

DBMS Assignment - 5

Text mining using apache spark



Problem statement

Implementation of various text mining paradigms on wiki pages.

The various paradigms which were implemented were:

1. Single word frequency
2. Two word frequency
3. Cosine similarity
4. POS tagging

1. Single word frequency

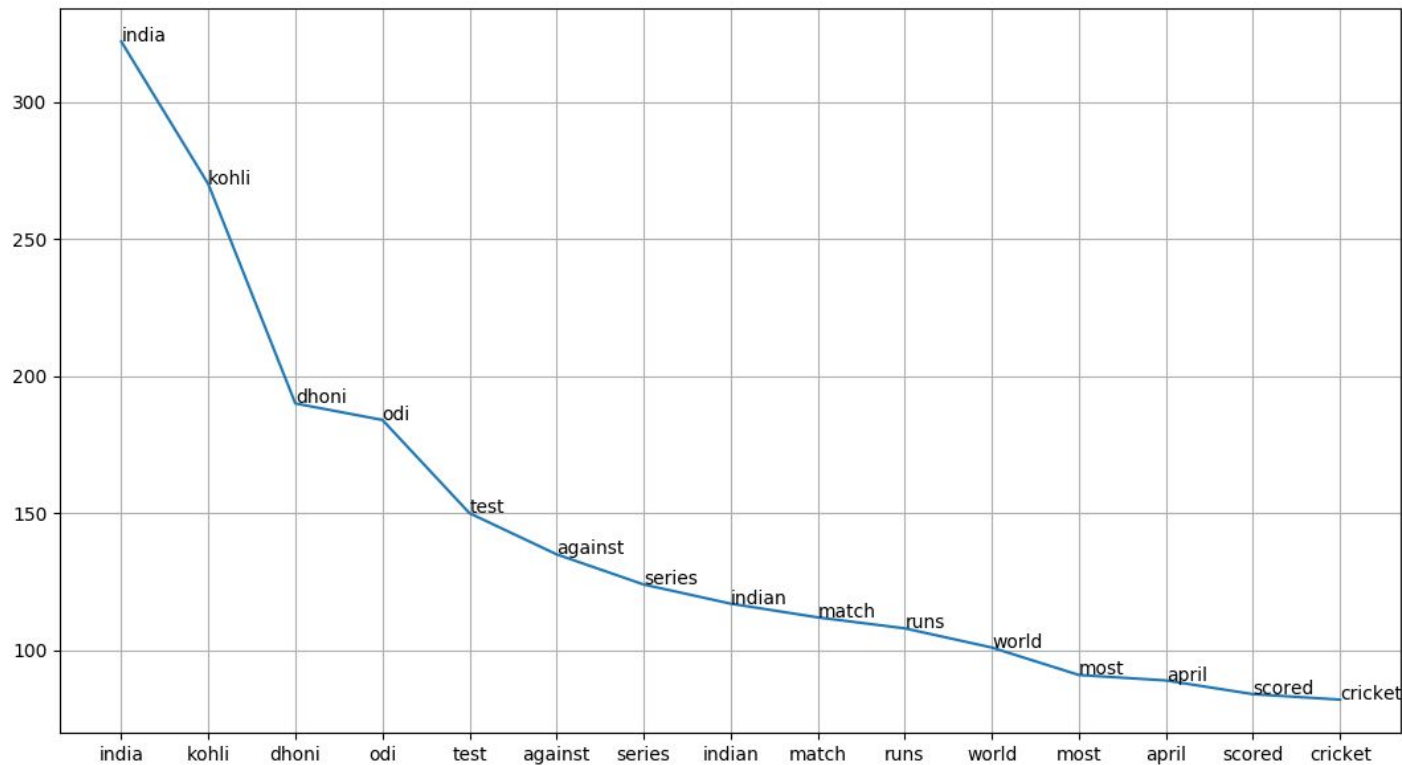
Simple word count was applied to the wiki corpus which was implemented using pyspark.

The graph of the most frequent words was plotted against their frequencies

Figure 1



Figure 1

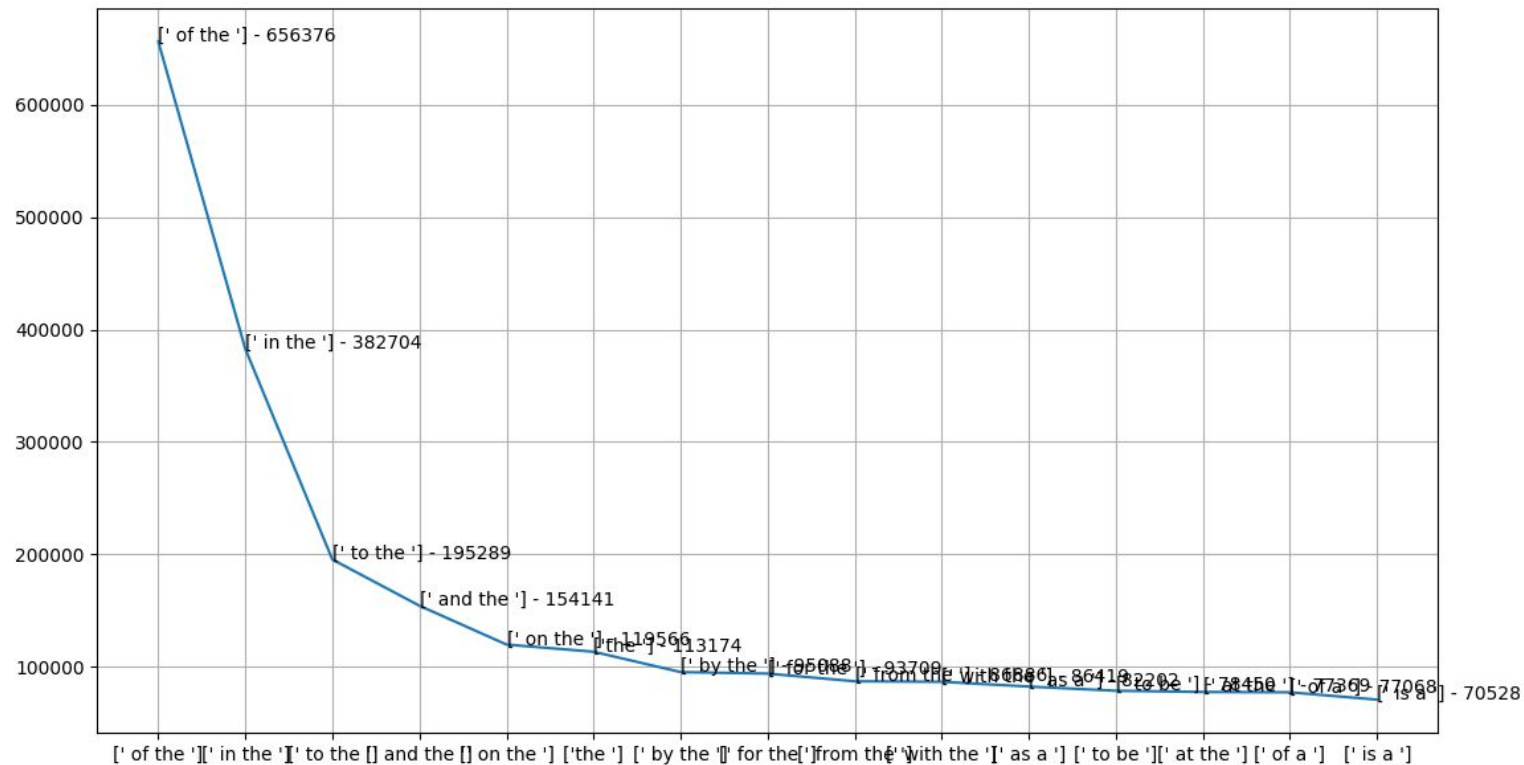


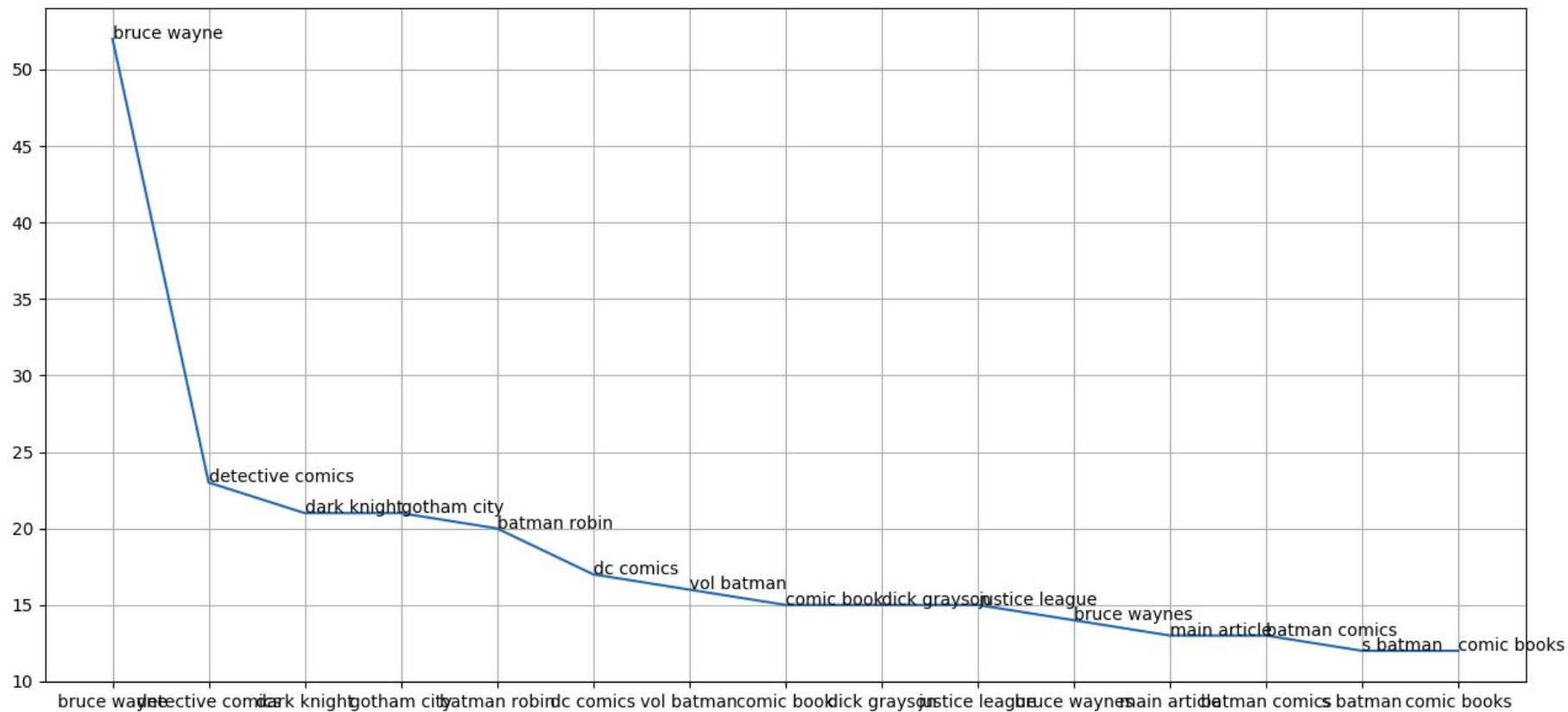
2. Two word frequency

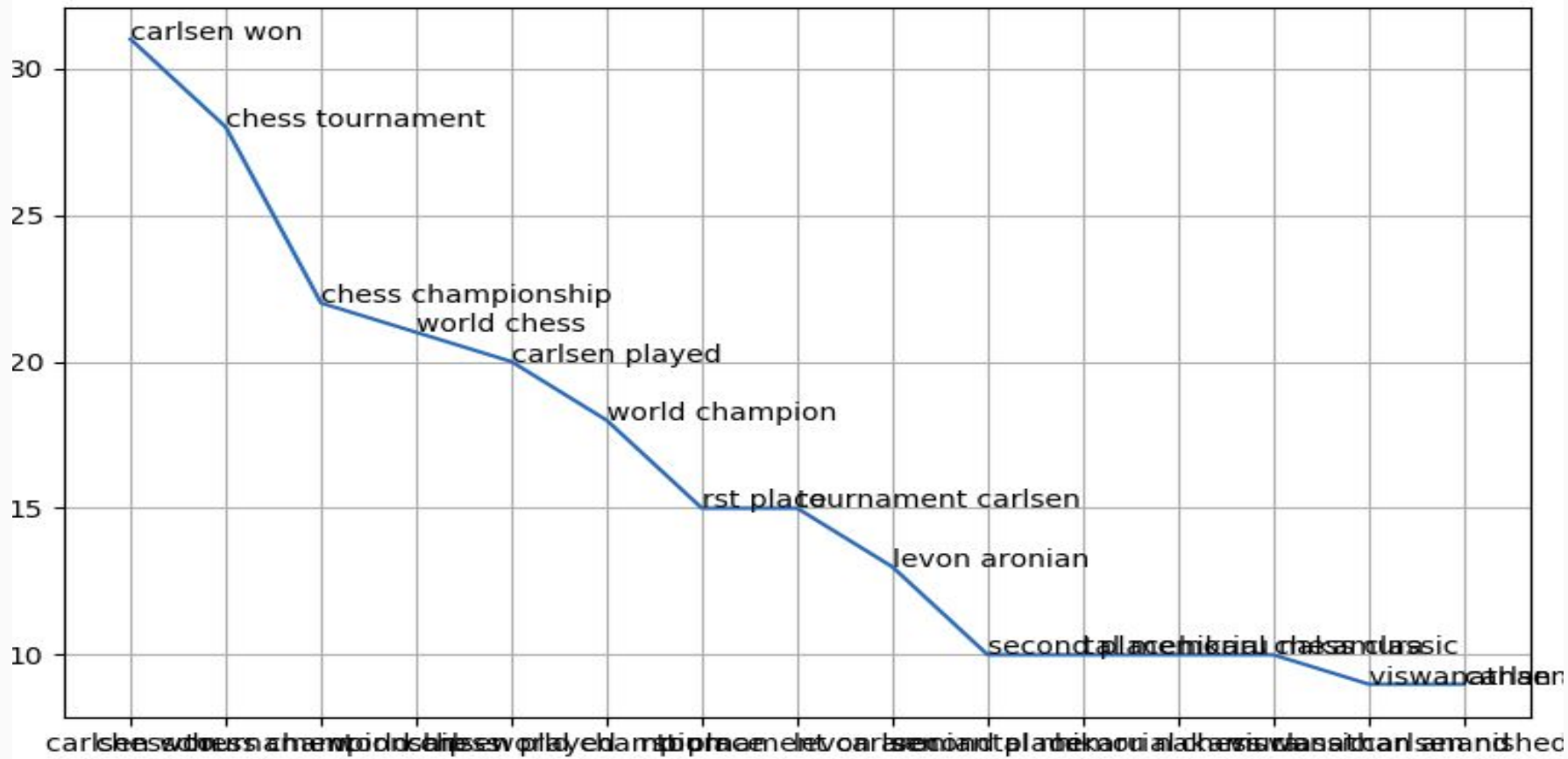
The most frequently used two word pairs were calculated using the RDDs available in pyspark. The words were split and necessary filters were applied to remove unnecessary words.

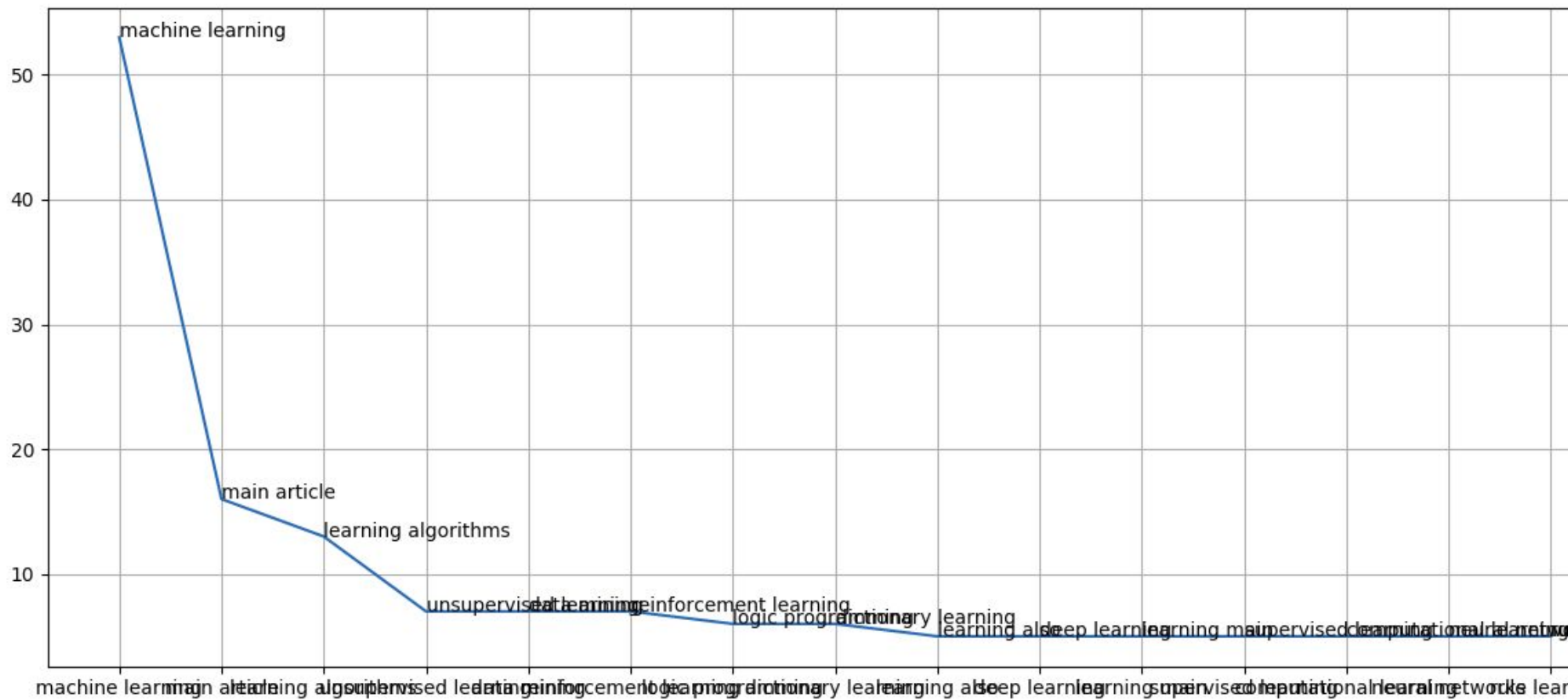
The graph of the most frequent words was plotted against their frequencies

Figure 1









3. Cosine similarity

Calculating the term frequency i.e number of times word occurs in a document. The frequency can be calculated by normalization (by dividing with the number of words in the document).

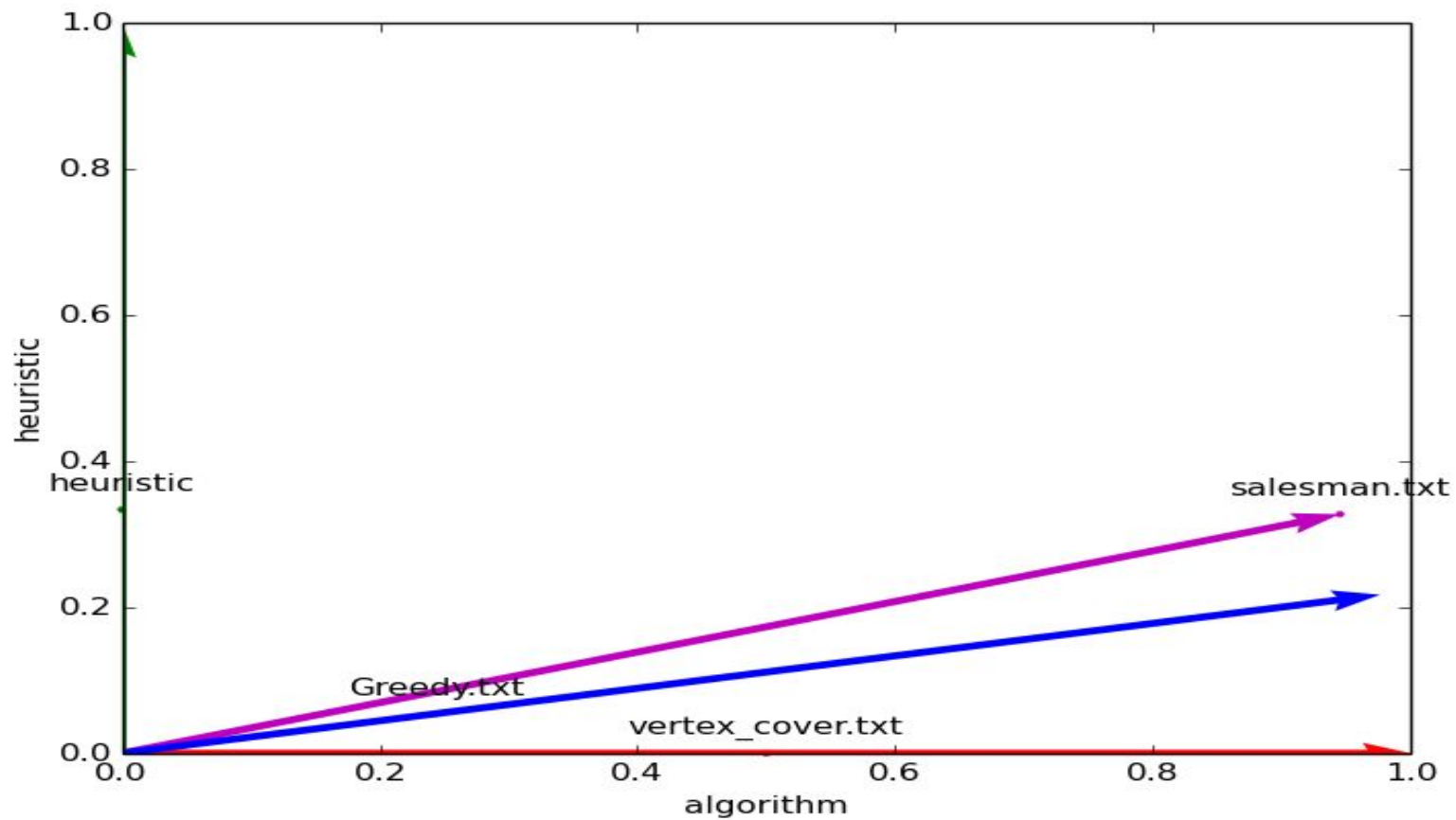
Then we find the Inverse Document frequency by:

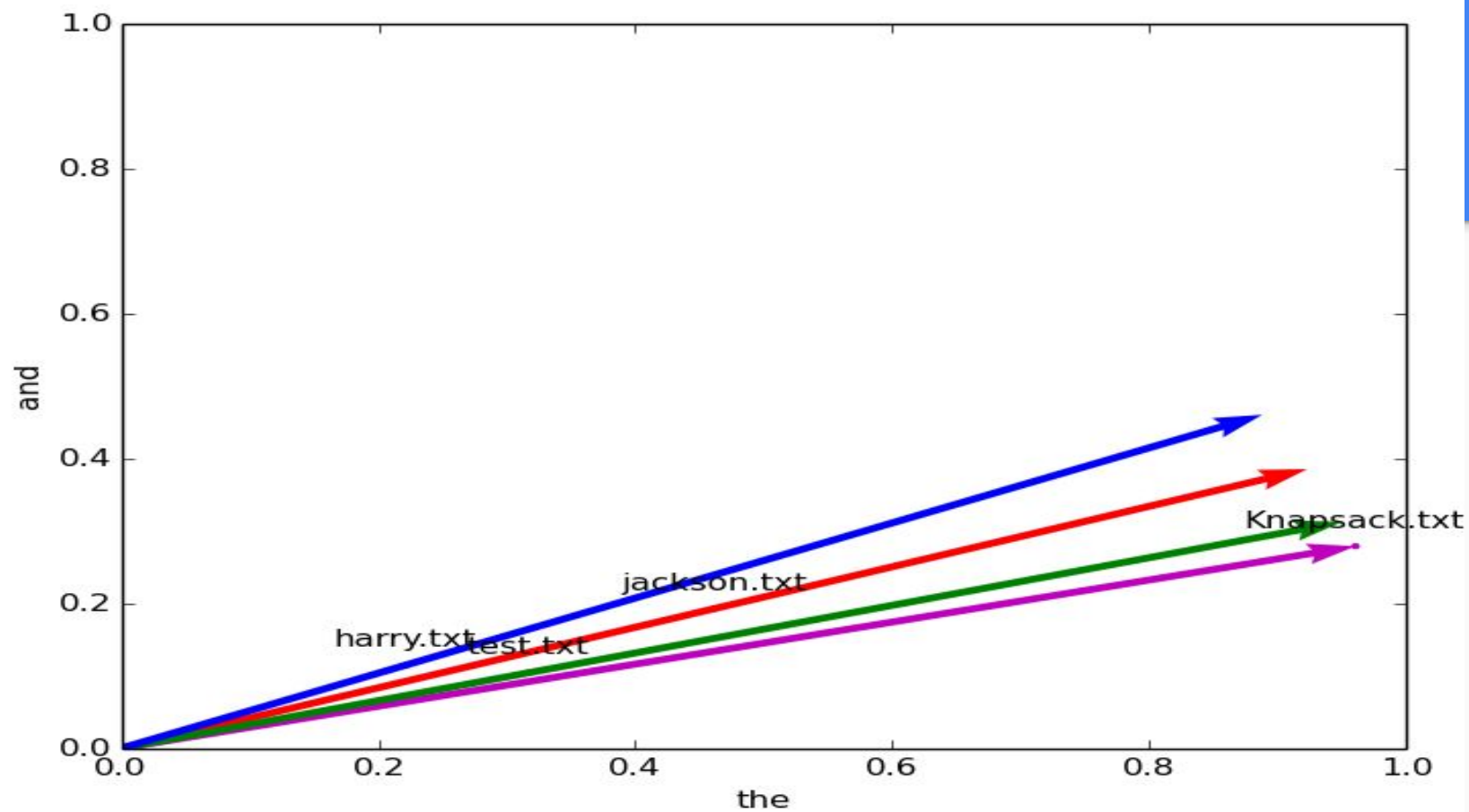
$$\text{IDF}(\mathbf{word}) = 1 + \log_e(\text{Total Number Of Documents} / \text{Number Of Documents with term } \mathbf{word} \text{ in it})$$

Since we are finding out the relevancy of each document with respect to the query of two words, we multiply and store the product of term frequency and Inverse document frequency of each word in query in each document.

Once the vectors are drawn, the angle between them can be computed by dot product, and the similarity between the lines is the measure of angle between them.

Corresponding graphs are drawn to represent this feature.





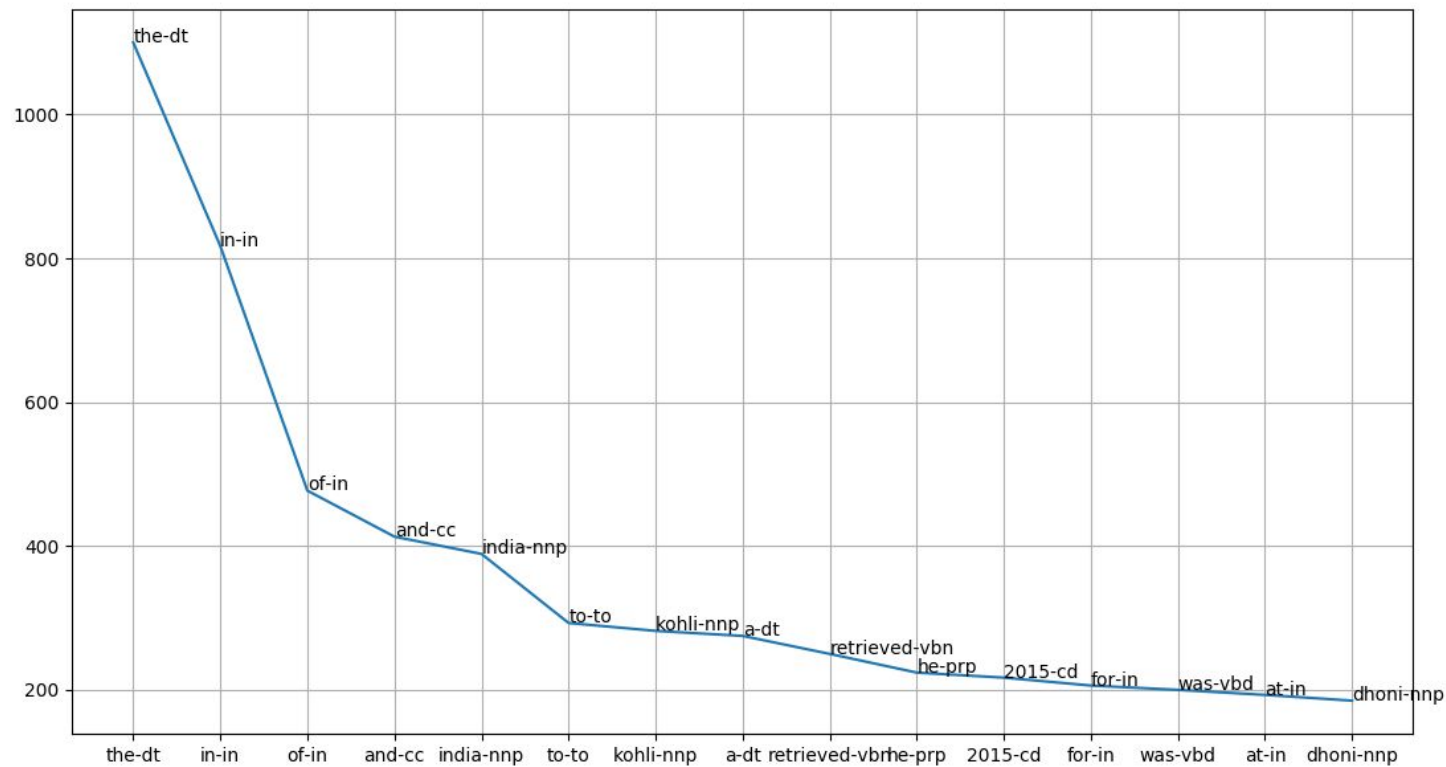
4. POS tagging

Each word in a sentence is tagged corresponding to its part of speech, based on its definition and context.

NLTK library of python was used in implementing this paradigm.

The graph of the most frequent words along with their POS tags were plotted against their frequencies

Figure 1



THANK YOU