

Text Mining

February 2018

Outline

- Text Mining - Introduction
- Text Mining Pipeline (Preprocessing)
 - Tokenisation
 - Stemming
 - Stopword Removal
 - Sentence Segmentation
- Part-Of-Speech (POS) tagging

Text Mining - Definition and Challenges

- Text mining
 - process of extracting interesting and non-trivial patterns or knowledge from unstructured text documents [Tan et al., 1999].
 - *a.k.a* text data mining [Hearst, 1997],
 - knowledge discovery from textual databases [Feldman and Dagan, 1995]
 - text analytics - application to solve business problems

Text Mining - Definition and Challenges

- Text mining
 - process of extracting interesting and non-trivial patterns or knowledge from unstructured text documents [Tan et al., 1999].
 - *a.k.a* text data mining [Hearst, 1997],
 - knowledge discovery from textual databases [Feldman and Dagan, 1995]
 - text analytics - application to solve business problems
- Challenges
 - deriving semantics from the content
 - combine information from multi-lingual texts
 - integrate domain knowledge

Text Mining - Framework

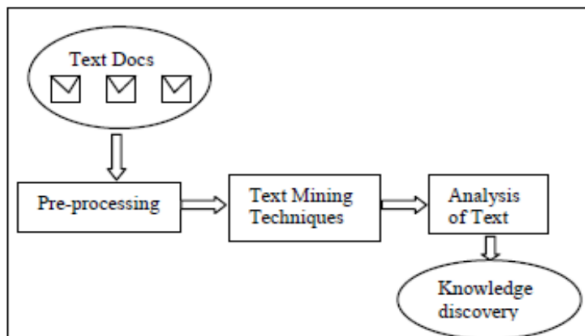


Figure1: Framework of Text Mining

Text Mining - Preprocessing

- Tokenisation
- Stemming
- Stopword Removal
- Sentence Segmentation

Tokenisation

- Process of splitting text into words
- What is a word?
 - string of contiguous alphanumeric characters with space on either side; may include hyphens and apostrophes, but no other punctuation marks [Kučera and Francis, 1967].
- Useful clue - space or tab (English)

Tokenisation - Problems

- Periods
 - usually helps if we remove them
 - but useful to retain in certain cases such as \$22.50; Ed.,
- Hyphenation
 - useful to retain in some cases e.g., state-of-the-art
 - better to remove in other cases e.g., gold-import ban, 50-year-old

Tokenisation - Problems

- Periods
 - usually helps if we remove them
 - but useful to retain in certain cases such as \$22.50; Ed.,
- Hyphenation
 - useful to retain in some cases e.g., state-of-the-art
 - better to remove in other cases e.g., gold-import ban, 50-year-old
- Single apostrophes
 - useful to remove them e.g., *is'nt*, *didn't*
- Space may not be a useful clue all the time
 - sometimes we want to use words separated by space as 'single' word
 - For example:
 - San Francisco
 - University of Liverpool
 - Danushka Bollegala

Regular Expressions for Tokenisation

- Regular Expressions Cheatsheet

REGEX	NOTE	EXAMPLE	EXPLANATION
<code>\s</code>	white space	<code>\d\s\d</code>	digit space digit
<code>\S</code>	not white space	<code>\d\S\d</code>	digit non-whitespace digit
<code>\d</code>	digit	<code>\d\d\d-\d\d-\d\d\d\d</code>	SSN
<code>\D</code>	not digit	<code>\D\D\D</code>	three non-digits
<code>\w</code>	word character (letter, number, or _)	<code>\w\w\w</code>	three word chars
<code>\W</code>	not a word character	<code>\W\W\W</code>	three non-word chars
<code>[...]</code>	any included character	<code>[a-z0-9#]</code>	any char that is a thru z, 0 thru 9, or #
<code>[^...]</code>	no included character	<code>[^xyz]</code>	any char but x, y, or z
<code>*</code>	zero or more	<code>\w*</code>	zero or more words chars
<code>+</code>	one or more	<code>\d+</code>	integer
<code>?</code>	zero or one	<code>\d\d\d-?\d\d-?\d\d\d\d</code>	SSN with dashes being optional
<code> </code>	or	<code>\w \d</code>	word or digit character

Regular Expressions for Tokenisation

```

1 raw = """'When I'M a Duchess,' she said to herself, (not in a very hopeful tone
2         though), 'I won't have any pepper in my kitchen AT ALL. Soup does very
3         well without--Maybe it's always pepper that makes people hot-tempered,'..."""
4
5 import re
6 print re.split(r' ', raw)
7 ['When', 'I'M', 'a', 'Duchess,', 'she', 'said', 'to', 'herself,', '(not', 'in', 'a',
8  'very', 'hopeful', 'tone\n\t\t', 'though)', 'I', 'won't', 'have', 'any',
9  'pepper', 'in', 'my', 'kitchen', 'AT', 'ALL.', 'Soup', 'does', 'very\n\t\t',
10 'well', 'without--Maybe', 'it's', 'always', 'pepper', 'that', 'makes', 'people',
11 'hot-tempered,...']
12
13 print re.split(r'[ \t\n]+', raw)
14 ['When', 'I'M', 'a', 'Duchess,', 'she', 'said', 'to', 'herself,', '(not', 'in', 'a',
15  'very', 'hopeful', 'tone', 'though)', 'I', 'won't', 'have', 'any', 'pepper', 'in',
16  'my', 'kitchen', 'AT', 'ALL.', 'Soup', 'does', 'very', 'well', 'without--Maybe',
17  'it's', 'always', 'pepper', 'that', 'makes', 'people', 'hot-tempered,...']
18
19 print re.findall(r"\w+ (?::-]\w+)*|[-.()+\S\w*", raw)
20 ['', 'When', 'I', 'M', 'a', 'Duchess', 'she', 'said', 'to',
21  'herself', '(', 'not', 'in', 'a', 'very', 'hopeful', 'tone', 'though',
22  ')', 'I', 'won', 't', 'have', 'any', 'pepper', 'in', 'my',
23  'kitchen', 'AT', 'ALL', 'Soup', 'does', 'very', 'well', 'without', '--',
24  'Maybe', 'it', 's', 'always', 'pepper', 'that', 'makes', 'people',
25  'hot', '-', 'tempered', '...', '']

```

Stanford Parser for Tokenisation

```
1 raw = """'When I'M a Duchess,' she said to herself, (not in a very hopeful tone
2         though), 'I won't have any pepper in my kitchen AT ALL. Soup does very
3         well without--Maybe it's always pepper that makes people hot-tempered,'...""
4
5 path_to_parser_jar = 'lib/stanford-parser.jar'
6 path_to_models_jar = 'lib/stanford-parser-3.5.1-models.jar'
7
8 # POS Tagger
9 from nltk.tokenize.stanford import StanfordTokenizer
10 tokenizer = StanfordTokenizer(path_to_parser_jar)
11
12 tokenized_text = tokenizer.tokenize(raw)
13 print tokenized_text
14
15 [u'\'', u'When', u'I', u'M', u'a', u'Duchess', u',', u'', u'she', u'said', u'to',
16 u'herself', u',', u'-LRB-', u'not', u'in', u'a', u'very', u'hopeful', u'tone', u'though',
17 u'-RRB-', u',', u'', u'I', u'wo', u'n't', u'have', u'any', u'pepper', u'in', u'my',
18 u'kitchen', u'AT', u'ALL', u'.', u'Soup', u'does', u'very', u'well', u'without', u'--',
19 u'Maybe', u'it', u's', u'always', u'pepper', u'that', u'makes', u'people', u'hot-tempered',
20 u',', u'', u'', u'...']]
```

- Stanford CoreNLP - <https://stanfordnlp.github.io/CoreNLP/>

Tokenisation

- Tokenisation turns out to be more difficult than one expects
- No single solution works well
- Decide what counts as a token depending on the application domain

Stemming

- Removal of inflectional ending from words (strip off any affixes)
 - connections, connecting, connect, connected → connect
- Problems
 - Can conflate semantically different words
 - *Gallery* and *gall* may both be stemmed to *gall*
- Lemmatization: a further step to ensure that the resulting form is a word present in a dictionary

Regular Expressions for Stemming

```
1 import re
2 print re.findall(r'^(.*)(ing|ly|ed|ions|ies|ive|es|s|ment)$', 'processing')
3 [('process', 'ing')]
4
5 import re
6 print re.findall(r'^(.*)(ing|ly|ed|ions|ies|ive|es|s|ment)$', 'processes')
7 [('processe', 's')]
8
```

- note that the star operator is “greedy”
- the .* part of expression tries to consume as much as the input as possible
- for non-greedy version of the star operator = *?

```
9 import re
10 print re.findall(r'^(.?*)(ing|ly|ed|ions|ies|ive|es|s|ment)$', 'processes')
11 [('process', 'es')]
12
13
```

Regular Expressions for Stemming

```
78 import nltk, re
79
80 def stem(word):
81     regexp = r'^(.*?)(ing|ly|ed|ions|ies|ive|es|s|ment)?$'
82     stem, suffix = re.findall(regexp, word)[0]
83     return stem
84
85 raw = """DENNIS: Listen, strange women lying in ponds distributing swords
86         is no basis for a system of government. Supreme executive power derives from
87         a mandate from the masses, not from some farcical aquatic ceremeony."""
88
89 tokens = nltk.word_tokenize(raw)
90 print [stem(t) for t in tokens]
91
92 ['DENNIS', ':', 'Listen', ',', 'strange', 'women', 'ly', 'in', 'pond', 'distribut', 'sword',
93  'i', 'no', 'basi', 'for', 'a', 'system', 'of', 'govern', 'Supreme', 'execut', 'power',
94  'deriv', 'from', 'a', 'mandate', 'from', 'the', 'mass', 'not', 'from', 'some',
95  'farcical', 'aquatic', 'ceremeony', '.']
96
97
```


Regular Expressions for Stemming

```
78 import nltk, re
79
80 def stem(word):
81     regexp = r'^(.?)(ing|ly|ed|ions|ies|ive|es|s|ment)?$'
82     stem, suffix = re.findall(regexp, word)[0]
83     return stem
84
85 raw = """DENNIS: Listen, strange women lying in ponds distributing swords
86         is no basis for a system of government. Supreme executive power derives from
87         a mandate from the masses, not from some farcical aquatic ceremeony."""
88
89 tokens = nltk.word_tokenize(raw)
90 print [stem(t) for t in tokens]
91
92 ['DENNIS', ':', 'Listen', ',', 'strange', 'women', 'ly', 'in', 'pond', 'distribut', 'sword',
93  'i', 'no', 'basi', 'for', 'a', 'system', 'of', 'govern', 'Supreme', 'execut', 'power',
94  'deriv', 'from', 'a', 'mandate', 'from', 'the', 'mass', 'not', 'from', 'some',
95  'farcical', 'aquatic', 'ceremeony', '.']
96
97
```

• Problems

- RE removes 's' from 'ponds', but also from 'is' and 'basis'
- produces some non-words like 'distribut', 'deriv'

NLTK Stemmers

- NLTK (<https://www.nltk.org/>) provides off-the-shelf stemmers
- Porter Stemmer (www.nltk.org/modules/nltk/stem/porter.html)
- Lancaster Stemmer (www.nltk.org/modules/nltk/stem/lancaster.html)
- Stemmers have their own rules for stripping affixes

```
1 import nltk, re
2
3 raw = """DENNIS: Listen, strange women lying in ponds distributing swords
4         is no basis for a system of government. Supreme executive power derives from
5         a mandate from the masses, not from some farcical aquatic ceremeony."""
6
7 porter = nltk.PorterStemmer()
8 lancaster = nltk.LancasterStemmer()
9 tokens = nltk.word_tokenize(raw)
10
11 print [porter.stem(t) for t in tokens]
12 [u'denni', ':', 'listen', ',', u'strang', 'women', u'lie', 'in', u'pond', u'distribut',
13  u'sword', 'is', 'no', u'basi', 'for', 'a', 'system', 'of', u'govern', '.', u'suprem',
14  u'execut', 'power', u'deriv', 'from', 'a', u'mandat', 'from', 'the', u'mass', ',', 'not',
15  'from', 'some', u'farcic', u'aquat', u'ceremeoni', '.']
16
17 print [lancaster.stem(t) for t in tokens]
18 ['den', ':', 'list', ',', 'strange', 'wom', 'lying', 'in', 'pond', 'distribut', 'sword',
19  'is', 'no', 'bas', 'for', 'a', 'system', 'of', 'govern', '.', 'suprem', 'execut', 'pow',
20  'der', 'from', 'a', 'mand', 'from', 'the', 'mass', ',', 'not', 'from', 'som', 'farc',
21  'aqu', 'ceremeony', '.']
```

Is stemming useful?

- Provides some improvement for IR performance (especially for smaller documents).
- Very useful for some queries, but on an average does not help much.
- Since improvement is very minimal, often IR engines does not use stemming.

Stemming vs. Lemmatization

- Stemming

- crude heuristic process that chops off the ends of words
- often includes the removal of derivational affixes
- Example: 'Gallery' , 'gall' → 'gall'

- Lemmatization

- doing things more properly, using a vocabulary and morphological analysis of words
- aims to remove inflectional endings only
- return the base or dictionary form of a word, known as lemma
- Example: 'better' → 'good'

Stopword Removal

- Removal of high frequency words
- Most common words such as articles, prepositions, and pronouns etc. does not help in identifying meaning

a	an	and	are	as	at	be	by	for	from
has	he	in	is	it	its	of	on	that	the
to	was	were	will	with					

Figure: Stop list of 25 semantically non-selective words common in Reuters-RCV1

- Uses
 - significantly reduce feature space
 - speed up learning

Methods for stopwords removal

- Classic Method
 - removing stop-words using pre-compiled lists
- Zipf's law (Z-methods)
 - frequency of a word is inversely proportional to its rank in the frequency table
 - remove most frequent words
- Mutual Information Method
 - supervised method that computes mutual information between a given term and a document class
 - low mutual information suggests low discrimination power of the term and hence should be removed

Sentence Segmentation

- Divide text into sentences
- Involves identifying **sentence boundaries** between words in different sentences
- *a.k.a* sentence boundary detection, sentence boundary disambiguation, sentence boundary recognition
- Useful and necessary for various NLP tasks such as
 - sentiment analysis
 - relation extraction
 - question answering systems
 - knowledge extraction

Sentence boundary detection algorithms

- Heuristic methods
- Statistical classification trees [Riley, 1989]
 - probability of a word occurring before or after a boundary, case and length of words
- Neural Networks [Palmer and Hearst, 1997]
 - POS distribution of preceding and following words
- Maximum entropy model [Mikheev 1998]

NLTK Tools for Sentence Segmentation:

Punkt Sentence Segmenter

```
import nltk, re

sent_tokenizer = nltk.data.load('tokenizers/punkt/english.pickle')
text = nltk.corpus.gutenberg.raw('chesterton-thursday.txt')
sents = sent_tokenizer.tokenize(text)
print sents[100:110]

[u'inquired Gregory sarcastically.',
 u'"I tell you," went on Syme with passion, "that every time a \
  train\ncomes in I feel that it has broken past batteries of besiegers, \
  and\nthat man has won a battle against chaos.',
 u'You say contemptuously\nthat when one has left Sloane Square one must come \
  to Victoria.', u'I\nsay that one might do a thousand things instead, and that \
  whenever\nI really come there I have the sense of hairbreadth escape.',
 u'And\nwhen I hear the guard shout out the word 'Victoria,' it is not \
  an\nunmeaning word.",
 u'It is to me the cry of a herald announcing\nconquest.',
 u'It is to me indeed \'Victoria\'; it is the victory of\nAdam."',
 u'Gregory wagged his heavy, red head with a slow and sad smile.',]
```

Part-of-Speech Tagging (POS)

- Task of tagging POS tags (Nouns, Verbs, Adjectives, Adverbs, ...) for words
- POS tags provide lot of information about a word
 - knowing whether a word is **noun** or **verb** gives information about neighbouring words
 - nouns are preceded by determiners; adjectives and verbs by nouns
 - provide useful features for **named entity recognition**

Part-of-Speech Tagging (POS)

- Task of tagging POS tags (Nouns, Verbs, Adjectives, Adverbs, ...) for words
- POS tags provide lot of information about a word
 - knowing whether a word is **noun** or **verb** gives information about neighbouring words
 - nouns are preceded by determiners; adjectives and verbs by nouns
 - provide useful features for **named entity recognition**
- Given a word, we assume it can belong to only of the POS tags.
- POS Tagging problem
 - Given a sentence $S = w_1 w_2 \dots w_n$ consisting of n words, determine the corresponding tag sequence $P = P_1 P_2 \dots P_n$

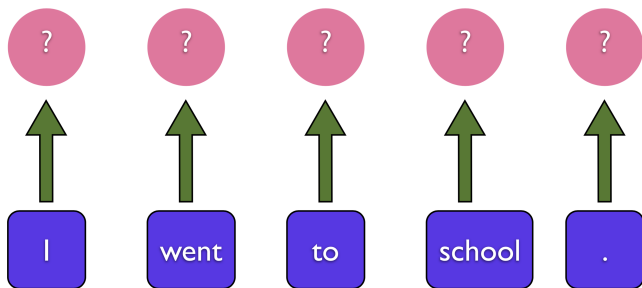
Part-of-Speech Tagging (POS)

Tag	Description	Example	Tag	Description	Example
CC	ordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &</i>
CD	cardinal number	<i>one, two</i>	TO	"to"	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential 'there'	<i>there</i>	VB	verb base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VBN	verb past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, sing.	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	"	left quote	<i>' or "</i>
POS	possessive ending	<i>'s</i>	"	right quote	<i>' or "</i>
PRP	personal pronoun	<i>I, you, he</i>	(left parenthesis	<i>[, (, {, <</i>
PRP\$	possessive pronoun	<i>your, one's</i>)	right parenthesis	<i>],), }, ></i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	<i>! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	<i>: ; ... --</i>
RP	particle	<i>up, off</i>			

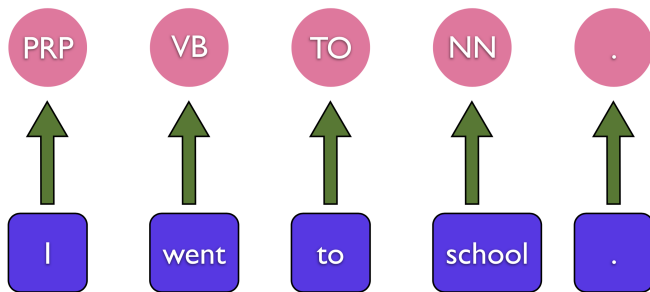
Figure: Penn Treebank POS Tags

Hidden Markov Model (HMM)

Given a sequence of words (observed states)
determine a sequence of state transitions (unobserved states)

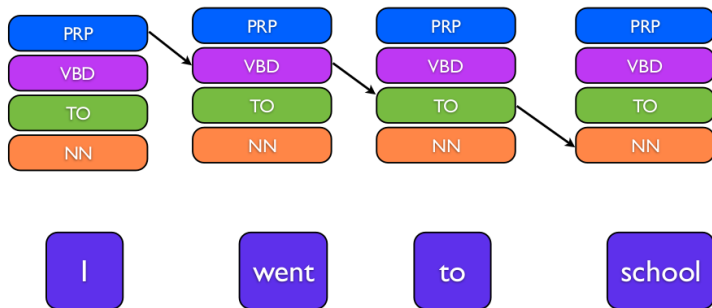


Hidden Markov Model (HMM)



Hidden Markov Model (HMM)

Why is it difficult?



Even if there were only four POS tags, then this is just one of $4 \times 4 \times 4 \times 4 = 256$ possible state sequences!

Hidden Markov Model (HMM)

Quiz 1

- If there are T number of POS tags, and you have a sentence with N number of words, how many different possible POS tag sequences can you get?

• Next Class

- Look at POS Tagging task in detail
- Hidden Markov Model
- Three Problems of HMMs
 - Decoding; Computing likelihood; Learning
- Decoding - Viterbi Algorithm



Cavnar, W. B., Trenkle, J. M., et al. (1994).

N-gram-based text categorization.

Ann Arbor MI, 48113(2):161–175.



Feldman, R. and Dagan, I. (1995).

Knowledge discovery in textual databases (kdt).

In *KDD*, volume 95, pages 112–117.



Hearst, M. A. (1997).

Texttiling: Segmenting text into multi-paragraph subtopic passages.

Computational linguistics, 23(1):33–64.



Kučera, H. and Francis, W. N. (1967).

Computational analysis of present-day American English.

Dartmouth Publishing Group.



Padró, M. and Padró, L. (2004).

Comparing methods for language identification.

Procesamiento del lenguaje natural, 33:155–162.





Cavnar, W. B., Trenkle, J. M., et al. (1994).

N-gram-based text categorization.

Ann Arbor MI, 48113(2):161–175.



Feldman, R. and Dagan, I. (1995).

Knowledge discovery in textual databases (kdt).

In *KDD*, volume 95, pages 112–117.



Hearst, M. A. (1997).

Texttiling: Segmenting text into multi-paragraph subtopic passages.

Computational linguistics, 23(1):33–64.



Kučera, H. and Francis, W. N. (1967).

Computational analysis of present-day American English.

Dartmouth Publishing Group.



Padró, M. and Padró, L. (2004).

Comparing methods for language identification.

Procesamiento del lenguaje natural, 33:155–162.

