

# COMP90042 Web Search & Text Analysis

## Workshop Week 2

---

Zenan Zhai

March 12, 2019

University of Melbourne

# Timetable

Mon	13:15	207-221-Bouverie-St-B113	Andrei
Mon	17:15	207-221-Bouverie-St-B116	Andrei
Thu	14:15	Alice-Hoy-211	Ekaterina
Thu	15:15	Alice-Hoy-211	Ekaterina
Mon	18:15	Old-Engineering-EDS4	Navnita
Wed	11:00	Elec.-Engineering-121	Navnita
Fri	10:00	221-Bouverie-St-B117	Nitika
Thu	17:15	221-Bouverie-St-B132	Nitika
Fri	15:15	Alice-Hoy-210	Shivashankar
Tue	18:15	Old-Engineering-EDS4	Shivashankar
Mon	11:00	Elec.-Engineering-121	Winn
Mon	9:00	Doug-McDonell-502	Xudong
Mon	17:15	221-Bouverie-St-B132	Xudong
Tue	16:15	221-Bouverie-St-B113	Zenan
Tue	17:15	221-Bouverie-St-B132	Zenan

- LMS - Discussion Board
- Subject Coordinator
  - A/Prof. Trevor Cohn
  - `t.cohn@unimelb.edu.au`
  - `https://trevorcohn.github.io/comp90042/`
- Me
  - Zenan Zhai
  - `zenan.zhai@unimelb.edu.au`

- Python 3
  - Virtualenv
  - Anaconda3
- Canopy EPD
  - License available when register with Unimelb Email.
- Packages
  - NLTK, gensim
  - Matplotlib, Numpy, Scipy
  - Scikit-learn

- Pre-processing
  - Pipeline
  - Lemmatisation/Stemming
- Vector Space Model
  - Document-Term Matrix/Inverted Index
  - TF-IDF/BM25
  - Exercise

## Pre-processing

---

# Pre-processing Pipeline

- Formatting
- Sentence Segmentation
- Tokenisation
- Normalisation
  - Lemmatisation
  - Stemming
- Remove Stopwords
  - May varies in different toolkit

[illegible]



# Sentence Segmentation & Tokenisation

'Ethiopian Airlines: Boeing faces questions after crash.'



['Ethiopian', 'Airlines', ':', 'Boeing', 'faces', 'questions', 'after', 'crash', '.']

- Sentence Segmentation / Tokenisation
  - Rule-based / Machine Learning
  - Varies in different languages/domains (e.g. Medicine Chemistry)
- Off-the-shelf implementations
  - NLTK  
<https://www.nltk.org/>
  - OpenNLP  
<https://opennlp.apache.org/>
  - StanfordNLP  
<https://stanfordnlp.github.io/stanfordnlp/>

- Inflectional Morphology
  - Grammatical variants
- Derivational morphology
  - Another word with different meaning

## Inflectional Morphology

airline → airlines

face → faces

question → questions

## Derivational morphology

Ethiopia → Ethiopian

# Lemmatisation & Stemming

## Lemmatisation

Remove all inflections

Matches with lexicons

Product: Lemma

## Stemming

Remove all suffixes

No matching required

Product: Stem

```
import nltk
nltk.download('wordnet')

sentence = ['Ethiopian', 'Airlines', ':', 'Boeing', 'faces', 'questions', 'after', 'crash', '.']
lemmatiser = nltk.stem.wordnet.WordNetLemmatizer()
stemmer = nltk.stem.porter.PorterStemmer()

# Code below from ...
def lemmatise(word):
    lemma = lemmatiser.lemmatize(word, 'v')
    if lemma == word:
        lemma = lemmatiser.lemmatize(word, 'n')
    return lemma
# End of copied code

lemmatised_sent = [lemmatise(word) for word in sentence]
stemmed_sent = [stemmer.stem(word) for word in sentence]

print('Sentence after lemmatisation: ' + lemmatised_sent)
print('Sentence after stemming: ', stemmed_sent)

['Ethiopian', 'Airlines', ':', 'Boeing', 'face', 'question', 'after', 'crash', '.']
['ethiopian', 'airlin', ':', 'boe', 'face', 'question', 'after', 'crash', '.']
```

## More word types $\Rightarrow$ Larger sparsity

- { 'apple': 1, 'apples':1, 'Apple': 1 }
- { 'apple': 3 }
- Stemming creates less sparsity than lemmatisation.
- When do we prefer smaller sparsity?
- Can we increase sparsity?

# Removing unwanted tokens

- Stopword

- Examples (NLTK): me, what, by, with, into, above ...

- Punctuation

- Examples: , . : ! ' " ...

## TF-IDF

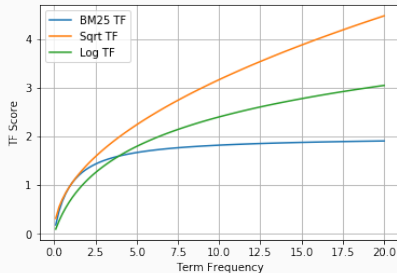
$$W_{d,t} = \underset{\text{(TF)}}{tf_{d,t}} \times \log \frac{N}{\underset{\text{(IDF)}}{df_t}}$$

## Okapi BM25

$$W_{d,t} = \frac{(k_1 + 1)tf_{d,t}}{k_1((1 - b) + b(\frac{L}{L_{avg}})) + tf_{d,t}} \times \log \frac{N - df_t + 0.5}{df_t + 0.5} \times \frac{(k_3 + 1)tf_{q,t}}{k_3 + tf_{q,t}}$$

(Document TF, document length)                      (IDF)                      (Query TF)

# TF smoothing



Raw TF

$$TF_{score} = tf_{d,t}$$

Square Root TF

$$TF_{score} = \sqrt{tf_{d,t}}$$

Log TF

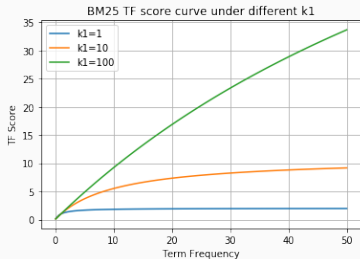
$$TF_{score} = \log(tf_{d,t})$$

BM25

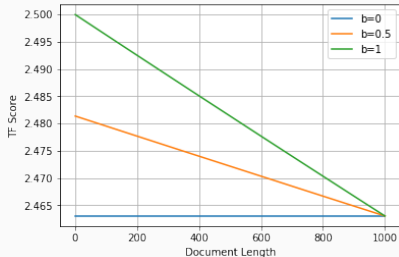
$$TF_{score} = \frac{(k+1)tf_{d,t}}{k + tf_{d,t}}$$

$$\lim_{tf_{d,t} \rightarrow \infty} \frac{(k+1)tf_{d,t}}{k + tf_{d,t}} = k+1$$

# Document TF and Document length



$$TF_{score} = f(tf_{d,t})$$



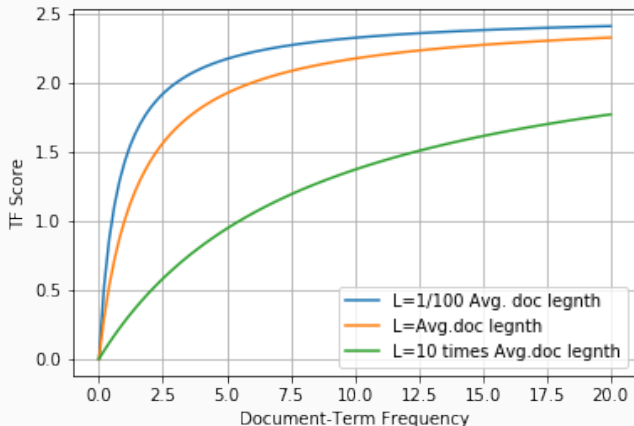
$$TF_{score} = f(L)$$

$$\frac{(k_1 + 1)tf_{d,t}}{k_1((1 - b) + b(\frac{L}{L_{avg}})) + tf_{d,t}}$$

- What does  $k_1$  controls?
- What happens when  $b = 0$ ?



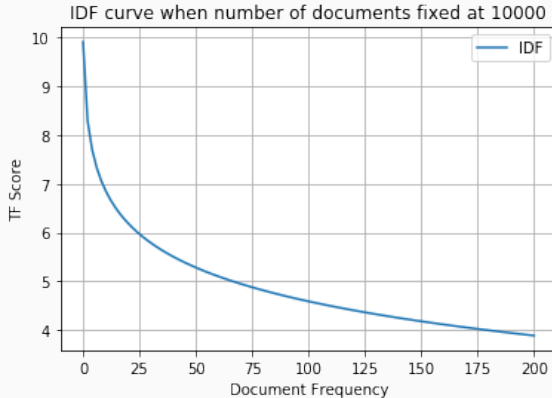
# Document length and growth of TF



$$\frac{(k_1 + 1)tf_{d,t}}{k_1((1 - b) + b(\frac{L}{L_{avg}})) + tf_{d,t}}$$

TF score grows faster when document length is short. Why?

# Inverted Document Frequency



$$IDF_{score} = \log \frac{N - df_t + 0.5}{df_t + 0.5}$$

What are “stop-words” and why are they often discarded in information retrieval? (Final exam, 2015)