

# GOAL

Ability to process and harness information from a large corpus of text with a very little human intervention

# Why it is tough?

- Multiple ways of representation of the same scenario
- Includes common sense and contextual representation
- Complex representation information (simple to hard vocabulary)
- Mixing of visual cues
- Ambiguous in nature
- Idioms, metaphors, sarcasm (Yeah! right), double negatives, etc. make it difficult for automatic processing
- Human language interpretation depends on real world, common sense, and contextual knowledge

# Ideal Properties of Corpus

- Collection of a written text in a digital form
- Useful to verify a hypothesis about a language
  - To determine how the usage of a particular sound, word, or syntactic construction varies in different contexts
  - The boys play cricket on the river bank. The boys play cricket by the side of a national bank
- Contains most of the words of a language
- Changes as a function of time - regular increase of corpus size with addition of new text samples
- Corpus is huge - Several billions of words [1]
- Even distribution of texts from all domains of language use
- Represents all areas of coverage of texts of a language
- Access of language data in an easy and simplified manner

# Lexical Resources

- A corpus is a collection of machine readable text collected according certain criteria
  - Representative collection of text
  - Used for statistical analysis and hypothesis testing
  - Used for validating linguistic rules within a specific language
- *Brown Corpus* contains a collection of written American English
  - *Sussane* is a subset of Brown, but is freely available
  - A bi-lingual parallel corpus, *Canadian Hansards*, contains French and English transcripts of the parliament
  - Penn-Treebank contains annotated text from the Wall Street journal
  - Most NLP software platforms such as *NLTK*, *Spacy* include several corpus for learning purposes

# Operations on Text Corpus

- Identify paragraphs, sentences
- Extract tokens
- Count the number of tokens/words in the corpus
- Find the vocabulary count
- Find patterns of words
- Find co-occurrence of words

The basic operation on text is *tokenization*. This is the process of dividing input text into tokens/words by identifying word boundary

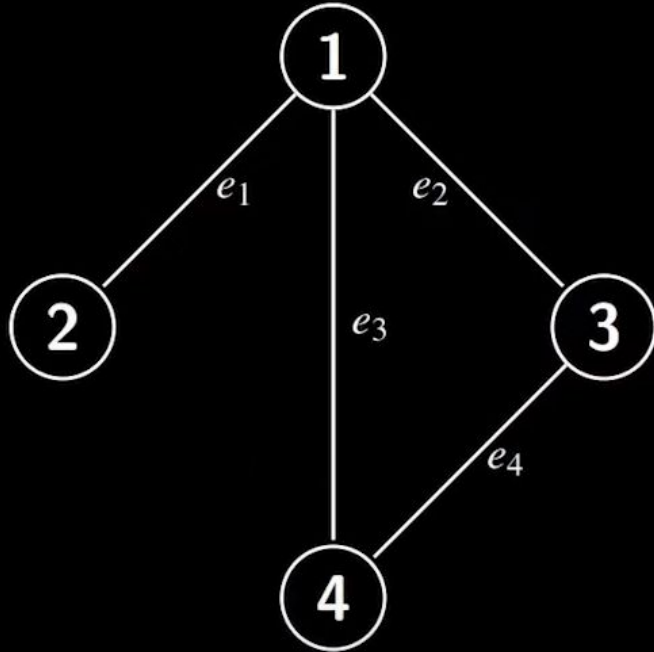
# Incidence Matrix

Let  $G$  be a graph with  $n$  vertices  $(v_1, v_2, \dots, v_n)$  and  $m$  edges  $(e_1, e_2, \dots, e_m)$ . Then *incidence matrix* defined of size  $n \times m$  is defined as

$$x_{ij} = \begin{cases} 1 & \text{if there is an edge connecting } i \text{ and } j \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

It is also called vertex-edge incidence matrix and is denoted by  $X(G)$

# Binary Incidence Matrix



The incidence matrix corresponding to the left is given below

	$e_1$	$e_2$	$e_3$	$e_4$
1	1	1	1	0
2	1	0	0	0
3	0	1	0	1
4	0	0	1	1

$$x_{ij} = \begin{cases} 1 & \text{if the edge } i \text{ connects the vertex } j \\ 0 & \text{otherwise} \end{cases}$$



# Term-Document Binary Incidence Matrix

To build a term-Document binary incidence matrix, let us consider Shakespeare's plays as our corpus. The terms are the vertices and the names of the plays are considered as edges. This incidence matrix does not represent any information related word order or its frequency[3]

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello
antony	1	1	0	0	0
brutus	1	1	0	1	0
caesar	1	1	0	1	1
calpurnia	0	1	0	0	0
cleopatra	1	0	0	0	0
...	...	...	...	...	...
...	...	...	...	...	...

$$x_{td} = \begin{cases} 1, & \text{if the word } t \in d \\ 0, & \text{if } t \notin d \end{cases} \quad (2)$$



# IR using Binary Incidence Matrix

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello
antony	1	1	0	0	0
brutus	1	1	0	1	0
caesar	1	1	0	1	1
calpurnia	0	1	0	0	0
cleopatra	1	0	0	0	0
...	...	...	...	...	...
...	...	...	...	...	...

To answer the query *Brutus AND Caesar AND NOT Calpurnia*, we take the vectors for Brutus, Caesar and Calpurnia, complement the last, and then do a bitwise AND:

$11010 \text{ AND } 11011 \text{ AND } 10111 = 10010$ . The answer for this query is found in the

# Words and Terms

The basic alphabet for the purpose of NLP is a **word**

The next logical step after the binary representation of words or terms  $t$ , is to assign weights to words

- The atomic unit for constructing a word in a language is its alphabet
- We use **Term** (co-located/co-occurring words) and **word** as atomic.
- It is necessary to consider the numerical representation of the word for computation purposes
- Vocabulary of size  $N = 1 \dots n$  is defined as  $V = w_1, w_2, w_3, \dots, w_n$  is the vocabulary containing unique words of a language
- Some words found in  $V$  appear in documents  $(D = D_1, D_2, D_3, \dots, D_m)$ , once or several times or may not appear at all.

# Term Frequency

## Term Frequency

For the given document, ***term frequency*** is defined as the number of occurrences of a term,  $t_i$ , in a document  $d_i$  belonging to a corpus  $(d_1, d_2, d_3, \dots, d_m)$ . This is denoted by  $tf_{t,d}$

# Term Frequency- Demo

```
import nltk
from nltk.probability import FreqDist
from nltk.corpus import stopwords

stop_words = set(stopwords.words('english'))

#read the corpus
words = nltk.Text(nltk.corpus.gutenberg.words('bryant-stories.txt'))
#convert to small letters
words=[word.lower() for word in words if word.isalpha() ]
words=[word.lower() for word in words if word not in stop_words ]

fDist = FreqDist(words)

#print(len(words)) #21718
#print(len(set(words))) #3688 - unique words
for x,v in fDist.most_common(10):
    print(x, v)

for x,v in fDist.most_common(10):
    print(x,v/len(fDist))
```

# Term Frequency

Raw count of words

little	597
said	453
came	191
one	183
could	158
king	141
went	122
would	112
great	110
day	107

Term frequency adjusted to document length

little	0.1618763557483731
said	0.12283080260303687
came	0.05178958785249458
one	0.04962039045553145
could	0.042841648590021694
king	0.038232104121475055
went	0.03308026030368764
would	0.03036876355748373
great	0.02982646420824295
day	0.02901301518438178

## Multiple Weighting Factors TF

*Boolean* – 0, 1 .

*RawCount* –  $tf_{i,d}$

Adjusted to document length –  $\frac{tf_{i,d}}{M}$

Log weighting –  $\begin{cases} f_{t,d} - 1 + \log tf_{i,d} & \text{if } tf_{i,d} > 0 \\ 0, & \text{otherwise} \end{cases}$

# Disadvantages of Raw Frequency

- All terms are given equal importance
- The Common term ***the*** has no relevance to the document, but gets high relevancy
- May not be suitable for classification when common words appear in documents



# Bag of Words

The collection  $[tf_1, tf_2, tf_5, tf_{15}, \dots, \dots, tf_n]$  is known as *bag of words*

- The ordering of the terms is not important
- Two documents with similar bag of words are similar in content
- It refers to the quantitative representation of the document

# Type Token Ratio

The lexical variety of the text is defined the ***Type Token Ratio (TTR)***. It can be used to measure the vocabulary variation or ***lexical density*** of the written text and speech. The ***type*** is the unique vocabulary in the text which is devoid of any repetitions.

$$TTR = \frac{V}{T_n}, \quad (7)$$

where  $V$  is the vocabulary and  $T_n$  is the number of tokens in the speech or written text.

# Python Code for TTR

```
import nltk
from nltk.corpus import stopwords
#get the stop words for English
stop_words = set(stopwords.words('english'))
words_bryant = nltk.Text(nltk.corpus.gutenberg.words('bryant-stories.txt'))
words_emma = nltk.Text(nltk.corpus.gutenberg.words('austen-emma.txt'))
#convert to small letters
words_bryant = [word.lower() for word in words_bryant if word.isalpha()]
words_emma = [word.lower() for word in words_emma if word.isalpha()]
#remove stop words
words_bryant = [word.lower() for word in words_bryant if word not in
                stop_words][:15000]
words_emma = [word.lower() for word in words_emma if word not in stop_words
             ][:15000]
TTR_bryant = len(set(words_bryant))/len(words_bryant)
TTR_emma = len(set(words_emma))/len(words_emma)
print('Number of tokens, Vocabulary, Type-token ratio (Bryant stories) = ',
      len(words_bryant), len(set(words_bryant)), TTR_bryant)
print('Number of tokens, Vocabulary, Type-token ratio (Jane Austen Emma) = ',
      len(words_emma), len(set(words_emma)), TTR_emma)
```

# TTR Demo Results

It is not reasonable to compare two unequal sized documents. A standardized TTR is used for fair comparison where the token size is restricted to the first 15000 tokens

	Number of tokens	Vocabulary	Type-token ratio
Bryant stories	15000	2796	0.19
Jane Austen (Emma)	15000	3274	0.22

# Applications of TTR

- Monitor the vocabulary usage
- Monitor child vocabulary development
- Estimate the vocabulary variation in the text

# Inverse Document Frequency

In order to attenuate the effect of frequently occurring terms, it is important to scale it down and at the same time it is necessary to increase the weight of terms that occur rarely.

Inverse document frequency (IDF) is defined as

$$IDF_t = \log \left( \frac{N}{D_{ft}} \right) \quad (8)$$

where  $N$  is the total number of documents in a collection, and  $D_{ft}$  is the count of documents containing the term  $t$

- Rare documents gets a significantly higher value
- Commonly occurring terms are attenuated
- It is a measure of informativeness
- Reduce the tf weight of a term by a factor that grows with its collection frequency.
- If a term appears in all the documents, then IDF is zero. This implies that the term is not important

# TF-IDF

Composition of TF and IDF produces a composite scaling for each term in the document

$$tf-idf_{t,d} = tf_{t,d} \times idf_{t,d}$$

- The value is high when  $t$  occurs many times within a few documents
- The value is very low when a term appears in all documents



# Inverse Document Frequency

IDF is the inverse frequency of the word 't' appearing in the corpus. It is computed as

$$IDF \text{ of a term } t = \log_{10} \left( \frac{\text{Total number of documents in a corpus}}{\text{Count of documents with term } t} \right)$$

IDF is the measure of *informativeness*

## Example:

Consider a corpus with 100000 documents. The word **moon** occurs in some documents (say, 100) with the following frequency:

$$TF_{d_1} = \frac{20}{427}, TF_{d_2} = \frac{30}{250}, TF_{d_3} = \frac{20}{250}, TF_{d_9} = \frac{5}{125} \text{ and } TF_{d_{1000}} = \frac{20}{1000}$$

The total number of words in the corpus = 100000

$$\therefore IDF_{d_1} = \log_{10} \left( \frac{100000}{100} \right)$$

$$TF_{d_1} * IDF = 0.141$$

If the word **Andromeda** appears only once  $d_1$ , then  $TF_{d_1} * IDF = 0.0117$ . If the word **the** appeared in every document and 45 times in  $d_1$ , then  $TF * IDF = 0$

So, in document one, it appears about 20 times and in a document 1000, it appears about 20 times and then rest of the numbers if you add up all those frequencies it may not be 100, in this case, rest of them probably would be appearing in other documents as well in small numbers. if you calculate the IDF for the term moon here it will be the logarithmic value of 100000 by 100 which means it is equal to 3 and then the term frequency and IDF would be equal to 0.141.

If you do not consider the TF IDF and only consider TF the word Andromeda if it appears only once in the entire corpus it will have no relevance at all. Because if you assume that it appears only once in one document let us say d 1 then it is going to be 1 by 427. So it is of very less significance, but when you consider the IDF for this it would be equal to 5 right. The IDF is going to be 5 and then the TF IDF is going to be 0.0117 it is here. if the word d appeared in every document then the TF IDF will be equal to 0, which means the entire word d will not be considered at all in the relevancy of documents

# Document Ranking using TF-IDF

Using the TF-IDF, the rank order for the documents can be determined for the documents for the term *moon*.

Document Name	tf	tf-idf	Rank
d1	0.047	0.14	3
d2	0.012	0.36	1
d3	0.08	0.24	2
d9	0.04	0.12	4
d1000	0.02	0.06	5

# ZIPF's Law

Zipf's law states that for a given some corpus, the frequency of any word is inversely proportional to its rank in the term frequency table[3]

$$f(r) \propto \frac{1}{r^\alpha}$$

where  $\alpha \approx 1$ ,  $r$  is the *frequency rank* of a word and  $f(r)$  is the frequency in the corpus. The most frequent word will have the value 1, the word ranked second in the frequency will have  $\frac{1}{2^\alpha}$ , the word ranked third in the frequency will have  $\frac{1}{3^\alpha}$ , etc

## Distribution of terms/words

This empirical law models the frequency distribution of words in languages. This distribution is observed across several languages with a large corpus. It may not be good enough to fit the frequency linearly, but enough to approximately model word frequencies

# Python's Code for ZIPF's Law

```
import re
from operator import itemgetter
import nltk
import pandas as pd

frequency = {}
words_emma = nltk.Text(nltk.corpus.gutenberg.words('austen-emma.txt'))

for word in words_emma:
    count = frequency.get(word, 0)
    frequency[word] = count + 1

rank = 1
column_header = ['Rank', 'Frequency', 'Frequency*Rank']
df = pd.DataFrame(columns=column_header)

for word, freq in reversed(sorted(frequency.items(), key=itemgetter(1))):
    df.loc[word] = [rank, freq, rank * freq]
    rank = rank + 1
print(df)
```

# ZIPF's Law Demo

Word	Frequency	Rank	Frequency*Rank
to	5183	3	15549
the	4844	4	19376
and	4672	5	23360
of	4279	6	25674
I	3178	7	22246
as	1387	21	29127
—	1382	22	30404
he	1365	23	31395
for	1321	24	31704
have	1301	25	32525
is	1220	26	31720
with	1187	27	32049
Mr	1153	28	32284
very	1151	29	33379
but	1148	30	34440