

# CS445 – Project 2

Project Title: Text Exploration

Written by: Alper Bingöl

Instructor: Reyhan Yeniterzi

Date: 10.12.2020



# 1 WordCloud

The very first function "*createWordCloud*" where frequencies of words in a given collection is represented. Some required setups should be done. These are not only for WordCloud but also for the whole project.

```
1  import nltk
2  nltk.download('stopwords')
3  from nltk.corpus import stopwords
4  stopWords = set(stopwords.words('turkish'))
5  # Start with loading all necessary libraries
6  import numpy as np
7  import pandas as pd
8  from os import path
9  from PIL import Image
10 from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator
11
12 import matplotlib.pyplot as plt
13 % matplotlib inline
14 from matplotlib import pylab
15
16 from sklearn.feature_extraction.text import CountVectorizer
17 from sklearn.feature_extraction.text import TfidfTransformer
18 from sklearn.feature_extraction.text import TfidfVectorizer
19 from sklearn.pipeline import Pipeline
20
21 from nltk.tokenize import sent_tokenize, word_tokenize
22 import warnings
23
24 #warnings.filterwarnings(action = 'ignore')
25
26 import gensim
27 from gensim.models import Word2Vec
28
29 import pickle
```

Figure 1: Setup for the Project

Sidenote: All the figures in the report are generated according to "Tsample5000.pkl" dataset.

The function initially checks if stopwords are wanted or not. If yes, it generates WordCloud after eliminating stopwords which is downloaded from `nltk.corpus`. The stopwords list is for Turkish. The width and height of wordcloud is multiplied with 100 intentionally. Otherwise the cloud would be too small. Background color is chosen as white. Design arrangements could be done by the users purpose. One important specification is using mask as a background shape. One example could be found in the continuation of the text. The second check is for weighting. It is either TDIDF or just TF. The term weight is calculated according to that. Another function is used for term weight. In this extra function, called `termweight`, `CountVectorizer` and `TfidfTransformer` are used. In `CountVectorizer` an attribute called `maxdf` is used to eliminate too common words in the collection in order to increase the amount of relevance. `maxdf` is chosen as 0.5. It is an hyperparameter.

## 1.1 Examples



Figure 2: mode="TFIDF", stopwords=False, maxdf = 0.5

Since the stopwords generated from nltk is not too much, the difference between is not huge at the first observation. But there are small differences between figure 2 and figure 3.

In the figure 4 the TFIDF values are not used. So the amount of words which includes information is less then figure 2 and 3.





## 2 Zipf's Law

This figure is drawn using the ranks and frequencies of words in the collection. It shows that term frequency decreases rapidly as a function of rank. The formula of the law:

$$f_t = \frac{k}{r_t}$$

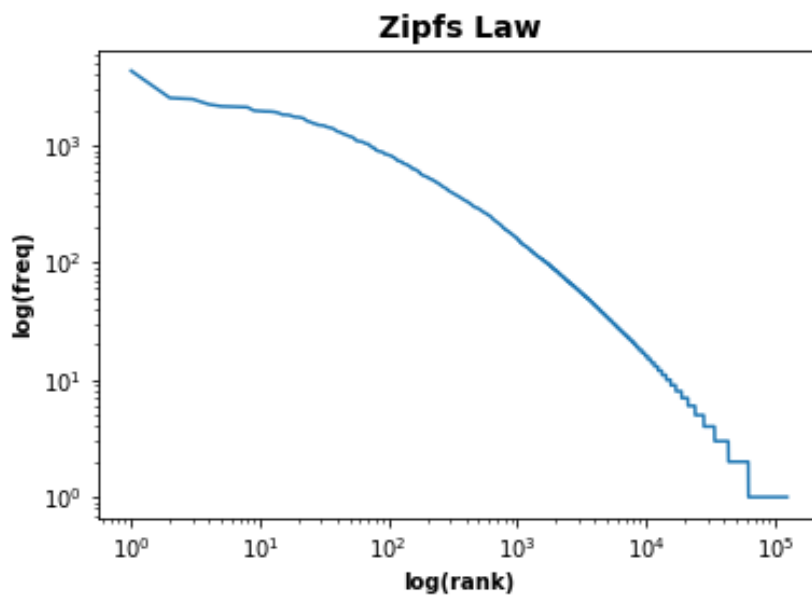


Figure 7: Zipf's plot

### 3 Heaps' Law

As the law indicates, the number of new words decreases as the size of the corpus increases. The formula of the law:

$$v = k \times n^{\beta}$$

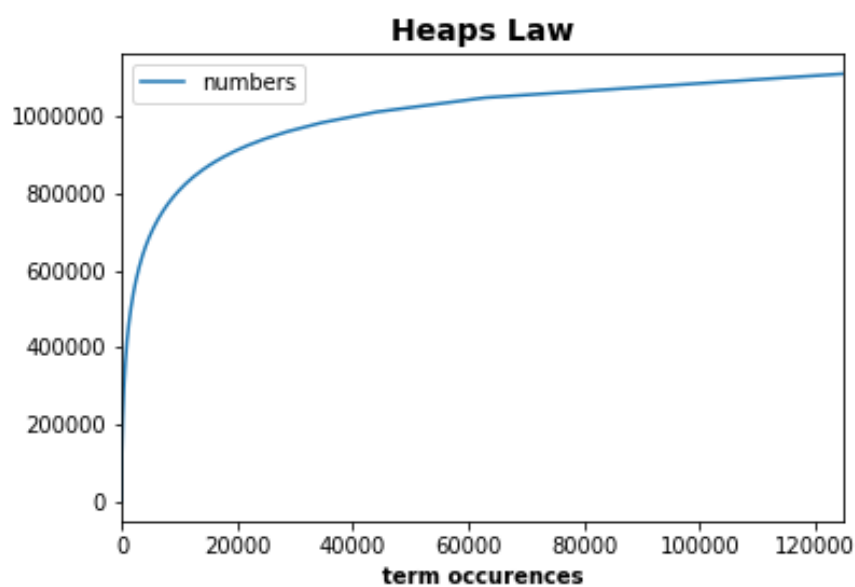


Figure 8: Heaps' plot

## 4 WordVectors

The function returns the models for trained word embeddings. The generated models are generated as the users desire. It is either CBOW or Skip-gram.

```
def create_WordVectors(Docs,dim,mode,size):

    counts = count(Docs)
    aa=counts.to_frame()
    cumulative_sum = aa.cumsum()
    cumulative_sum.rename(columns = {0:'numbers'}, inplace = True)
    cumulative_sum=cumulative_sum.reset_index()
    cumulative_sum.rename(columns = {"index":'words'}, inplace = True)
    unique_words=cumulative_sum['words']
    unique_list = unique_words.tolist()

    if mode == "cbow":
        model = gensim.models.Word2Vec(unique_list, min_count = 1, size = dim, window = size)
    else:
        model = gensim.models.Word2Vec(unique_list, min_count = 1, size = dim, window = size, sg = 1)

    return model
```

Figure 9: WordVectors Fuction



## 5 References

<https://kavita-ganesan.com/tfidftransformer-tfidfvectorizer-usage-differences/.X9KBsdgzYvj>

<https://towardsdatascience.com/simple-wordcloud-in-python-2ae54a9f58e5>

[https://amueller.github.io/wordcloud/auto\\_examples/frequency.html](https://amueller.github.io/wordcloud/auto_examples/frequency.html)

<https://www.datacamp.com/community/tutorials/wordcloud-python>

[https://scikitlearn.org/stable/modules/generated/sklearn.feature\\_extraction.text.TfidfVectorizer.html](https://scikitlearn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html)

[https://scikitlearn.org/stable/modules/generated/sklearn.feature\\_extraction.text.CountVectorizer.html](https://scikitlearn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html)

[https://matplotlib.org/api/as\\_gen/matplotlib.pyplot.html](https://matplotlib.org/api/as_gen/matplotlib.pyplot.html)

<https://pypi.org/project/wordcloud/>

<https://www.nltk.org/api/nltk.lm.html>

<https://pypi.org/project/gensim/>