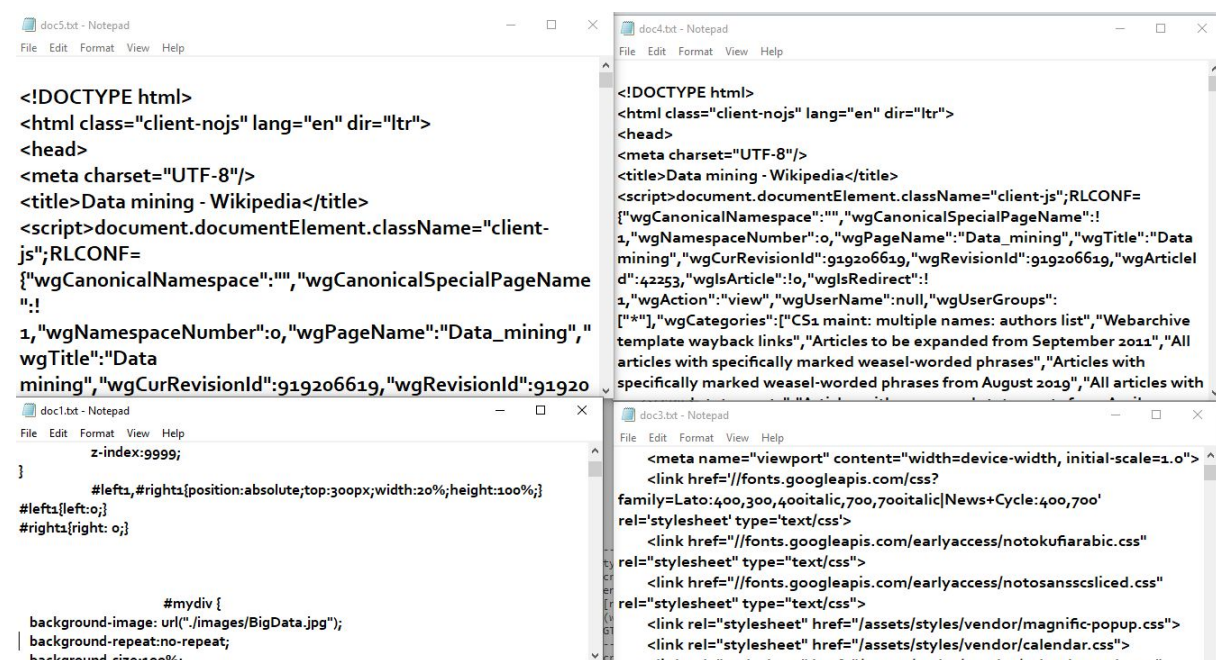**Part 1: Preprocessing to Build Document Vectors for Web Page Content Analysis**

## Step 1: Analysing Dataset

The first step in any of the Machine Learning tasks is to analyse the data. So if we look at the dataset, the data is of html source code.

Now one of the important tasks is to identify the body(important text), if we analyse the documents, there are different patterns of alignment. Now we need to figure out a way to extract the Content of the page.



## Step 2: Preprocessing

Preprocessing is one of the major steps when we are dealing with any kind of text models. During this stage we have to look at the distribution of our data, what techniques are needed and how deep we should clean.

This step totally depends on the problem statement. Few mandatory preprocessing are converting to lowercase, removing punctuation, removing stop words and lemmatization/stemming. In this problem removing <...> tags is also our main preprocessing step.

**Tag Removal**

Removing all the tags from the text file which are between < >, and storing the content back to textfiles.

tag_removal = re.compile(r'<[^>]+>')

def remove_tags(text):

```
        return tag_removal.sub(",text)
```

'''The process of removing html tags from the text by replacing with null'''

```
for i in range(1,6):
        string = open('doc'+str(i)+'.txt').read()
        open('doc'+str(i)+'.txt','w').write(remove_tags(string))
```

## Special Symbols

''' #Removing all the symbols like punctuation marks ,question marks etc...

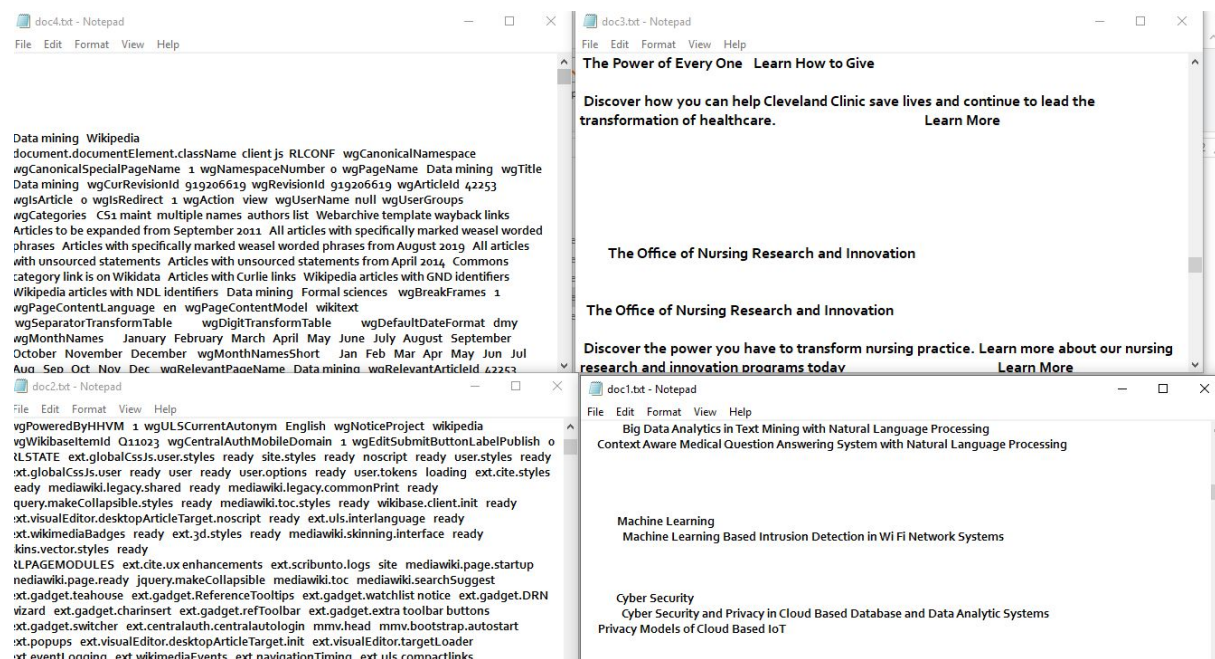      and replacing by '  '(space)'''

```
for i in range(1,6):
   string1 = open('doc'+str(i)+'.txt').read()
   new_str = re.sub('[^a-zA-Z0-9\n\.]',' ',string1)
   string1 = re.sub("\S*\d\S*", "", new_str).strip()
   open('doc'+str(i)+'.txt','w').write(string1)
```

After removing all the tags and special symbols from the text files.



## Stop Words

Stop words are the most commonly occurring words which don't give any additional value to the document vector. in-fact removing these will increase computation and space efficiency. nltk library has a method to download the stopwords, so instead of explicitly mentioning all the stopwords ourselves we can just use the nltk library and iterate over all the words and remove the stop words.
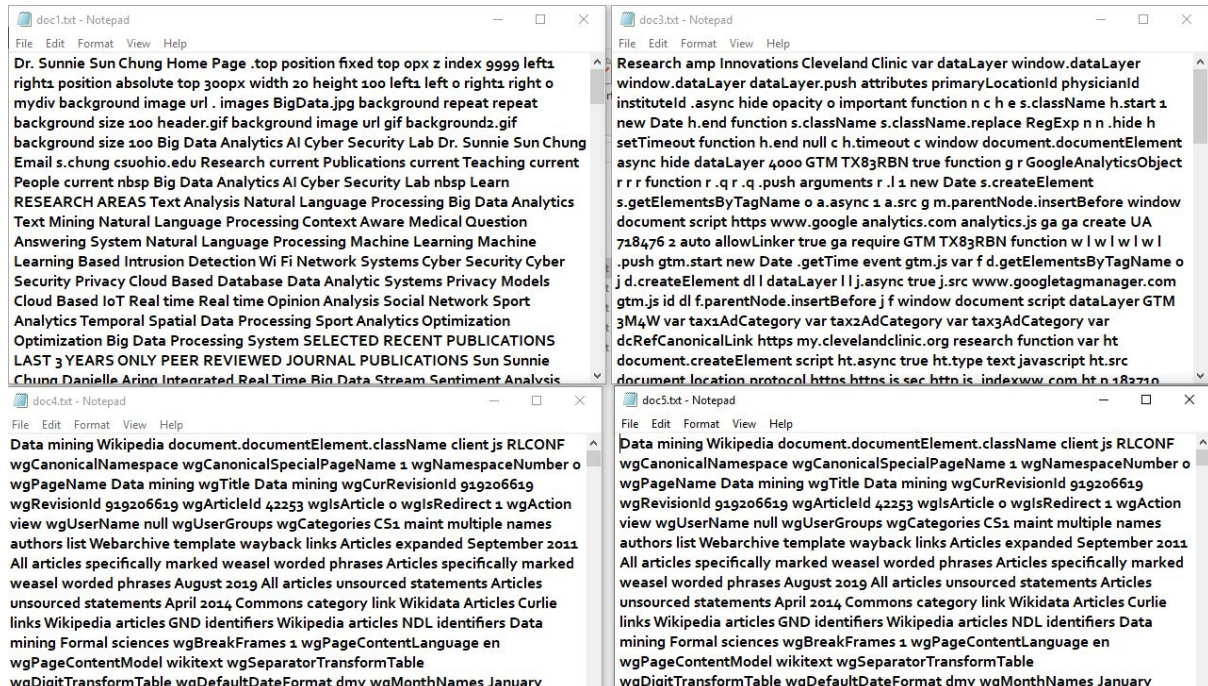
'''Removing stop words using tokenize and stop_words in python '''

```
stop_words = stopwords.words("english")

for i in range(1,6):

        string2= open('doc'+str(i)+'.txt').read()

        open('doc'+str(i)+'.txt','w').write(' '.join([word for word in string2.split() if word not in
stop_words]).lower()+' ')
```



'''#Stemming data

```
stemmer  = PorterStemmer()

temp_string = ""

for i in range(1,6):

        string = open('doc'+str(i)+'.txt').read()

        tokens = word_tokenize(string)

        for w in tokens:

        temp_string = temp_string+" "+stemmer.stem(w)

        open('doc'+str(i)+'.txt','w').write(temp_string)
```

'''

**Creating Tokens**

We have the preprocessed text files with only content after removing all the stop words,html tags,punctuation and converting to lowercase .We create a List in the python file which saves the tokens that is each word of the text.

'''Creating tokens '''

processed_text = []

for i in range(1,6):

      string3 = open('doc'+str(i)+'.txt').read()

      processed_text.append(word_tokenize(string3))

**Creating Bigrams**

A bigram or digram is a sequence of two adjacent elements from a string of tokens, which are typically letters, syllables, or words. A bigram is an *n*-gram for *n*=2. The frequency distribution of every bigram in a string is commonly used for simple statistical analysis of text in many applications, including in computational linguistics, cryptography, speech recognition, and so on.

'''function to create ngrams from a list of words '''

def generate_ngrams(words_list,n):

      ngrams_list = []

      for num in range(0,len(words_list)):

      ngram = ' '.join(words_list[num:num + n])

      ngrams_list.append(ngram)

      return ngrams_list

'''Replacing data mining and machine learning'''

bigrams=[]

for i in range(0,5):

      bigrams.append(generate_ngrams(processed_text[i],2))

data=[]

'''Total N '''

for i in range(0,5):

      processed_text[i] = processed_text[i] + bigrams[i]

      data.extend(processed_text[i])

      #data.extend(bigrams[i])

N = len(data)

After creating bigrams we append that to the already existing unigram list of lists.

**Extra Credit: Inverted Index (Term Dictionary) Construction with TF-IDF**

TF-IDF = Term Frequency (TF) * Inverse Document Frequency (IDF)

**Creating tables for storing tf,df and tf_idf :**

With coulumns as DOC#,Term and Freq(for tf,tf_idf)

'''DataBase connection'''

# Connect

db = pymysql.connect(host="localhost", user="root",passwd="abc5s3",db="mydb")

cursor = db.cursor()

cursor.execute("DELETE FROM tf")

cursor.execute("DELETE FROM df")

cursor.execute("DELETE FROM tf_idf")


# Terminology

- t — term (word)
- d — document (set of words)
- N — count of corpus
- corpus — the total document set

 TF-IDF stands for **"Term Frequency — Inverse Document Frequency"**. This is a technique to quantify a word in documents, we generally compute a weight to each word which signifies the importance of the word in the document and corpus. This method is a widely used technique in Information Retrieval and Text Mining.


**Term Frequency**

This measures the frequency of a word in a document. This highly depends on the length of the document and the generality of word , for example a very common word such as "was" can appear multiple times in a document.

TF is individual to each document and word, hence we can formulate TF as follows.

tf(t,d) = count of t in d / number of words in d


**Document Frequency**

This measures the importance of document in whole set of corpus, this is very similar to TF. The only difference is that TF is frequency counter for a term t in document d, where as DF is the count of **occurrences** of term t in the document set N. In other words, DF is the number of documents in which the word is present. We consider one occurrence if the term consists

in the document at least once, we do not need to know the number of times the term is present.

df(t) = occurrence of t in documents

To keep this also in a range, we normalize by dividing with the total number of documents. Our main goal is to know the informativeness of a term, and DF is the exact inverse of it. that is why we inverse the DF

**Inverse Document Frequency**

IDF is the inverse of the document frequency which measures the informativeness of term t. When we calculate IDF, it will be very low for the most occurring words such as stop words (because stop words such as "is" is present in almost all of the documents, and N/df will give a very low value to that word). This finally gives what we want, a relative weightage.

idf(t) = N/df

During the query time, when a word which is not in vocab occurs, the df will be 0. As we cannot divide by 0, we smoothen the value by adding 1 to the denominator.

**tf-idf(t, d) = tf(t, d) * log(N/(df + 1))**

'''Calculating DF for all words'''

DF = {}

for j in range(0,5):

tokens = processed_text[j]

for w in tokens:

try:

DF[w].add(j)

except:

DF[w] = {j}

for i in DF:

DF[i] = len(DF[i])

total_vocab_size = len(DF)


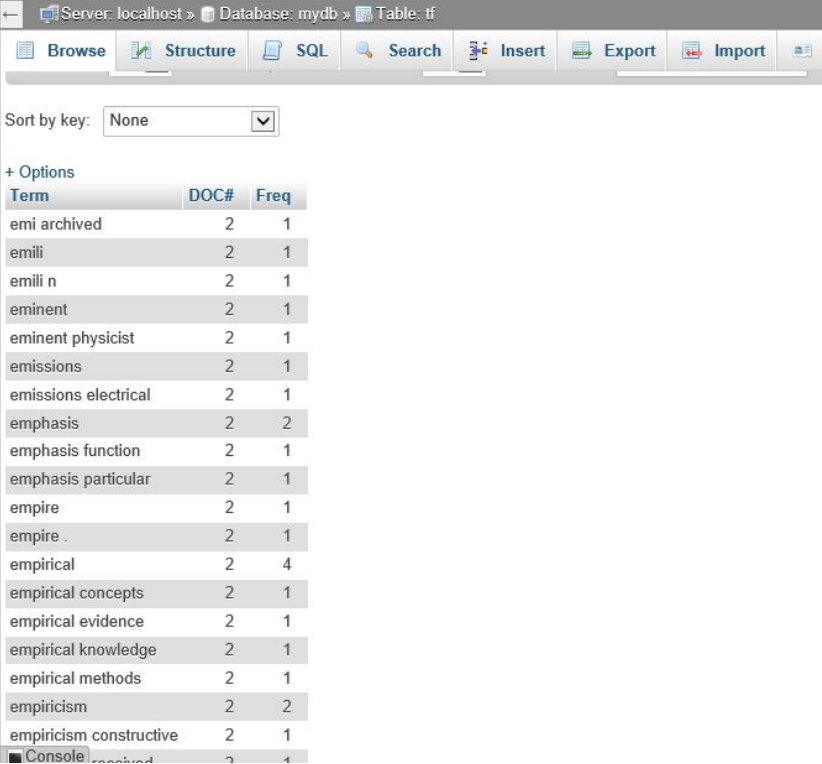total_vocab = [x for x in DF]


def doc_freq(word):

c =0

```python
        try:
        c = DF[word]
        except:
        pass
        return c
'''Calculating TF-IDF '''
doc = 1
tf_idf = {}
vtf = {}
for j in processed_text:
        tokens = j
        counter = Counter(tokens)
words_count = len(tokens)
        for token in np.unique(tokens):
        tf = counter[token]/words_count
        df = doc_freq(token)
        idf = np.log((N+1)/(df+1))
        L = [str(token),int(doc),int(counter[token])]
        cursor.execute("INSERT INTO tf VALUES(%s,%s,%s)",L)
        #vtf[doc,token]=counter[token]
        L1 = [str(token),doc_freq(token)]
        cursor.execute("INSERT INTO df VALUES(%s,%s)",L1)
        tf_idf[doc,token] = tf*idf
        L2=[int(doc),str(token),float(tf_idf[doc,token])]
        cursor.execute("INSERT INTO TF_IDF VALUES(%s,%s,%s)",L2)
        doc +=1
```

After storing tf,df,tf_idf for all the Terms

We have tf as



We have df as :

We have tf_idf as :

**Part 2: Data Transformation for Topic Analysis of Documents (Webpages)**

Building document vectors using if_idf because it already uses weight of each term in the document and will be in normalized way .

d_f=pd.DataFrame(index=[1,2,3,4,5],columns=['engineering','research','data','mi

ning','data mining','machine learning'])

for i in d_f.index:

       for j in d_f.columns:

       te = [str(j),i]

       cursor.execute("SELECT tfidf FROM tf_idf WHERE `TERM` =%s AND

`DOC#` =%s",te)

       _tfidf = cursor.fetchall()

       if len(_tfidf) > 0:

       d_f.at[i,j] = float(_tfidf[0][0])

       else :

       d_f.at[i,j] = 0

Finding cosine similarity of the vectors built.

print('------------------------TF-IDF----------------------------')

print(d_f)

cosinesimilarity = pd.DataFrame(index=[0,1,2,3,4],columns=[0,1,2,3,4])

a =cosine_similarity(d_f)

print('-----------------Cosine Similarity------------------')

cosinesimilarity=pd.DataFrame(data=a,index=['doc1','doc2','doc3','doc4','doc5'],

columns=['doc1','doc2','doc3','doc4','doc5'])

print(cosinesimilarity)


Final Output:

```
==== RESTART: C:\Users\BHUVAN\Desktop\study zone\DataMinig\Lab 2\lab2.py ====
------------------------TF-IDF-----------------------------
  engineering   research       data     mining data mining machine learning
1   0.0379362  0.0675626   0.151745 0.00689749  0.00353517         0.0353517
2    0.195894  0.0113659 0.00204768 0.00409535           0                 0
3           0   0.115418          0          0           0                 0
4   0.0053046   0.011258   0.236939   0.152949    0.132313         0.0290002
5   0.0053046   0.011258   0.236939   0.152949    0.132313         0.0290002
------------------Cosine Similarity------------------
           doc1       doc2       doc3       doc4       doc5
doc1   1.000000   0.249744   0.387880   0.723602   0.723602
doc2   0.249744   1.000000   0.057907   0.037078   0.037078
doc3   0.387880   0.057907   1.000000   0.035956   0.035956
doc4   0.723602   0.037078   0.035956   1.000000   1.000000
doc5   0.723602   0.037078   0.035956   1.000000   1.000000
```

Analysis and Discussion Discuss briefly about your topic analysis from your cosine similarity focusing on whether the indications by the values of your Cosine Sim are all correct?

The Topics of Doc1 is similar to the Topics of Doc 4 and 5? Explain Why or Why Not in terms of 6 TFs? If not, what are the reasons?

- The documents 4 and 5 are similar

Doc4: https://en.wikipedia.org/wiki/Data_mining

Doc5: https://en.wikipedia.org/wiki/Data_mining#Data_mining

So we can see the cosine similarity matrix that they both are similar.

- Document 1 is a website of prof. SS Chung which is related to data mining and machine learning because of her content in the webpage.And related document or close to document 4 and 5 because they are about data mining.
- Document 3 was about research and cleveland clinic which has nothing to do with data mining and machine learning so we can see the results in both the matrices .