

## Assignment No-9

Implement Document Retrieval and ranking using tf-idf algorithm. Also implement the process of the document.

Software Requirement- Windows 10,11 Python 3.8 or higher

Theory -

what is text tokenisation and lower casing?

The process of converting a sequence of text into smaller parts, known as tokens. Tokenization breaks text into smaller parts for easier machine analysis, helping machines understand human language. Lower casing: Converting a word to lower case

Removing special characters

Removing special character is necessary because characters does not add value to the text. One way to utilize it is to parse text upon the occur of any special character or indicate the need for expansion.

contraction expansion

Contractions are words or combinations of words that are shortened by dropping letters and replacing them by an apostrophe.

Stemming and Lemmanization

What is Stemming?

Stemming is a technique used to extract the base form of the words by

removing affixes from them. It is just like cutting down the branches of a tree to its stems. For example, the stem of the words eating, eats, eaten is eat.

Search engines use stemming for indexing the words. That's why rather than storing all forms of a word, a search engine can store only the stems. In this way, stemming reduces the size of the index and increases retrieval accuracy.

What is Lemmatization?

Lemmatization technique is like stemming. The output we will get after lemmatization is called 'lemma', which is a root word rather than root stem, the output of stemming. After lemmatization, we will be getting a valid word that means the same thing.

NLP based Features

NLP is used to understand the structure and meaning of human language by analysing different aspects like syntax, semantics etc. NLP is beneficial

TF-IDF stands for Term Frequency Inverse Document Frequency of records. It can be defined as the calculation of how relevant a word in a series or corpus is to a text. The meaning increases proportionally to the number of times in the text a word appears but is compensated by the word frequency in the corpus (data-set).

TF-IDF

Term Frequency: In document  $d$ , the frequency represents the number of instances of a given word  $t$ . Therefore, we can see that it becomes more relevant when a word appears in the text, which is rational. Since

the ordering of terms is not significant, we can use a vector to describe the text in the bag of term models. For each specific term in the paper, there is an entry with the value being the term frequency.

The weight of a term that occurs in a document is simply proportional to the term frequency.

$$tf(t,d) = \text{count of } t \text{ in } d / \text{number of words in } d$$

Document Frequency: This tests the meaning of the text, which is very similar to TF, in the whole corpus collection. The only difference is that in document  $d$ , TF is the frequency counter for a term  $t$ , while  $df$  is the number of occurrences in the document set  $N$  of the term  $t$ . In other words, the number of papers in which the word is present is DF.

$$df(t) = \text{occurrence of } t \text{ in documents}$$

Inverse Document Frequency: Mainly, it tests how relevant the word is. The key aim of the search is to locate the appropriate records that fit the demand. Since  $tf$  considers all terms equally significant, it is therefore not only possible to use the term frequencies to measure the weight of the term in the paper. First, find the document frequency of a term  $t$  by counting the number of documents containing the term:

$$df(t) = N(t)$$

where

$$df(t) = \text{Document frequency of a term } t$$

$$N(t) = \text{Number of documents containing the term } t$$

Term frequency is the number of instances of a term in a single

document only; although the frequency of the document is the number of separate documents in which the term appears, it depends on the entire corpus. Now let's look at the definition of the frequency of the inverse paper. The IDF of the word is the number of documents in the corpus separated by the frequency of the text.

Conclusion-Hence in this experiment we have implemented the sub-process for page ranking and document retrieval using python libraries.