanguages

*Query the languages supported in this package*


---

**Description**

This dynamically determines the names of the languages for which stemming is supported by this package. This is controlled when the package is created (not installed) by downloading the stemming algorithms for the different languages.

This language support requires more support for Unicode and more complex text than simple strings.

**Usage**

```
getStemLanguages()
```

**Details**

This queries the C code for the list of languages that were compiled when the package was installed which in turn is determined by the code that was included in the distributed package itself.

**Value**

A character vector giving the names of the languages.

**Author(s)**

Duncan Temple Lang <duncan@wald.ucdavis.edu>

**References**

See <http://snowball.tartarus.org/>

**See Also**

[wordStem](#) inst/scripts/download in the source of the Rstem package.

**Examples**

```
getStemLanguages()
```

---

matrix\_container-class

*an S4 class containing the training and classification matrices.*

---

**Description**

An S4 class containing all information necessary to train, classify, and generate analytics for a dataset.

**Objects from the Class**

Objects could in principle be created by calls of the form `new("matrix_container", ...)`. The preferred form is to have them created via a call to `create_container`.

**Slots**

`training_matrix` Object of class "matrix.csr": stores the training set of the DocumentTermMatrix created by `create_matrix`

`training_codes` Object of class "factor": stores the training labels for each document in the `training_matrix` slot of `matrix_container-class`

`classification_matrix` Object of class "matrix.csr": stores the classification set of the DocumentTermMatrix created by `create_matrix`

`testing_codes` Object of class "factor": if `virgin=FALSE`, stores the labels for each document in `classification_matrix`

`column_names` Object of class "vector": stores the column names of the DocumentTermMatrix created by `create_matrix`

`virgin` Object of class "logical": boolean specifying whether the classification set is virgin data (TRUE) or not (FALSE).

**Author(s)**

Timothy P. Jurka <tpjurka@ucdavis.edu>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100, size=100, replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"], data["Subject"]), language="english",
  removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix, data$Topic.Code, trainSize=1:75, testSize=76:100,
  virgin=FALSE)

container@training_matrix
container@training_codes
container@classification_matrix
container@testing_codes
container@column_names
container@virgin
```

---

NYTimes	<i>a sample dataset containing labeled headlines from The New York Times.</i>
---------	---

---

**Description**

A sample dataset containing labeled headlines from The New York Times, compiled by Professor Amber E. Boydston at the University of California, Davis.

**Usage**

```
data(NYTimes)
```

**Format**

A data.frame containing five columns.

1. Article\_ID - A unique identifier for the headline from The New York Times.
2. Date - The date the headline appeared in The New York Times.
3. Title - The headline as it appeared in The New York Times.
4. Subject - A manually classified subject of the headline.
5. Topic.Code - A manually labeled topic code corresponding to the subject.

**Source**

<http://www.amberboystun.com/>

**Examples**

```
data(NYTimes)
```

---

```
print_algorithms      prints available algorithms for train_model() and train_models().
```

---

**Description**

An informative function that displays options for the algorithms parameter in `train_model` and `train_models`.

**Usage**

```
print_algorithms()
```

**Value**

Prints a list of available algorithms.

**Author(s)**

Timothy P. Jurka <tpjurka@ucdavis.edu>

**Examples**

```
library(RTextTools)
print_algorithms()
```



---

read_data	<i>reads data from files into an R data frame.</i>
-----------	--

---

### Description

Reads data from several types of data storage types into an R data frame.

### Usage

```
read_data(filepath, type=c("csv","delim","folder"), index=NULL, ...)
```

### Arguments

filepath	Character string of the name of the file or folder, include path if the file is not located in the working directory.
type	Character vector specifying the file type. Options include <code>csv</code> , <code>delim</code> , and <code>folder</code> to denote <code>.csv</code> files, delimited files ( <code>tab</code> , <code>pipe</code> , etc.) files, or folders of text files. If using the <code>delim</code> option, be sure to pass in a separate <code>sep</code> parameter to indicate how the file is delimited.
index	The path to a CSV file specifying the training label of each file in the folder of text files, one per line. An example of one line would be <code>1.txt,1</code> . Do not include the full file path for each file, that will be handled automatically using the folder location passed into <code>filepath</code> . This index file must be located outside the folder of files.
...	Other arguments passed to R's <code>read.csv</code> function.

### Value

An `data.frame` object is returned with the contents of the file.

### Author(s)

Loren Collingwood <lorenc2@uw.edu>, Timothy P. Jurka <tpjurka@ucdavis.edu>

### Examples

```
library(RTextTools)
data <- read_data(system.file("data/NYTimes.csv.gz", package="RTextTools"), type="csv", sep=";")
```

---

recall_accuracy	<i>calculates the recall accuracy of the classified data.</i>
-----------------	---

---

### Description

Given the true labels to compare to the labels predicted by the algorithms, calculates the recall accuracy of each algorithm.

### Usage

```
recall_accuracy(true_labels, predicted_labels)
```

### Arguments

`true_labels` A vector containing the true labels, or known values for each document in the classification set.

`predicted_labels` A vector containing the predicted labels, or classified values for each document in the classification set.

### Author(s)

Loren Collingwood <lorenc2@uw.edu>, Timothy P. Jurka <tpjurka@ucdavis.edu>

### Examples

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100, size=100, replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"], data["Subject"]), language="english",
  removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix, data$Topic.Code, trainSize=1:75, testSize=76:100,
  virgin=FALSE)
models <- train_models(container, algorithms=c("MAXENT", "SVM"))
results <- classify_models(container, models)
analytics <- create_analytics(container, results)
recall_accuracy(analytics@document_summary$MANUAL_CODE,
  analytics@document_summary$GLMNET_LABEL)
recall_accuracy(analytics@document_summary$MANUAL_CODE,
  analytics@document_summary$MAXENTROPY_LABEL)
recall_accuracy(analytics@document_summary$MANUAL_CODE,
  analytics@document_summary$SVM_LABEL)
```

---

summary.analytics      *summarizes the [analytics-class](#) class*

---

### Description

Returns a summary of the contents within an object of class [analytics-class](#).

### Usage

```
## S3 method for class 'analytics'  
summary(object, ...)
```

### Arguments

object      An object of class [analytics-class](#) containing the output of the [create\\_analytics](#) function.

...      Additional parameters to be passed onto the summary function.

### Author(s)

Timothy P. Jurka <tpjurka@ucdavis.edu>

### Examples

```
library(RTextTools)  
data(NYTimes)  
data <- NYTimes[sample(1:3100, size=100, replace=FALSE),]  
matrix <- create_matrix(cbind(data["Title"], data["Subject"]), language="english",  
  removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)  
container <- create_container(matrix, data$Topic.Code, trainSize=1:75, testSize=76:100,  
  virgin=FALSE)  
models <- train_models(container, algorithms=c("MAXENT", "SVM"))  
results <- classify_models(container, models)  
analytics <- create_analytics(container, results)  
  
summary(analytics)
```

---

summary.analytics\_virgin      *summarizes the [analytics\\_virgin-class](#) class*

---

### Description

Returns a summary of the contents within an object of class [analytics\\_virgin-class](#).

**Usage**

```
## S3 method for class 'analytics_virgin'
summary(object, ...)
```

**Arguments**

**object** An object of class `analytics_virgin-class` containing the output of the `create_analytics` function.

**...** Additional parameters to be passed onto the summary function.

**Author(s)**

Timothy P. Jurka <tpjurka@ucdavis.edu>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100,size=100,replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"],data["Subject"]), language="english",
removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix,data$Topic.Code,trainSize=1:75, testSize=76:100,
virgin=TRUE)
models <- train_models(container, algorithms=c("MAXENT","SVM"))
results <- classify_models(container, models)
analytics <- create_analytics(container, results)

summary(analytics)
```

---

train\_model

*makes a model object using the specified algorithm.*

---

**Description**

Creates a trained model using the specified algorithm.

**Usage**

```
train_model(container, algorithm=c("SVM","SLDA","BOOSTING","BAGGING",
"RF","GLMNET","TREE","NNET","MAXENT"), method = "C-classification",
cross = 0, cost = 100, kernel = "radial", maxitboost = 100,
maxitglm = 10^5, size = 1, maxitnnet = 1000, MaxNWts = 10000,
rang = 0.1, decay = 5e-04, trace=FALSE, ntree = 200,
l1_regularizer = 0, l2_regularizer = 0, use_sgd = FALSE,
set_heldout = 0, verbose = FALSE,
...)
```

**Arguments**

container	Class of type <code>matrix_container-class</code> generated by the <code>create_container</code> function.
algorithm	Character vector (i.e. a string) specifying which algorithm to use. Use <code>print_algorithms</code> to see a list of options.
method	Method parameter for SVM implementation. See <b>e1071</b> documentation for more details.
cross	Cross parameter for SVM implementation. See <b>e1071</b> documentation for more details.
cost	Cost parameter for SVM implementation. See <b>e1071</b> documentation for more details.
kernel	Kernel parameter for SVM implementation. See <b>e1071</b> documentation for more details.
maxitboost	Maximum iterations parameter for boosting implementation. See <b>caTools</b> documentation for more details.
maxitglm	Maximum iterations parameter for glmnet implementation. See <b>glmnet</b> documentation for more details.
size	Size parameter for neural networks implementation. See <b>nnet</b> documentation for more details.
maxitnnet	Maximum iterations for neural networks implementation. See <b>nnet</b> documentation for more details.
MaxNWts	Maximum number of weights parameter for neural networks implementation. See <b>nnet</b> documentation for more details.
rang	Range parameter for neural networks implementation. See <b>nnet</b> documentation for more details.
decay	Decay parameter for neural networks implementation. See <b>nnet</b> documentation for more details.
trace	Trace parameter for neural networks implementation. See <b>nnet</b> documentation for more details.
nntree	Number of trees parameter for RandomForests implementation. See <b>randomForest</b> documentation for more details.
l1_regularizer	An numeric turning on L1 regularization and setting the regularization parameter. A value of 0 will disable L1 regularization. See <b>maxent</b> documentation for more details.
l2_regularizer	An numeric turning on L2 regularization and setting the regularization parameter. A value of 0 will disable L2 regularization. See <b>maxent</b> documentation for more details.
use_sgd	A logical indicating that SGD parameter estimation should be used. Defaults to FALSE. See <b>maxent</b> documentation for more details.
set_heldout	An integer specifying the number of documents to hold out. Sets a held-out subset of your data to test against and prevent overfitting. See <b>maxent</b> documentation for more details.

`verbose` A logical specifying whether to provide descriptive output about the training process. Defaults to FALSE, or no output. See **maxent** documentation for more details.

`...` Additional arguments to be passed on to algorithm function calls.

### Details

Only one algorithm may be selected for training. See [train\\_models](#) and [classify\\_models](#) to train and classify using multiple algorithms.

### Value

Returns a trained model that can be subsequently used in [classify\\_model](#) to classify new data.

### Author(s)

Timothy P. Jurka <tpjurka@ucdavis.edu>, Loren Collingwood <lorenc2@uw.edu>

### Examples

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100,size=100,replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"],data["Subject"]), language="english",
  removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix,data$Topic.Code,trainSize=1:75, testSize=76:100,
  virgin=FALSE)
maxent_model <- train_model(container,"MAXENT")
svm_model <- train_model(container,"SVM")
```

---

`train_models` *makes a model object using the specified algorithms.*

---

### Description

Creates a trained model using the specified algorithms.

### Usage

```
train_models(container, algorithms, ...)
```

### Arguments

`container` Class of type [matrix\\_container-class](#) generated by the [create\\_container](#) function.

`algorithms` List of algorithms as a character vector (e.g. `c("SVM", "MAXENT")`).

`...` Other parameters to be passed on to [train\\_model](#).

**Details**

Calls the `train_model` function for each algorithm you list.

**Value**

Returns a list of trained models that can be subsequently used in `classify_models` to classify new data.

**Author(s)**

Wouter Van Atteveldt <wouter@vanatteveldt.com>

**Examples**

```
library(RTextTools)
data(NYTimes)
data <- NYTimes[sample(1:3100,size=100,replace=FALSE),]
matrix <- create_matrix(cbind(data["Title"],data["Subject"]), language="english",
  removeNumbers=TRUE, stemWords=FALSE, weighting=tm::weightTfIdf)
container <- create_container(matrix,data$Topic.Code,trainSize=1:75, testSize=76:100,
  virgin=FALSE)
models <- train_models(container, algorithms=c("MAXENT","SVM"))
```

---

USCongress	<i>a sample dataset containing labeled bills from the United State Congress.</i>
------------	--

---

**Description**

A sample dataset containing labeled bills from the United States Congress, compiled by Professor John D. Wilkerson at the University of Washington, Seattle and E. Scott Adler at the University of Colorado, Boulder.

**Usage**

```
data(USCongress)
```

**Format**

A data.frame containing five columns.

1. ID - A unique identifier for the bill.
2. cong - The session of congress that the bill first appeared in.
3. billnum - The number of the bill as it appears in the congressional docket.
4. h\_or\_sen - A field specifying whether the bill was introduced in the House (HR) or the Senate (S).
5. major - A manually labeled topic code corresponding to the subject of the bill.

**Source**

<http://www.congressionalbills.org/>

**Examples**

```
data(USCongress)
```

---

wordStem

*Get the common root/stem of words*

---

**Description**

This function computes the stems of each of the given words in the vector. This reduces a word to its base component, making it easier to compare words like win, winning, winner. See <http://snowball.tartarus.org/> for more information about the concept and algorithms for stemming.

**Usage**

```
wordStem(words, language = character(), warnTested = FALSE)
```

**Arguments**

words	a character vector of words whose stems are to be computed.
language	the name of a recognized language for the package. This should either be a single string which is an element in the vector returned by <a href="#">getStemLanguages</a> , or alternatively a character vector of length 3 giving the names of the routines for creating and closing a Snowball SN_env environment and performing the stem (in that order). See the example below.
warnTested	an option to control whether a warning is issued about languages which have not been explicitly tested as part of the unit testing of the code. For the most part, one can ignore these warnings and so they are turned off. In the future, we might consider controlling this with a global option, but for now we suppress the warnings by default.

**Details**

This uses Dr. Martin Porter's stemming algorithm and the interface generated by Snowball <http://snowball.tartarus.org/>.

**Value**

A character vector with as many elements as there are in the input vector with the corresponding elements being the stem of the word.

**Author(s)**

Duncan Temple Lang <duncan@wald.ucdavis.edu>



## References

See <http://snowball.tartarus.org/>

## Examples

```
# Simple example
# "win"      "win"      "winner"
wordStem(c("win", "winning", 'winner'))

# test the supplied vocabulary.
testWords = readLines(system.file("words", "english", "voc.txt", package = "RTextTools"))
validate = readLines(system.file("words", "english", "output.txt", package = "RTextTools"))

## Not run:
# Read the test words directly from the snowball site over the Web
testWords = readLines(url("http://snowball.tartarus.org/english/voc.txt"))

## End(Not run)

testOut = wordStem(testWords)
all(validate == testOut)

# Specify the language from one of the built-in languages.
testOut = wordStem(testWords, "english")
all(validate == testOut)

# To illustrate using the dynamic lookup of symbols that allows one
# to easily add new languages or create and close environment
# routines (for example, to manage pools if this were an efficiency
# issue!)
testOut = wordStem(testWords, c("testDynCreate", "testDynClose", "testDynStem"))
```

# Index

## \*Topic **IO**

getStemLanguages, 13  
wordStem, 24

## \*Topic **classes**

analytics-class, 2  
analytics\_virgin-class, 3  
matrix\_container-class, 14

## \*Topic **datasets**

NYTimes, 15  
USCongress, 23

## \*Topic **method**

classify\_model, 4  
classify\_models, 5  
create\_analytics, 6  
create\_container, 7  
create\_ensembleSummary, 8  
create\_matrix, 9  
create\_precisionRecallSummary, 10  
create\_scoreSummary, 11  
cross\_validate, 12  
print\_algorithms, 16  
read\_data, 17  
recall\_accuracy, 18  
summary.analytics, 19  
summary.analytics\_virgin, 19  
train\_model, 20  
train\_models, 22

## \*Topic **utilities**

getStemLanguages, 13  
wordStem, 24

analytics-class, 2, 19  
analytics\_virgin-class, 3, 19

classify\_model, 4, 5–7, 10, 11, 22  
classify\_models, 4, 5, 6, 7, 10, 11, 22, 23  
create\_analytics, 2, 3, 6, 7, 8, 19, 20  
create\_container, 2–6, 7, 9–12, 14, 21, 22  
create\_ensembleSummary, 8  
create\_matrix, 7, 9, 15

create\_precisionRecallSummary, 10  
create\_scoreSummary, 11  
cross\_validate, 12

getStemLanguages, 13, 24

matrix\_container-class, 14

NYTimes, 15

print\_algorithms, 12, 16, 21

read\_data, 17  
recall\_accuracy, 18

summary.analytics, 19  
summary.analytics\_virgin, 19

train\_model, 4, 7, 16, 20, 22, 23  
train\_models, 4, 5, 7, 16, 22, 22

USCongress, 23

wordStem, 14, 24