

Machine Learning: Natural Language Processing

1 / 133

Speaker Qualifications

- Specialize in next-generation technologies
- Author of O'Reilly Videos on Hypermedia, Linking Data, Security and Encryption
- Author of 'Resource-Oriented Architecture Patterns for Webs of Data'
- Teaches and speaks internationally about REST, Semantic Web, Data Science, Machine Learning, GPU Computing, Security, Visualization, Architecture
- Worked in Defense, Finance, Retail, Hospitality, Video Game, Health Care, Telecommunications and Publishing Industries
- International Pop Recording Artist

2 / 133

Agenda

- Introduction
- Vector Space Model
- Word2Vec
- Naive Bayes

3 / 133

Introduction

4 / 133

Natural Language Processing (NLP) Goals

- Search and Retrieval

5 / 133

Natural Language Processing (NLP) Goals

- Search and Retrieval
- Entity and Relationship Extraction

6 / 133

Natural Language Processing (NLP) Goals

- Search and Retrieval
- Entity and Relationship Extraction
- Linguistic structure

7 / 133

Natural Language Processing (NLP) Goals

- Search and Retrieval
- Entity and Relationship Extraction
- Linguistic structure
- Machine Translation

8 / 133

Natural Language Processing (NLP) Goals

- Search and Retrieval
- Entity and Relationship Extraction
- Linguistic structure
- Machine Translation
- Generative Content

9 / 133

Natural Language Processing (NLP) Goals

- Search and Retrieval
- Entity and Relationship Extraction
- Linguistic structure
- Machine Translation
- Generative Content
- Question Answering

10 / 133

NLP History

- Early successes in the 1950s in automatic translation

11 / 133

NLP History

- Early successes in the 1950s in automatic translation
- SHRDLU

12 / 133

NLP History

- Early successes in the 1950s in automatic translation
- SHRDLU
- ELIZA

13 / 133

NLP History

- Early successes in the 1950s in automatic translation
- SHRDLU
- ELIZA
- Shift to statistical models in the 1980s

14 / 133

NLP Difficulties

15 / 133

NLP Difficulties

- Character encoding

16 / 133

NLP Difficulties

- Character encoding
- Tokenization

17 / 133

NLP Difficulties

- Character encoding
- Tokenization
- Part of Speech labeling

18 / 133

NLP Difficulties

- Character encoding
- Tokenization
- Part of Speech labeling
- Stemming (e.g. "walking", "walked", "walks")

19 / 133

NLP Difficulties

- Character encoding
- Tokenization
- Part of Speech labeling
- Stemming (e.g. "walking", "walked", "walks")
- Lemmatization (e.g. "operating")

20 / 133

NLP Difficulties

- Character encoding
- Tokenization
- Part of Speech labeling
- Stemming (e.g. "walking", "walked", "walks")
- Lemmatization (e.g. "operating")
- Sentence/paragraph detection

21 / 133

NLP Difficulties

- Character encoding
- Tokenization
- Part of Speech labeling
- Stemming (e.g. "walking", "walked", "walks")
- Lemmatization (e.g. "operating")
- Sentence/paragraph detection
- Coreference resolution

22 / 133

Some Open Source NLP Frameworks

API	URL
GATE	http://gate.ac.uk
LingPipe	http://alias-i.com/lingpipe
Apache OpenNLP	http://opennlp.apache.org
UIMA	http://uima.apache.org
Stanford Parser	http://nlp.stanford.edu/software
Mallet	http://mallet.cs.umass.edu

23 / 133

Vector Space Model

24 / 133

Bag of Words

- (1) John likes to watch movies. Mary likes movies too.
- (2) John also likes to watch football games.

https://en.wikipedia.org/wiki/Bag-of-words_model

25 / 133

Bag of Words

- (1) John likes to watch movies. Mary likes movies too.
- (2) John also likes to watch football games.

```
[  
  "John",  
  "likes",  
  "to",  
  "watch",  
  "movies",  
  "also",  
  "football",  
  "games",  
  "Mary",  
  "too"  
]
```

https://en.wikipedia.org/wiki/Bag-of-words_model

26 / 133

Bag of Words

- (1) John likes to watch movies. Mary likes movies too.
(2) John also likes to watch football games.

```
[  
  "John",  
  "likes",  
  "to",  
  "watch",  
  "movies",  
  "also",  
  "football",  
  "games",  
  "Mary",  
  "too"  
]
```

```
(1) [1, 2, 1, 1, 2, 0, 0, 0, 1, 1]  
(2) [1, 1, 1, 1, 0, 1, 1, 1, 0, 0]
```

https://en.wikipedia.org/wiki/Bag-of-words_model

27 / 133

N-gram model

- (1) John likes to watch movies. Mary likes movies too.
(2) John also likes to watch football games.

https://en.wikipedia.org/wiki/Bag-of-words_model

28 / 133

N-gram model

- (1) John likes to watch movies. Mary likes movies too.
- (2) John also likes to watch football games.

```
[  
  "John likes",  
  "likes to",  
  "to watch",  
  "watch movies",  
  "Mary likes",  
  "likes movies",  
  "movies too",  
]
```

https://en.wikipedia.org/wiki/Bag-of-words_model

29 / 133

Vector Space Model

$$d_j = (w_{1,j}, w_{2,j}, \dots w_{t,j})$$

30 / 133

Vector Space Model

$$d_j = (w_{1,j}, w_{2,j}, \dots w_{t,j})$$

$$q = (w_{1,q}, w_{2,q}, \dots w_{t,q})$$

31 / 133

Euclidean Dot Product

$$a \cdot b = \|a\| \|b\| \cos \theta$$

32 / 133

Euclidean Dot Product

$$a \cdot b = \|a\| \|b\| \cos \theta$$

$$\cos \theta = \frac{A \cdot B}{\|A\| \|B\|}$$

33 / 133

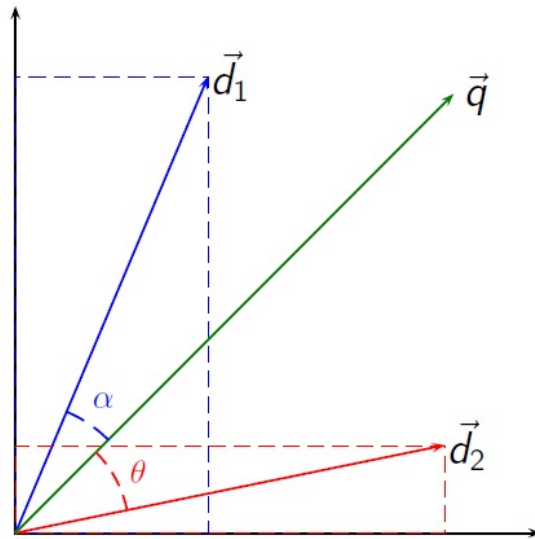
Euclidean Dot Product

$$a \cdot b = \|a\| \|b\| \cos \theta$$

$$\cos \theta = \frac{A \cdot B}{\|A\| \|B\|}$$

$$= \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

34 / 133



By Riclas - Own work, CC BY 3.0, [Link](#)

35 / 133

Cosine Similarity

$$\cos\theta = \frac{d_2 \cdot q}{\|d_2\| \|q\|}$$

36 / 133

Cosine Similarity

$$\cos\theta = \frac{d_2 \cdot q}{\|d_2\| \|q\|}$$

$$\cos\theta = 1 \implies \textit{Identical}$$

37 / 133

Cosine Similarity

$$\cos\theta = \frac{d_2 \cdot q}{\|d_2\| \|q\|}$$

$$\cos\theta = 1 \implies \textit{Identical}$$

$$\cos\theta = 0 \implies \textit{Orthogonal}$$

38 / 133

Term Frequency

$$tf(t, d) = 1$$

<https://en.wikipedia.org/wiki/Tf-idf>

39 / 133

Term Frequency

$$tf(t, d) = 1$$

$$tf(t, d) = f_{t,d}$$

<https://en.wikipedia.org/wiki/Tf-idf>

40 / 133

Term Frequency

$$tf(t, d) = 1$$

$$tf(t, d) = f_{t,d}$$

$$tf(t, d) = 0.5 + 0.5 \cdot \frac{f_{t,d}}{\max\{f_{t',d} : t' \in d\}}$$

<https://en.wikipedia.org/wiki/Tf-idf>

41 / 133

Inverse Document Frequency

$$idf(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

<https://en.wikipedia.org/wiki/Tf-idf>

42 / 133

Term Frequency-Inverse Document Frequency

$$tfidf(t, d, D) = tf(t, d) \cdot idf(t, D)$$

<https://en.wikipedia.org/wiki/Tf-idf>

43 / 133

Document 1

Term	Term Count
this	1
is	1
a	2
sample	1

Document 2

Term	Term Count
this	1
is	1
another	2
example	3

44 / 133

Document 1

Term	Term Count
this	1
is	1
a	2
sample	1

Term Frequency

$$tf("this", d_1) = \frac{1}{5} = 0.2$$

$$tf("this", d_2) = \frac{1}{7} \approx 0.14$$

Document 2

Term	Term Count
this	1
is	1
another	2
example	3

45 / 133

Document 1

Term	Term Count
this	1
is	1
a	2
sample	1

Term Frequency

$$tf("this", d_1) = \frac{1}{5} = 0.2$$

$$tf("this", d_2) = \frac{1}{7} \approx 0.14$$

Document 2

Term	Term Count
this	1
is	1
another	2
example	3

Inverse Document Frequency

$$idf("this", D) = \log\left(\frac{2}{2}\right) = 0$$

46 / 133

Document 1

Term	Term Count
this	1
is	1
a	2
sample	1

Term Frequency

$$tf("this", d_1) = \frac{1}{5} = 0.2$$

$$tf("this", d_2) = \frac{1}{7} \approx 0.14$$

Document 2

Term	Term Count
this	1
is	1
another	2
example	3

Inverse Document Frequency

$$idf("this", D) = \log\left(\frac{2}{2}\right) = 0$$

Term Frequency-Inverse Document Frequency

$$tfidf("this", d_1) = 0.2 \times 0 = 0$$

$$tfidf("this", d_2) = 0.14 \times 0 = 0$$

47 / 133

Document 1

Term	Term Count
this	1
is	1
a	2
sample	1

Document 2

Term	Term Count
this	1
is	1
another	2
example	3

48 / 133

Document 1

Term	Term Count
this	1
is	1
a	2
sample	1

Term Frequency

$$tf(\text{" example "}, d_1) = \frac{0}{5} = 0$$

$$tf(\text{" example "}, d_2) = \frac{3}{7} \approx 0.429$$

Document 2

Term	Term Count
this	1
is	1
another	2
example	3

49 / 133

Document 1

Term	Term Count
this	1
is	1
a	2
sample	1

Term Frequency

$$tf(\text{" example "}, d_1) = \frac{0}{5} = 0$$

$$tf(\text{" example "}, d_2) = \frac{3}{7} \approx 0.429$$

Document 2

Term	Term Count
this	1
is	1
another	2
example	3

Inverse Document Frequency

$$idf(\text{" example "}, D) = \log \left(\frac{2}{1} \right) = 0.301$$

50 / 133

Document 1

Term	Term Count
this	1
is	1
a	2
sample	1

Term Frequency

$$tf("example", d_1) = \frac{0}{5} = 0$$

$$tf("example", d_2) = \frac{3}{7} \approx 0.429$$

Document 2

Term	Term Count
this	1
is	1
another	2
example	3

Inverse Document Frequency

$$idf("example", D) = \log\left(\frac{2}{1}\right) = 0.301$$

Term Frequency-Inverse Document Frequency

$$tfidf("example", d_1) = 0 \times 0.301 = 0$$

$$tfidf("example", d_2) = 0.429 \times 0.301 = 0.13$$

51 / 133

Lonely Words

- Many NLP systems treated words as isolated indices

52 / 133

Lonely Words

- Many NLP systems treated words as isolated indices
- Words are connected

53 / 133

Lonely Words

- Many NLP systems treated words as isolated indices
- Words are connected
- Similar words are used similarly

54 / 133

Types of Similarity

- Topical similarity based upon textual regions

55 / 133

Types of Similarity

- Topical similarity based upon textual regions
- Paradigmatic similarity based upon co-occurrence

56 / 133

Types of Similarity

- Topical similarity based upon textual regions
- Paradigmatic similarity based upon co-occurrence
- Syntagmatic similarity by examining the vectors

57 / 133

t-SNE Visualizations of Word Embeddings

<http://tinyurl.com/hpf24ox>

58 / 133

Advantages of VSM

- Simple models using linear algebra

https://en.wikipedia.org/wiki/Vector_space_model

59 / 133

Advantages of VSM

- Simple models using linear algebra
- Non-binary weights

https://en.wikipedia.org/wiki/Vector_space_model

60 / 133

Advantages of VSM

- Simple models using linear algebra
- Non-binary weights
- Continuous ranges of similarity between documents and queries

https://en.wikipedia.org/wiki/Vector_space_model

61 / 133

Advantages of VSM

- Simple models using linear algebra
- Non-binary weights
- Continuous ranges of similarity between documents and queries
- Ranking possibilities based on relevance

https://en.wikipedia.org/wiki/Vector_space_model

62 / 133

Advantages of VSM

- Simple models using linear algebra
- Non-binary weights
- Continuous ranges of similarity between documents and queries
- Ranking possibilities based on relevance
- Partial matching

https://en.wikipedia.org/wiki/Vector_space_model

63 / 133

Limitations of VSM

- Long documents are poorly represented

https://en.wikipedia.org/wiki/Vector_space_model

64 / 133

Limitations of VSM

- Long documents are poorly represented
- Search keywords must generally be exact

https://en.wikipedia.org/wiki/Vector_space_model

65 / 133

Limitations of VSM

- Long documents are poorly represented
- Search keywords must generally be exact
- Documents with similar contexts but different vocabularies yield missing results (synonymy)

https://en.wikipedia.org/wiki/Vector_space_model

66 / 133

Limitations of VSM

- Long documents are poorly represented
- Search keywords must generally be exact
- Documents with similar contexts but different vocabularies yield missing results (synonymy)
- Words with different meaning can yield bad results (polysemy)

https://en.wikipedia.org/wiki/Vector_space_model

67 / 133

Limitations of VSM

- Long documents are poorly represented
- Search keywords must generally be exact
- Documents with similar contexts but different vocabularies yield missing results (synonymy)
- Words with different meaning can yield bad results (polysemy)
- Term ordering is lost

https://en.wikipedia.org/wiki/Vector_space_model

68 / 133

Distributional Semantics

- **Distributional hypothesis** is based upon the semantic theory of language*.

*https://en.wikipedia.org/wiki/Distributional_semantics

69 / 133

Distributional Semantics

- **Distributional hypothesis** is based upon the semantic theory of language*.
- "i.e. words that are used and occur in the same contexts tend to purport similar meanings" (Harris, Z.)

*https://en.wikipedia.org/wiki/Distributional_semantics

70 / 133

Distributional Semantics

- **Distributional hypothesis** is based upon the semantic theory of language*.
- "i.e. words that are used and occur in the same contexts tend to purport similar meanings" (Harris, Z.)
- "a word is characterized by the company it keeps" (Firth, John R.)

*https://en.wikipedia.org/wiki/Distributional_semantics

71 / 133

Distributional Semantics

- **Distributional hypothesis** is based upon the semantic theory of language*.
- "i.e. words that are used and occur in the same contexts tend to purport similar meanings" (Harris, Z.)
- "a word is characterized by the company it keeps" (Firth, John R.)
- Uses linear algebra for computational and representational purposes

*https://en.wikipedia.org/wiki/Distributional_semantics

72 / 133

Distributional Semantics

- **Distributional hypothesis** is based upon the semantic theory of language*.
- "i.e. words that are used and occur in the same contexts tend to purport similar meanings" (Harris, Z.)
- "a word is characterized by the company it keeps" (Firth, John R.)
- Uses linear algebra for computational and representational purposes
- High-dimensional vectors are built from a corpus

*https://en.wikipedia.org/wiki/Distributional_semantics

73 / 133

Latent Semantic Analysis

74 / 133

Latent Semantic Analysis

- Construct a term-document matrix

75 / 133

Latent Semantic Analysis

- Construct a term-document matrix
- Find a low-rank approximation

76 / 133

Latent Semantic Analysis

- Construct a term-document matrix
- Find a low-rank approximation
 - Matrix is too big

77 / 133

Latent Semantic Analysis

- Construct a term-document matrix
- Find a low-rank approximation
 - Matrix is too big
 - Matrix is too noisy

78 / 133

Latent Semantic Analysis

- Construct a term-document matrix
- Find a low-rank approximation
 - Matrix is too big
 - Matrix is too noisy
 - Matrix is too sparse

79 / 133

LSA Uses

80 / 133

LSA Uses

- Document classification

81 / 133

LSA Uses

- Document classification
- Cross language retrieval

82 / 133

LSA Uses

- Document classification
- Cross language retrieval
- Solutions to synonymy and polysemy

83 / 133

LSA Uses

- Document classification
- Cross language retrieval
- Solutions to synonymy and polysemy
- Convert queries into the low-dimensional matrix to find documents

84 / 133

Word2Vec

85 / 133

http://mccormickml.com/assets/word2vec/Alex_Minnaar_Word2Vec_Tutorial_Part_I_of-Words_Model.pdf

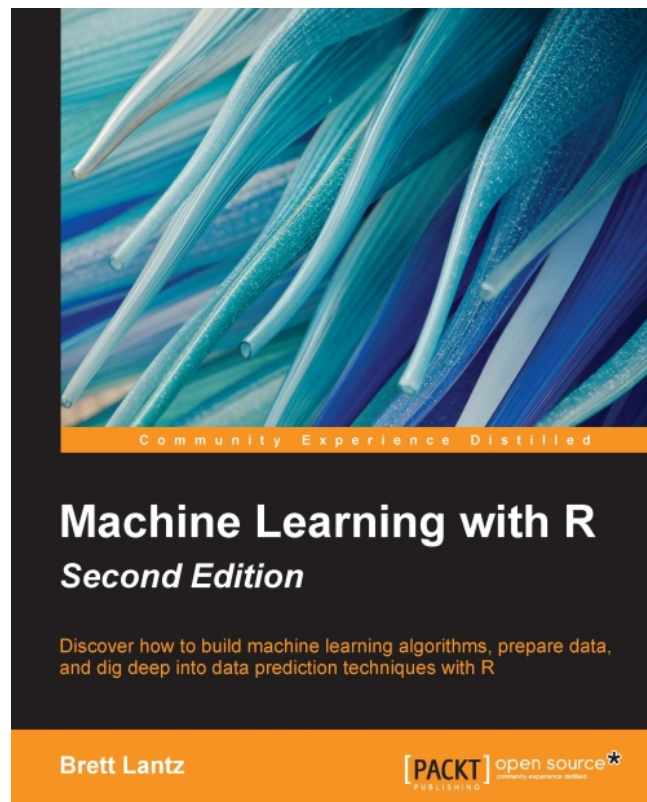
86 / 133

<https://deeplearning4j.org/word2vec.html>

87 / 133

NaiveBayes

88 / 133



89 / 133

<https://github.com/dataspelunking/MLwR/>

90 / 133

```
# read the sms data into the sms data frame
> sms_raw <- read.csv("sms_spam.csv", stringsAsFactors = FALSE)
```

91 / 133

```
# read the sms data into the sms data frame
> sms_raw <- read.csv("sms_spam.csv", stringsAsFactors = FALSE)
```

```
# examine the structure of the sms data
> str(sms_raw)
'data.frame': 5559 obs. of 2 variables:
 $ type: chr "ham" "ham" "ham" "spam" ...
 $ text: chr "Hope you are having a good week. Just checking in"
 "K..give back my thanks."
 "Am also doing in cbe only. But have to pay."
 "complimentary 4 STAR Ibiza Holiday or £10,000 cash needs your URGENT
 collection. 09066364349 NOW from Landline not to lose out!"
 |__truncated__ ...
```

92 / 133

```
# read the sms data into the sms data frame
> sms_raw <- read.csv("sms_spam.csv", stringsAsFactors = FALSE)
```

```
# examine the structure of the sms data
> str(sms_raw)
'data.frame': 5559 obs. of 2 variables:
 $ type: chr "ham" "ham" "ham" "spam" ...
 $ text: chr "Hope you are having a good week. Just checking in"
 "K..give back my thanks."
 "Am also doing in cbe only. But have to pay."
 "complimentary 4 STAR Ibiza Holiday or £10,000 cash needs your URGENT
 collection. 09066364349 NOW from Landline not to lose out!"
 |__truncated__ ...
```

```
# convert spam/ham to factor.
> sms_raw$type <- factor(sms_raw$type)
```

93 / 133

```
# read the sms data into the sms data frame
> sms_raw <- read.csv("sms_spam.csv", stringsAsFactors = FALSE)
```

```
# examine the structure of the sms data
> str(sms_raw)
'data.frame': 5559 obs. of 2 variables:
 $ type: chr "ham" "ham" "ham" "spam" ...
 $ text: chr "Hope you are having a good week. Just checking in"
 "K..give back my thanks."
 "Am also doing in cbe only. But have to pay."
 "complimentary 4 STAR Ibiza Holiday or £10,000 cash needs your URGENT
 collection. 09066364349 NOW from Landline not to lose out!"
 |__truncated__ ...
```

```
# convert spam/ham to factor.
> sms_raw$type <- factor(sms_raw$type)
```

```
# examine the type variable more carefully
> str(sms_raw$type)
Factor w/ 2 levels "ham","spam": 1 1 1 2 2 1 1 1 2 1 ...
```

94 / 133

```
# read the sms data into the sms data frame
> sms_raw <- read.csv("sms_spam.csv", stringsAsFactors = FALSE)
```

```
# examine the structure of the sms data
> str(sms_raw)
'data.frame': 5559 obs. of 2 variables:
 $ type: chr "ham" "ham" "ham" "spam" ...
 $ text: chr "Hope you are having a good week. Just checking in"
 "K..give back my thanks."
 "Am also doing in cbe only. But have to pay."
 "complimentary 4 STAR Ibiza Holiday or £10,000 cash needs your URGENT
 collection. 09066364349 NOW from Landline not to lose out!"
 |__truncated__ ...
```

```
# convert spam/ham to factor.
> sms_raw$type <- factor(sms_raw$type)
```

```
# examine the type variable more carefully
> str(sms_raw$type)
Factor w/ 2 levels "ham","spam": 1 1 1 2 2 1 1 1 2 1 ...
```

```
> table(sms_raw$type)
```

```
ham spam
4812 747
```

95 / 133

```
# build a corpus using the text mining (tm) package
> library(tm)
Loading required package: NLP
```

96 / 133


```
# build a corpus using the text mining (tm) package
> library(tm)
Loading required package: NLP
```

```
> sms_corpus <- VCorpus(VectorSource(sms_raw$text))
```

97 / 133

```
# build a corpus using the text mining (tm) package
> library(tm)
Loading required package: NLP
```

```
> sms_corpus <- VCorpus(VectorSource(sms_raw$text))
```

```
# examine the sms corpus
> print(sms_corpus)
<<VCorpus>>
Metadata: corpus specific: 0, document level (indexed): 0
Content: documents: 5559
```

98 / 133

```
# build a corpus using the text mining (tm) package
> library(tm)
Loading required package: NLP
```

```
> sms_corpus <- VCorpus(VectorSource(sms_raw$text))
```

```
# examine the sms corpus
> print(sms_corpus)
<<VCorpus>>
Metadata: corpus specific: 0, document level (indexed): 0
Content: documents: 5559
```

```
> inspect(sms_corpus[1:2])
<<VCorpus>>
Metadata: corpus specific: 0, document level (indexed): 0
Content: documents: 2
```

```
[[1]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 49
```

```
[[2]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 23
```

99 / 133

```
> as.character(sms_corpus[[1]])
[1] "Hope you are having a good week. Just checking in"
```

100 / 133

```
> as.character(sms_corpus[[1]])  
[1] "Hope you are having a good week. Just checking in"
```

```
> lapply(sms_corpus[1:2], as.character)  
$`1`  
[1] "Hope you are having a good week. Just checking in"  
  
$`2`  
[1] "K..give back my thanks."
```

101 / 133

```
> as.character(sms_corpus[[1]])  
[1] "Hope you are having a good week. Just checking in"
```

```
> lapply(sms_corpus[1:2], as.character)  
$`1`  
[1] "Hope you are having a good week. Just checking in"  
  
$`2`  
[1] "K..give back my thanks."
```

```
> sms_corpus_clean <- tm_map(sms_corpus, content_transformer(tolower))
```

102 / 133

```
> as.character(sms_corpus[[1]])  
[1] "Hope you are having a good week. Just checking in"
```

```
> lapply(sms_corpus[1:2], as.character)  
$`1`  
[1] "Hope you are having a good week. Just checking in"  
  
$`2`  
[1] "K..give back my thanks."
```

```
> sms_corpus_clean <- tm_map(sms_corpus, content_transformer(tolower))
```

```
# show the difference between sms_corpus and corpus_clean  
> as.character(sms_corpus[[1]])  
[1] "Hope you are having a good week. Just checking in"  
> as.character(sms_corpus_clean[[1]])  
[1] "hope you are having a good week. just checking in"
```

103 / 133

```
> as.character(sms_corpus[[1]])  
[1] "Hope you are having a good week. Just checking in"
```

```
> lapply(sms_corpus[1:2], as.character)  
$`1`  
[1] "Hope you are having a good week. Just checking in"  
  
$`2`  
[1] "K..give back my thanks."
```

```
> sms_corpus_clean <- tm_map(sms_corpus, content_transformer(tolower))
```

```
# show the difference between sms_corpus and corpus_clean  
> as.character(sms_corpus[[1]])  
[1] "Hope you are having a good week. Just checking in"  
> as.character(sms_corpus_clean[[1]])  
[1] "hope you are having a good week. just checking in"
```

```
# remove numbers  
> sms_corpus_clean <- tm_map(sms_corpus_clean, removeNumbers)  
# remove stop words  
> sms_corpus_clean <- tm_map(sms_corpus_clean, removeWords, stopwords())  
# remove punctuation  
> sms_corpus_clean <- tm_map(sms_corpus_clean, removePunctuation)
```

104 / 133

```
# tip: create a custom function to replace (rather than remove) punctuation
> removePunctuation("hello...world")
[1] "helloworld"
```

105 / 133

```
# tip: create a custom function to replace (rather than remove) punctuation
> removePunctuation("hello...world")
[1] "helloworld"
```

```
> replacePunctuation <- function(x) { gsub("[[:punct:]]+", " ", x) }
> replacePunctuation("hello...world")
[1] "hello world"
```

106 / 133

```
# tip: create a custom function to replace (rather than remove) punctuation
> removePunctuation("hello...world")
[1] "helloworld"
```

```
> replacePunctuation <- function(x) { gsub("[:punct:]]+", " ", x) }
> replacePunctuation("hello...world")
[1] "hello world"
```

```
# illustration of word stemming
> library(SnowballC)
> wordStem(c("learn", "learned", "learning", "learns"))
[1] "learn" "learn" "learn" "learn"
```

107 / 133

```
# tip: create a custom function to replace (rather than remove) punctuation
> removePunctuation("hello...world")
[1] "helloworld"
```

```
> replacePunctuation <- function(x) { gsub("[:punct:]]+", " ", x) }
> replacePunctuation("hello...world")
[1] "hello world"
```

```
# illustration of word stemming
> library(SnowballC)
> wordStem(c("learn", "learned", "learning", "learns"))
[1] "learn" "learn" "learn" "learn"
```

```
> sms_corpus_clean <- tm_map(sms_corpus_clean, stemDocument)
# eliminate unneeded whitespace
> sms_corpus_clean <- tm_map(sms_corpus_clean, stripWhitespace)
```

108 / 133

```
# examine the final clean corpus
> lapply(sms_corpus[1:3], as.character)
$`1`
[1] "Hope you are having a good week. Just checking in"

$`2`
[1] "K..give back my thanks."

$`3`
[1] "Am also doing in cbe only. But have to pay."
```

109 / 133

```
# examine the final clean corpus
> lapply(sms_corpus[1:3], as.character)
$`1`
[1] "Hope you are having a good week. Just checking in"

$`2`
[1] "K..give back my thanks."

$`3`
[1] "Am also doing in cbe only. But have to pay."
```

```
> lapply(sms_corpus_clean[1:3], as.character)
$`1`
[1] "hope you are have a good week. just check in"

$`2`
[1] "k..give back my thanks."

$`3`
[1] "am also do in cbe only. but have to pay."
```

110 / 133

```

# create a document-term sparse matrix
> sms_dtm <- DocumentTermMatrix(sms_corpus_clean)
> sms_dtm
<<DocumentTermMatrix (documents: 5559, terms: 10856)>>
Non-/sparse entries: 59144/60289360
Sparsity           : 100%
Maximal term length: 48
Weighting           : term frequency (tf)

```

111 / 133

```

# create a document-term sparse matrix
> sms_dtm <- DocumentTermMatrix(sms_corpus_clean)
> sms_dtm
<<DocumentTermMatrix (documents: 5559, terms: 10856)>>
Non-/sparse entries: 59144/60289360
Sparsity           : 100%
Maximal term length: 48
Weighting           : term frequency (tf)

```

```

# creating training and test datasets
> sms_dtm_train <- sms_dtm[1:4169, ]
> sms_dtm_test  <- sms_dtm[4170:5559, ]

```

112 / 133


```
# create a document-term sparse matrix
> sms_dtm <- DocumentTermMatrix(sms_corpus_clean)
> sms_dtm
<<DocumentTermMatrix (documents: 5559, terms: 10856)>>
Non-/sparse entries: 59144/60289360
Sparsity           : 100%
Maximal term length: 48
Weighting          : term frequency (tf)
```

```
# creating training and test datasets
> sms_dtm_train <- sms_dtm[1:4169, ]
> sms_dtm_test  <- sms_dtm[4170:5559, ]
```

```
# also save the labels
sms_train_labels <- sms_raw[1:4169, ]$type
sms_test_labels  <- sms_raw[4170:5559, ]$type
```

113 / 133

```
# create a document-term sparse matrix
> sms_dtm <- DocumentTermMatrix(sms_corpus_clean)
> sms_dtm
<<DocumentTermMatrix (documents: 5559, terms: 10856)>>
Non-/sparse entries: 59144/60289360
Sparsity           : 100%
Maximal term length: 48
Weighting          : term frequency (tf)
```

```
# creating training and test datasets
> sms_dtm_train <- sms_dtm[1:4169, ]
> sms_dtm_test  <- sms_dtm[4170:5559, ]
```

```
# also save the labels
sms_train_labels <- sms_raw[1:4169, ]$type
sms_test_labels  <- sms_raw[4170:5559, ]$type
```

```
# check that the proportion of spam is similar
> prop.table(table(sms_train_labels))
sms_train_labels
   ham   spam
0.8647158 0.1352842

> prop.table(table(sms_test_labels))
sms_test_labels
   ham   spam
0.8683453 0.1316547
```

114 / 133

```
> # word cloud visualization
> library(wordcloud)
Loading required package: RColorBrewer
```

115 / 133

```
> # word cloud visualization
> library(wordcloud)
Loading required package: RColorBrewer
```

```
> wordcloud(sms_corpus_clean, min.freq = 50, random.order = FALSE)
```

116 / 133

```
> # word cloud visualization
> library(wordcloud)
Loading required package: RColorBrewer
```

```
> wordcloud(sms_corpus_clean, min.freq = 50, random.order = FALSE)
```



117 / 133

```
# subset the training data into spam and ham groups
> spam <- subset(sms_raw, type == "spam")
> ham <- subset(sms_raw, type == "ham")
```

118 / 133

```
# subset the training data into spam and ham groups
> spam <- subset(sms_raw, type == "spam")
> ham <- subset(sms_raw, type == "ham")
```

```
> wordcloud(spam$text, max.words = 40, scale = c(3, 0.5))
```

119 / 133

```
# subset the training data into spam and ham groups
> spam <- subset(sms_raw, type == "spam")
> ham <- subset(sms_raw, type == "ham")
```

```
> wordcloud(spam$text, max.words = 40, scale = c(3, 0.5))
```

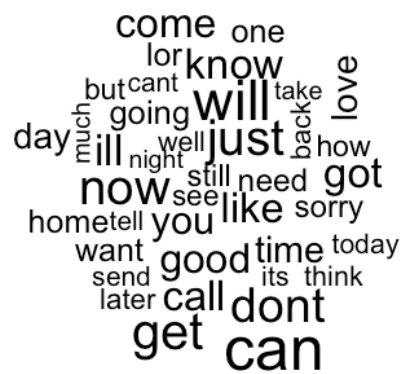


120 / 133

```
> wordcloud(ham$text, max.words = 40, scale = c(3, 0.5))
```

121 / 133

```
> wordcloud(ham$text, max.words = 40, scale = c(3, 0.5))
```



122 / 133

```
> sms_dtm_freq_train <- removeSparseTerms(sms_dtm_train, 0.999)
```

123 / 133

```
> sms_dtm_freq_train <- removeSparseTerms(sms_dtm_train, 0.999)
```

```
> sms_dtm_freq_train
<<DocumentTermMatrix (documents: 4169, terms: 1306)>>
Non-/sparse entries: 33250/5411464
Sparsity           : 99%
Maximal term length: 24
Weighting          : term frequency (tf)
```

124 / 133

```
> sms_dtm_freq_train <- removeSparseTerms(sms_dtm_train, 0.999)
```

```
> sms_dtm_freq_train
<<DocumentTermMatrix (documents: 4169, terms: 1306)>>
Non-/sparse entries: 33250/5411464
Sparsity           : 99%
Maximal term length: 24
Weighting          : term frequency (tf)
```

```
# save frequently-appearing terms to a character vector
> sms_freq_words <- findFreqTerms(sms_dtm_train, 5)
```

125 / 133

```
> sms_dtm_freq_train <- removeSparseTerms(sms_dtm_train, 0.999)
```

```
> sms_dtm_freq_train
<<DocumentTermMatrix (documents: 4169, terms: 1306)>>
Non-/sparse entries: 33250/5411464
Sparsity           : 99%
Maximal term length: 24
Weighting          : term frequency (tf)
```

```
# save frequently-appearing terms to a character vector
> sms_freq_words <- findFreqTerms(sms_dtm_train, 5)
```

```
> str(sms_freq_words)
chr [1:1332] ":-(" ":-)" "!!'.' "..." "....." ...
```

126 / 133

```
> sms_dtm_freq_train <- removeSparseTerms(sms_dtm_train, 0.999)
```

```
> sms_dtm_freq_train
<<DocumentTermMatrix (documents: 4169, terms: 1306)>>
Non-/sparse entries: 33250/5411464
Sparsity           : 99%
Maximal term length: 24
Weighting          : term frequency (tf)
```

```
# save frequently-appearing terms to a character vector
> sms_freq_words <- findFreqTerms(sms_dtm_train, 5)
```

```
> str(sms_freq_words)
chr [1:1332] ":-(" ":-)" "!!'." "... " "....." ...
```

```
# create DTMs with only the frequent terms
> sms_dtm_freq_train <- sms_dtm_train[, sms_freq_words]
> sms_dtm_freq_test <- sms_dtm_test[, sms_freq_words]
```

127 / 133

```
> # convert counts to a factor
> convert_counts <- function(x) {
+   x <- ifelse(x > 0, "Yes", "No")
+ }
```

128 / 133


```
> # convert counts to a factor
> convert_counts <- function(x) {
+   x <- ifelse(x > 0, "Yes", "No")
+ }
```

```
# apply() convert_counts() to columns of train/test data
> sms_train <- apply(sms_dtm_freq_train, MARGIN = 2, convert_counts)
> sms_test  <- apply(sms_dtm_freq_test, MARGIN = 2, convert_counts)
```

129 / 133

```
> # convert counts to a factor
> convert_counts <- function(x) {
+   x <- ifelse(x > 0, "Yes", "No")
+ }
```

```
# apply() convert_counts() to columns of train/test data
> sms_train <- apply(sms_dtm_freq_train, MARGIN = 2, convert_counts)
> sms_test  <- apply(sms_dtm_freq_test, MARGIN = 2, convert_counts)
```

```
## Step 3: Training a model on the data ----
> library(e1071)
> sms_classifier <- naiveBayes(sms_train, sms_train_labels)
```

130 / 133

```
> # convert counts to a factor
> convert_counts <- function(x) {
+   x <- ifelse(x > 0, "Yes", "No")
+ }
```

```
# apply() convert_counts() to columns of train/test data
> sms_train <- apply(sms_dtm_freq_train, MARGIN = 2, convert_counts)
> sms_test <- apply(sms_dtm_freq_test, MARGIN = 2, convert_counts)
```

```
## Step 3: Training a model on the data ----
> library(e1071)
> sms_classifier <- naiveBayes(sms_train, sms_train_labels)
```

```
## Step 4: Evaluating model performance ----
> sms_test_pred <- predict(sms_classifier, sms_test)
```

131 / 133

```
> library(gmodels)
> CrossTable(sms_test_pred, sms_test_labels,
+   prop.chisq = FALSE, prop.t = FALSE, prop.r = FALSE,
+   dnn = c('predicted', 'actual'))
```

```
Cell Contents
|-----|
|              N |
| N / Col Total |
|-----|
```

Total Observations in Table: 1390

predicted	actual		Row Total
	ham	spam	
ham	1201 0.995	31 0.169	1232
spam	6 0.005	152 0.831	158
Column Total	1207 0.868	183 0.132	1390

132 / 133

A large golden Buddha statue is shown from the chest up, positioned on the left side of the frame. The statue has a serene expression and is set against a background of a bright blue sky filled with scattered white clouds. The lighting suggests a sunny day.

Questions

✉ brian@bosatsu.net

🐦 [@bsletten](https://twitter.com/bsletten)

🔄 [bsletten](https://www.youtube.com/channel/UCbsletten)