

# Package ‘wordvector’

December 11, 2024

**Type** Package

**Title** Word and Document Vector Models

**Version** 0.1.0

**Maintainer** Kohei Watanabe <watanabe.kohei@gmail.com>

**Description** Create dense vector representation of words and documents using 'quanteda'. Currently implements Word2vec (Mikolov et al., 2013) <[doi:10.48550/arXiv.1310.4546](https://doi.org/10.48550/arXiv.1310.4546)> and Latent Semantic Analysis (Deerwester et al., 1990) <[doi:10.1002/\(SICI\)1097-4571\(199009\)41:6%3C391::AID-AS11%3E3.0.CO;2-9](https://doi.org/10.1002/(SICI)1097-4571(199009)41:6%3C391::AID-AS11%3E3.0.CO;2-9)>.

**URL** <https://github.com/koheiw/wordvector>

**License** Apache License (>= 2.0)

**Encoding** UTF-8

**RoxygenNote** 7.3.2

**Depends** R (>= 3.5.0)

**Imports** quanteda (>= 4.1.0), methods, stringi, Matrix, proxyC, RSpectra, irlba, rsvd

**Suggests** testthat, word2vec, spelling

**LinkingTo** Rcpp, quanteda

**Language** en-US

**LazyData** true

**NeedsCompilation** yes

**Author** Kohei Watanabe [aut, cre, cph]  
(<<https://orcid.org/0000-0001-6519-5265>>),  
Jan Wijffels [aut] (Original R code),  
BNOSAC [cph] (Original R code),  
Max Fomichev [ctb, cph] (Original C++ code)

**Repository** CRAN

**Date/Publication** 2024-12-11 17:00:02 UTC

## Contents

analogy	2
as.matrix.textmodel_wordvector	3
data_corpus_news2014	3
doc2vec	4
lsa	4
similarity	6
word2vec	6
<b>Index</b>	<b>9</b>

---

analogy	<i>[experimental] Find analogical relationships between words</i>
---------	---

---

### Description

[experimental] Find analogical relationships between words

### Usage

```
analogy(x, formula, n = 10, exclude = TRUE, type = c("word", "simil"))
```

### Arguments

x	a textmodel_wordvector object.
formula	a <a href="#">formula</a> object that defines the relationship between words using + or - operators.
n	the number of words in the resulting object.
exclude	if TRUE, words in formula are excluded from the result.
type	specify the type of vectors to be used. "word" is word vectors while "simil" is similarity vectors.

### Value

a data.frame with the words sorted and their cosine similarity sorted in descending order.

### References

Mikolov, T., Sutskever, I., Chen, K., Corrado, G., & Dean, J. (2013). Distributed Representations of Words and Phrases and their Compositionality. <http://arxiv.org/abs/1310.4546>.

### Examples

```
# from Mikolov et al. (2023)
analogy(wdv, ~ berlin - germany + france)
analogy(wdv, ~ quick - quickly + slowly)
```

---

```
as.matrix.textmodel_wordvector
```

*Extract word vectors*

---

**Description**

Extract word vectors from a textmodel\_wordvector or textmodel\_docvector object.

**Usage**

```
## S3 method for class 'textmodel_wordvector'  
as.matrix(x, ...)
```

**Arguments**

x	a textmodel_wordvector or textmodel_docvector object.
...	not used

**Value**

a matrix that contain the word vectors in rows

---

```
data_corpus_news2014
```

*Yahoo News summaries from 2014*

---

**Description**

A corpus object containing 2,000 news summaries collected from Yahoo News via RSS feeds in 2014. The title and description of the summaries are concatenated.

**Usage**

```
data_corpus_news2014
```

**Format**

An object of class corpus (inherits from character) of length 20000.

**Source**

<https://www.yahoo.com/news/>

**References**

Watanabe, K. (2018). Newsmap: A semi-supervised approach to geographical news classification. *Digital Journalism*, 6(3), 294–309. <https://doi.org/10.1080/21670811.2017.1293487>

---

doc2vec	<i>Create distributed representation of documents</i>
---------	---

---

**Description**

Create distributed representation of documents

**Usage**

```
doc2vec(x, model = NULL, ...)
```

**Arguments**

x	a <a href="#">quanteda::tokens</a> object.
model	a textmodel_wordvector object.
...	passed to [word2vec] when model = NULL.

**Value**

Returns a textmodel\_docvector object with elements inherited from model or passed via ... plus:

vectors	a matrix for document vectors.
call	the command used to execute the function.

---

lsa	<i>Latent Semantic Analysis model</i>
-----	---------------------------------------

---

**Description**

Train a Latent Semantic Analysis model (Deerwester et al., 1990) on a [quanteda::tokens](#) object.

**Usage**

```
lsa(  
  x,  
  dim = 50,  
  min_count = 5L,  
  engine = c("RSpectra", "irlba", "rsvd"),  
  weight = "count",  
  verbose = FALSE,  
  ...  
)
```

**Arguments**

<code>x</code>	a <code>quanteda::tokens</code> object.
<code>dim</code>	the size of the word vectors.
<code>min_count</code>	the minimum frequency of the words. Words less frequent than this in <code>x</code> are removed before training.
<code>engine</code>	select the engine perform SVD to generate word vectors.
<code>weight</code>	weighting scheme passed to <code>quanteda::dfm_weight()</code> .
<code>verbose</code>	if TRUE, print the progress of training.
<code>...</code>	additional arguments.

**Value**

Returns a `textmodel_wordvector` object with the following elements:

<code>vectors</code>	a matrix for word vectors.
<code>frequency</code>	the frequency of words in <code>x</code> .
<code>engine</code>	the SVD engine used.
<code>weight</code>	weighting scheme.
<code>concatenator</code>	the concatenator in <code>x</code> .
<code>call</code>	the command used to execute the function.
<code>version</code>	the version of the wordvector package.

**References**

Deerwester, S. C., Dumais, S. T., Landauer, T. K., Furnas, G. W., & Harshman, R. A. (1990). Indexing by latent semantic analysis. *JASIS*, 41(6), 391–407.

**Examples**

```
library(quanteda)
library(wordvector)

# pre-processing
corp <- corpus_reshape(data_corpus_news2014)
toks <- tokens(corp, remove_punct = TRUE, remove_symbols = TRUE) %>%
  tokens_remove(stopwords("en", "marimo"), padding = TRUE) %>%
  tokens_select("^[a-zA-Z-]+$", valuetype = "regex", case_insensitive = FALSE,
               padding = TRUE) %>%
  tokens_tolower()

# train LSA
lsa <- lsa(toks, dim = 50, min_count = 5, verbose = TRUE, )
head(similarity(lsa, c("berlin", "germany", "france"), mode = "word"))
analogy(lsa, ~ berlin - germany + france)
```

---

similarity	<i>Compute similarity between word vectors</i>
------------	--

---

### Description

Compute similarity between word vectors

### Usage

```
similarity(x, words, mode = c("simil", "word"))
```

### Arguments

x	a <code>textmodel_wordvector</code> object.
words	words for which similarity is computed.
mode	specify the type of resulting object.

### Value

a matrix of cosine similarity scores when `mode = "simil"` or of words sorted by the similarity scores when `mode = "word"`.

---

word2vec	<i>Word2vec model</i>
----------	-----------------------

---

### Description

Train a Word2vec model (Mikolov et al., 2023) in different architectures on a `quanteda::tokens` object.

### Usage

```
word2vec(
  x,
  dim = 50,
  type = c("cbow", "skip-gram"),
  min_count = 5L,
  window = ifelse(type == "cbow", 5L, 10L),
  iter = 10L,
  alpha = 0.05,
  use_ns = TRUE,
  ns_size = 5L,
  sample = 0.001,
  verbose = FALSE,
  ...
)
```

**Arguments**

x	a <code>quanteda::tokens</code> object.
dim	the size of the word vectors.
type	the architecture of the model; either "cbow" (continuous back of words) or "skip-gram".
min_count	the minimum frequency of the words. Words less frequent than this in x are removed before training.
window	the size of the word window. Words within this window are considered to be the context of a target word.
iter	the number of iterations in model training.
alpha	the initial learning rate.
use_ns	if TRUE, negative sampling is used. Otherwise, hierarchical softmax is used.
ns_size	the size of negative samples. Only used when use_ns = TRUE.
sample	the rate of sampling of words based on their frequency. Sampling is disabled when sample = 1.0
verbose	if TRUE, print the progress of training.
...	additional arguments.

**Details**

User can changed the number of processors used for the parallel computing via `options(wordvector_threads)`.

**Value**

Returns a `textmodel_wordvector` object with the following elements:

vectors	a matrix for word vectors.
dim	the size of the word vectors.
type	the architecture of the model.
frequency	the frequency of words in x.
window	the size of the word window.
iter	the number of iterations in model training.
alpha	the initial learning rate.
use_ns	the use of negative sampling.
ns_size	the size of negative samples.
concatenator	the concatenator in x.
call	the command used to execute the function.
version	the version of the wordvector package.

**References**

Mikolov, T., Sutskever, I., Chen, K., Corrado, G., & Dean, J. (2013). Distributed Representations of Words and Phrases and their Compositionality. <https://arxiv.org/abs/1310.4546>.

**Examples**

```
library(quanteda)
library(wordvector)

# pre-processing
corp <- data_corpus_news2014
toks <- tokens(corp, remove_punct = TRUE, remove_symbols = TRUE) %>%
  tokens_remove(stopwords("en", "marimo"), padding = TRUE) %>%
  tokens_select("[a-zA-Z-]+$", valuetype = "regex", case_insensitive = FALSE,
                padding = TRUE) %>%
  tokens_tolower()

# train word2vec
w2v <- word2vec(toks, dim = 50, type = "cbow", min_count = 5, sample = 0.001)
head(similarity(w2v, c("berlin", "germany", "france"), mode = "word"))
analogy(w2v, ~ berlin - germany + france)
```



# Index

## \* datasets

data\_corpus\_news2014, 3

analogy, 2

as.matrix.textmodel\_wordvector, 3

data\_corpus\_news2014, 3

doc2vec, 4

formula, 2

lsa, 4

quanteda::dfm\_weight(), 5

quanteda::tokens, 4-7

similarity, 6

word2vec, 6