

# Text Pre-Processing

In [1]:

```
1 import pandas as pd
```

In [2]:

```
1 f = open("allen.txt")
```

In [3]:

```
1 # print(f.read())  
2
```

In [4]:

```
1 data = f.read()
```

In [5]:

```
1 print(data)
```

Alan Mathison Turing OBE FRS (/ˈtʃʊɹɪŋ/; 23 June 1912 – 7 June 1954) was an English mathematician, computer scientist, logician, cryptanalyst, philosopher, and theoretical biologist.[6] Turing was highly influential in the development of theoretical computer science, providing a formalisation of the concepts of algorithm and computation with the Turing machine, which can be considered a model of a general-purpose computer.[7][8][9] He is widely considered to be the father of theoretical computer science and artificial intelligence.

## Pre-Processing

1. LowerCase Conversion
2. Sentence Tokenization
3. Stemming

## LowerCase

In [6]:

```
1 data.lower()
```

Out[6]:

```
'alan mathison turing obe frs (/ë^tjêšë™réå</; 23 june 1912 â€“ 7 june 19
54) was an english mathematician, computer scientist, logician, cryptanaly
st, philosopher, and theoretical biologist.[6] turing was highly influenti
al in the development of theoretical computer science, providing a formali
sation of the concepts of algorithm and computation with the turing machin
e, which can be considered a model of a general-purpose computer.[7][8][9]
he is widely considered to be the father of theoretical computer science a
nd artificial intelligence.'
```

## Stemming

In [7]:

```
1 import re
2
3 def remove_punctuation(text):
4     return re.sub(r'^\w\s', '', text)
```

In [8]:

```
1 clean_text = remove_punctuation(data)
```

In [9]:

```
1 print(clean_text)
```

```
Alan Mathison Turing OBE FRS Ė^tjÊŠÉrÉå 23 June 1912 â 7 June 1954 was an
English mathematician computer scientist logician cryptanalyst philosopher
and theoretical biologist6 Turing was highly influential in the developmen
t of theoretical computer science providing a formalisation of the concept
s of algorithm and computation with the Turing machine which can be consid
ered a model of a generalpurpose computer789 He is widely considered to be
the father of theoretical computer science and artificial intelligence
```

## Tokenization

In [10]:

```
1 from nltk.tokenize import sent_tokenize
```

In [11]:

```
1 sent_tokenize(clean_text)
```

Out[11]:

```
['Alan Mathison Turing OBE FRS 23 June 1912 â 7 June 1954 was an English mathematician computer scientist logician cryptanalyst philosopher and theoretical biologist6 Turing was highly influential in the development of theoretical computer science providing a formalisation of the concepts of algorithm and computation with the Turing machine which can be considered a model of a generalpurpose computer789 He is widely considered to be the father of theoretical computer science and artificial intelligence']
```

In [12]:

```
1 from nltk.tokenize import word_tokenize
```

In [17]:

```
1 tokens = word_tokenize(clean_text)
```

In [18]:

1	tokens
---	--------

Out[18]:

```
[ 'Alan',
  'Mathison',
  'Turing',
  'OBE',
  'FRS',
  'ĖˆtjĖŠĖrĖāĀ',
  '23',
  'June',
  '1912',
  'ā',
  '7',
  'June',
  '1954',
  'was',
  'an',
  'English',
  'mathematician',
  'computer',
```

```
In [19]: print(len(tokens))
'scientist',
'logician',
```

```
75cryptanalyst',
'philosopher',
'and',
```

```
'theoretical',
'wednesday',
'Turing',
```

```
'was',
In [14]: highly,
```

```
'influential'
In [14]: import nltk
In [14]: from nltk.corpus import stopwords
'the',
```

```
'development',
In [15]: of,
```

```
'theoretical',
In [15]: nltk.download('stopwords')
'computer',
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\MSCIT\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
'formalisation',
```

```
Out[15]:
```

```
'the',
True 'concepts',
'of',
```

```
In [16]: algorithm',
```

```
'and',
In [16]: stop_words = set(stopwords.words('english'))
'computation',
```

```
'with',
```

```
In [21]: the
```

```
'Turing',
In [21]: machine = list(stop_words)
```

```
'which',
```

```
'can',
```

```
'be',
```

```
'considered',
```

```
'a',
```

```
'model',
```

```
'of',
```

```
'a',
```

```
'generalpurpose',
```

```
'computer789',
```

```
'He',
```

## Removing Stop Words

```

'is',
'widely',
'considered',
'to',
'best',
'best',
'father',
'you',
'theoretical',
'computer',
'balance',
'bad',
'wouldn't',
'shouldn't']
'no',
'mustn',
'under',
'only',
'what',
"aren't",
'whom',
'from',
'isn',
'how'

```

In [25]:

```

1 filtered_words = [word for word in tokens if word not in stp]

```

In [26]:

1	filtered_words
---	----------------

Out[26]:

```
['Alan',  
'Mathison',  
'Turing',  
'OBE',  
'FRS',  
'ĖˆtjĖŠĖrĖÅ',  
'23',  
'June',  
'1912',  
'â',  
'7',  
'June',  
'1954',  
'English',  
'mathematician',  
'computer',  
'scientist',  
'logician',  
'cryptanalyst',  
'philosopher',  
'theoretical',  
'biologist6',  
'Turing',  
'highly',  
'influential',  
'development',  
'theoretical',  
'computer',  
'science',  
'providing',  
'formalisation',  
'concepts',  
'algorithm',  
'computation',  
'Turing',  
'machine',  
'considered',  
'model',  
'generalpurpose',  
'computer789',  
'He',  
'widely',  
'considered',  
'father',  
'theoretical',  
'computer',  
'science',  
'artificial',  
'intelligence']
```



In [29]:

```
1 print("Before Removing Stop Words : ",len(tokens))
2
3
4 print("After Removing Stop Words : ",len(filtered_words))
```

Before Removing Stop Words : 75

After Removing Stop Words : 49

## Stemmimng

In [28]:

```
1 l1 = ['walk', 'walking', 'walked']
```

In [30]:

```
1 from nltk.stem import PorterStemmer
```

In [31]:

```
1 ps = PorterStemmer()
```

In [32]:

```
1 for w in l1:
2     print(w, " : ", ps.stem(w))
```

walk : walk

walking : walk

walked : walk

In [35]:

```
1 l2 = ['Probability', 'Probable', 'Probably']
```

In [36]:

```
1 for w in l2:
2     print(w, " : ", ps.stem(w))
```

Probability : probabl

Probable : probabl

Probably : probabl

## Lemmatization

In [37]:

```
1 from nltk.stem import WordNetLemmatizer
2
3 lemmatizer = WordNetLemmatizer()
```

In [44]:

```
1 l3 = ['rocks', 'apples', 'fruits']
```

In [45]:

```
1 for w in l3:  
2     print(w, " : ", lemmatizer.lemmatize(w))
```

```
rocks : rock  
apples : apple  
fruits : fruit
```

In [41]:

```
1 print("rocks :", lemmatizer.lemmatize("rocks"))
```

```
rocks : rock
```

In [ ]:

```
1
```