



TEXT AND WEB INTELLIGENCE ANALYTICS

LAB MANUAL

SUBMITTED BY	Shruti Jain
ROLL NUMBER	17CSU186
CLASS	CSE-VIII-D
GROUP	C3
SESSION	2020-2021
FACULTY	DR. VAISHALI KALRA

Department of Computer Science and Engineering The
Northcap University, Sector-23, Gurugram, Haryana

INDEX

COURSE-CURRICULUM EXPERIMENTS				
S.NO.	Experiment	Date	Grade	Signature
1	Frequency Distribution	30-01-2021		
2	Stopwords	06/02/2021		
3	Exploring Names Corpus	13/02/2021		
4	Similarity metrices	27/02/2021		
5	Lesk Algorithm for Word Sense Disambiguation	06/03/2021		
6	Implement LDA with BOW and TF-IDF features and compare the results	01/04/2021		
7	Implementation of KNN, Naive Bayes and Multinomial Naive Bayes.	08/04/2021		
8	Implementation of K means, K Medoids and Hierarchical Clustering algorithms.	02/05/2021		
9	Implement homophily for social media analysis.	09/05/2021		
VALUE ADDED EXPERIMENTS				
1	Chunking, Chinking and Named Entity Recognition	03/03/2021		
2	Implementation of Fatman Evolutionary Network on jupyter notebook	09/05/2021		

EXPERIMENT NO. 1

Student Name and Roll Number: SHRUTI JAIN (17CSU186)
Semester /Section: VIII-D
Date: 30 January 2021
Link to Code: https://github.com/shruti17csu186/TWIA-LAB-EXPERIMENTS/blob/main/stopwords%20and%20synset%20implementation.ipynb
Faculty Signature:
Grade:

Objectives:

1. Import the corpus Shakespeare and find the frequency of each word in the file dream.xml.
2. Find 5 most frequently occurring words from the file dream.xml.
3. Import wordnet corpus from the available nltk corpus list and find out the sysnset of word bank. Also find the definition and example of first sysnset in the list.

Background Study:

nltk.FreqDist()

A frequency distribution records the number of times each outcome of an experiment has occurred. For example, a frequency distribution could be used to record the frequency of each word type in a document. Formally, a frequency distribution can be defined as a function mapping from each sample to the number of times that sample occurred as an outcome.

Wordnet

Wordnet is a lexical database of semantic relations between words in more than 200 languages. WordNet links words into semantic relations including synonyms, hyponyms, and meronym.

Outcome: Students will be able to learn the concepts of nltk.freqdist(), collections in python and sysnets in wordnet library.

Problem Statement:

1. Import the corpus Shakespeare and find the frequency of each word in the file dream.xml.

CODE AND OUTPUT:

```
In [16]: 1 #Import a file and find the frequency of each word in a file.  
2 import nltk  
3 import os  
4 import nltk.corpus  
5 from nltk.probability import FreqDist  
6 from nltk.corpus import brown  
7 text=(brown.words(categories='news'))  
8
```

```
In [22]: 1 #cleaning data of punctuations  
2 import string  
3 l=string.punctuation.split()  
4 no_punct_dream=[words for words in text if words not in string.punctuation]  
5 no_punct_dream  
6 #finding frequency  
7 fdist=nltk.FreqDist(w.lower() for w in no_punct_dream)  
8 fdist
```

```
Out[22]: FreqDist({'the': 6386, 'of': 2861, 'and': 2186, 'to': 2144, 'a': 2130, 'in': 2020, 'for': 969, 'that': 829, 'is': 733, "'": 732, ...})
```

```
In [8]: 1 #Find the five most frequently occurring words from the document  
2 fdist.most_common(5)
```

```
Out[8]: [('the', 5580), (',', 5188), ('.', 4030), ('of', 2849), ('and', 2146)]
```

2. Find 5 most frequently occurring words from the file dream.xml.

CODE AND OUTPUT:

Q2. Finding most frequently occurring words from the file dream.xml.

```
[ ] Dictionary=dict(fdist)
```

```
[ ] from collections import Counter  
dict(Counter(Dictionary).most_common(5))
```

```
{'and': 574, 'i': 470, 'the': 563, 'to': 340, 'you': 274}
```

- 3. Import wordnet corpus from the available nltk corpus list and find out the synset of word bank. Also find the definition and example of first synset in the list.**

CODE AND OUTPUT:

```
nltk.download('wordnet')
from nltk.corpus import wordnet
syns=wordnet.synsets("Bank")

[nltk_data]  Downloading package wordnet to /root/nltk_data...
[nltk_data]    Package wordnet is already up-to-date!
```

```
print("Definition of the word Bank:")
print(syns[0].definition())
print("\nExamples of the word Bank:")
print(syns[0].examples())
```

Definition of the word Bank:
sloping land (especially the slope beside a body of water)

Examples of the word Bank:
['they pulled the canoe up on the bank', 'he sat on the bank of the river and watched the currents']

EXPERIMENT NO. 2

Student Name and Roll Number: SHRUTI JAIN (17CSU186)

Semester /Section: VIII-D

Date: 6 February 2021

Link to Code:

<https://github.com/shruti17csu186/TWIA-LAB-IMENTS/blob/main/stopwords%20and%20synset%20implementation.ipynb>

Faculty Signature:

Grade:

Objective:

1. Print all the Arabic Stopwords.
2. Omit a given list of stop words from the total stopwords list of English language.

Background Study:

Stopwords

A stop word is a commonly used word (such as “the”, “a”, “an”, “in”) that a search engine has been programmed to ignore, both when indexing entries for searching and when retrieving them as the result of a search query. We would not want these words to take up space in our database or taking up valuable processing time. For this, we can remove them easily, by storing a list of words that are considered stopping words.

NLTK (Natural Language Toolkit) in python has a list of stopwords stored in 16 different languages

Outcome: Students will be able to understand the concept of stopwords in nltk and list comprehensions in python.

Problem Statement:

1. Print all the Arabic Stopwords. CODE AND OUTPUT:

```
In [10]: 1 #Print all the Arabic Stopwords.  
2 from nltk.corpus import stopwords  
3 stopwords=stopwords.words("Arabic")  
4 print(stopwords)
```

2. Omit a given list of stop words from the total stopwords list of English language

CODE AND OUTPUT:

```
In [13]: M
1 from nltk.corpus import stopwords
2 stopwords=stopwords.words("english")
3 l=[ "not", "nor", "very", "most", "shouldn't", "couldn't" ]
4 for i in l :
5     if i in stopwords:
6         stopwords.remove(i)
7 stopwords

"she's",
'her',
'hers',
'herself',
'it',
"it's",
'its',
'itself',
'they',
'them',
'their',
'theirs',
'themselves',
'what',
'which',
'who',
'whom',
'this',
'that',
```

EXPERIMENT NO. 3

Student Name and Roll Number: SHRUTI JAIN (17CSU186)		
Semester /Section: VIII-D		
Date: 13 February 2021		
Faculty Signature:		
Link	to	Code: https://github.com/shruti17csu186/TWIA-LAB-EXPERIMENTS/blob/main/stopwords%20and%20synset%20implementation.ipynb
Grade:		

Objectives:

1. Print the total number of male and female names in the names corpus. Then, Print the first 15 male and female names.
2. From the names corpus, combine all the labelled male and female names and print any 20.
3. Print the definition and examples of any one English language word using WordNet corpus.

Background Study:**Text Corpus**

A text corpus is a large and structured set of texts (nowadays usually electronically stored and processed). Text corpora are used to do statistical analysis and hypothesis testing, checking occurrences, or validating linguistic rules within a specific language territory.

Wordnet Corpus

Wordnet is a large lexical database of English. Nouns, verbs, adjectives and adverbs are grouped into sets of cognitive synonyms (synsets), each expressing a distinct concept. Synsets are interlinked by means of conceptual-semantic and lexical relations.

WordNet superficially resembles a thesaurus, in that it groups words together based on their meanings

Outcome: Students will be able to explore names corpus, understand the concept of labelling the data and learn about synsets in Wordnet corpus.

Problem Statement:

- 1. Print the total number of male and female names in the names corpus. Then, Print the first 15 male and female names.**

CODE AND OUTPUT:

```
[ ] nltk.download('names')
from nltk.corpus import names
print("\nNumber of male names:")
print(len(names.words('male.txt')))
print("Number of female names:")
print(len(names.words('female.txt')))

[nltk_data] Downloading package names to /root/nltk_data...
[nltk_data]   Package names is already up-to-date!

Number of male names:
2943
Number of female names:
5001
```

- 2. From the names corpus, combine all the labelled male and female names and print any 20.**

CODE AND OUTPUT:

```
[ ] Male_names = names.words('male.txt')
Female_names = names.words('female.txt')
print("\nFirst 15 male names:")
print(male_names[0:15])
print("\nFirst 15 female names:")
print(female_names[0:15])
```

First 15 male names:

['Aamir', 'Aaron', 'Abbey', 'Abbie', 'Abbot', 'Abbott', 'Abby', 'Abdel', 'Abdul', 'Abdulkarim', 'Abdullah', 'Abe', 'Abel', 'Abelard', 'Abner']

First 15 female names:

['Abagael', 'Abagail', 'Abbe', 'Abbey', 'Abbi', 'Abbie', 'Abby', 'Abigael', 'Abigail', 'Abigale', 'Abra', 'Acacia', 'Ada', 'Adah', 'Adaline']

3. Print the definition and examples of any one English language word using WordNet corpus.

CODE AND OUTPUT:

```
nltk.download('wordnet')
from nltk.corpus import wordnet
syns=wordnet.synsets("Telephone")
```

[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Package wordnet is already up-to-date!

```
print("Definition of the word Telephone:")
print(syns[0].definition())
print("\nExamples of the word Telephone:")
print(syns[0].examples())
```

Definition of the word Telephone:
electronic equipment that converts sound into electrical signals that can be transmitted over distances and then converts received signals back into sounds

Examples of the word Telephone:
['I talked to him on the telephone']

2020-21

EXPERIMENT NO. 4

Student Name and Roll Number: SHRUTI JAIN (17CSU186)
Semester /Section: VIII-D
Date: 27 February 2021
Link to Code: https://github.com/shruti17csu186/TWIA-LAB-EXPERIMENTS/blob/main/similarity%20matrices%20and%20tfidf%20implementation.ipynb
Faculty Signature:
Grade:

Objective: To implement Levenshtein Edit Distance , Jaccard similarity , Cosine Similarity using both TF-IDF and count vectorizer

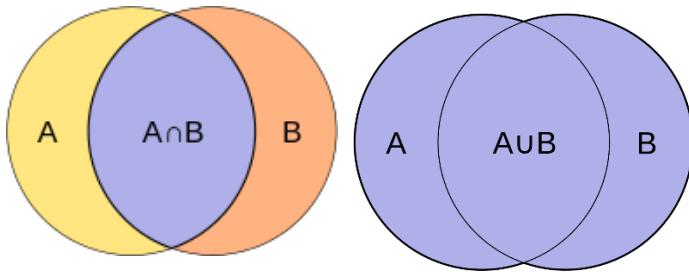
Background Study:

Levenshtein distance

The Levenshtein distance is a string metric for measuring difference between two sequences. Informally, the Levenshtein distance between two words is the minimum number of single-character edits (i.e. insertions, deletions or substitutions) required to change one word into the other.

Jaccard similarity

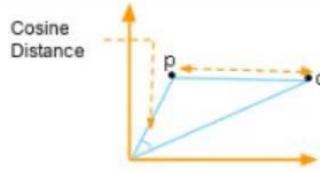
The Jaccard index, also known as the Jaccard similarity coefficient, is a statistic used for gauging the similarity and diversity of sample sets. The Jaccard coefficient measures similarity between finite sample sets, and is defined as the size of the intersection divided by the size of the union of the sample sets.



Cosine similarity

Cosine similarity measures the similarity between two vectors of an inner product space. It is measured by the cosine of the angle between two vectors and determines whether two vectors are pointing in roughly the same direction. It is often used to measure document similarity in text analysis.

$$d = \frac{\sum_{i=0}^{n-1} q_i - p_r}{\sum_{i=0}^{n-1} (q_i)^2 \times \sum_{i=0}^{n-1} (p_i)^2}$$



TF-IDF Vectorizer

TF-IDF is a statistical measure that evaluates how relevant a word is to a document in a collection of documents. This is done by multiplying two metrics: how many times a word appears in a document, and the inverse document frequency of the word across a set of documents.

TF-IDF for a word in a document is calculated by multiplying two different metrics:

- The term frequency of a word in a document. There are several ways of calculating this frequency, with the simplest being a raw count of instances a word appears in a document. Then, there are ways to adjust the frequency, by length of a document, or by the raw frequency of the most frequent word in a document.
- The inverse document frequency of the word across a set of documents. This means, how common or rare a word is in the entire document set. The closer it is to 0, the more common a word is. This metric can be calculated by taking the total number of documents, dividing it by the number of documents that contain a word, and calculating the logarithm.
- So, if the word is very common and appears in many documents, this number will approach 0. Otherwise, it will approach 1.

CountVectorizer

CountVectorizer is used to convert a collection of text documents to a vector of term/token counts. It also enables the pre-processing of text data prior to generating the vector representation. This functionality makes it a highly flexible feature representation module for text.

```
Data = ['The', 'quick', 'brown', 'fox', 'jumps', 'over', 'the', 'lazy', 'dog']
```



	The	quick	brown	fox	jumps	over	lazy	dog
Data	2	1	1	1	1	1	1	1

Outcome: Students will be able to demonstrate Levenshtein Edit Distance , Jaccard similarity , Cosine Similarity using both TF-IDF and count vectorizer

Problem Statement:

1. Demonstrate the computation of Similarity Metrics such as Jaccard, Levenshtein and Cosine.
CODE AND OUTPUT :

```
#Jaccard Similarity - 

In [1]: 
 1 import nltk
 2 from nltk.metrics import *
 3 doc_1 = "Data is the new oil of the digital economy"
 4 doc_2 = "Data is a new oil"
 5 words_doc1=doc_1.lower().split()
 6 words_doc1
 7 words_doc2=doc_2.lower().split()
 8 words_doc2
 9 set_doc1=set(words_doc1)
10 set_doc2=set(words_doc2)
11 print(set_doc1)
12 intersection = set_doc1.intersection(set_doc2)
13 union = set_doc1.union(set_doc2)
14 jaccard=float(len(intersection)) / len(union)
15 jaccard

{'of', 'data', 'is', 'oil', 'digital', 'the', 'new', 'economy'}

Out[1]: 0.4444444444444444
```

LEVENSHTEIN DISTANCE

```
#LEVENSHTEIN DISTANCE

In [8]: 
 1 #LEVENSHTEIN DISTANCE
 2
 3 import enchant
 4 data1= "shruti jain is a good girl."
 5 data2 = "shruti jain"
 6 enchant.utils.levenshtein(data1,data2)

Out[8]: 16
```

COSINE SIMILARITY

```
#COSINE SIMILARITY

In [9]: 
 1 #COSINE SIMILARITY
 2 import nltk
 3 nltk.download( ' punkt ' )
 4 nltk.download('stopwords')
 5
 6 from nltk.corpus import stopwords
 7 from nltk.tokenize import word_tokenize
 8
 9
10 I = input("Enter first string: ").lower()
11 II = input("Enter second string: ").lower()
12

[nltk_data] Error loading punkt : Package 'punkt' not found in
[nltk_data]     index
[nltk_data] Downloading package stopwords to
[nltk_data]     C:\Users\Shruti\AppData\Roaming\nltk_data...
[nltk_data]     Package stopwords is already up-to-date!

Enter first string: shruti is a good girl
Enter second string: shruti studies at ncu
```

```
In [11]: 1 # tokenization
2 I_list = word_tokenize(I)
3 II_list = word_tokenize(II)
4 print('First List', I_list)
5 print('Second Lzst', II_list)
6
7
```

First List ['shruti', 'is', 'a', 'good', 'girl']
 Second Lzst ['shruti', 'studies', 'at', 'ncu']

```
In [14]: 1 #removing stopwords
2
3 sw = stopwords.words('english')
4 l1 = []
5 l2 = []
6 I_set = {w for w in I_list if not w in sw}
7 II_set = {w for w in II_list if not w in sw}
8
9 print (' First Set ',I_set)
10 print (' second Set',II_set )
11
12
13
```

First Set {'shruti', 'good', 'girl'}
 second Set {'shruti', 'ncu', 'studies'}

```
In [24]: 1 # Forming a set containing keywords of both strings
2
3 rvector = I_set.union(II_set)
4 for w in rvector:
5     if w in I_set:
6         l1.append(1)
7     else:
8         l1.append(0)
9     if w in II_set:
10        l2.append(1)
11    else:
12        l2.append(0)
13 c = 0
14
```

COSINE SIMILARITY MATRIX

```
In [35]: 1 #Cosine similarity matrix
2 I = input("Enter first string: ").lower()
3 II = input("Enter second string: ").lower()
4
```

Enter first string: shruti jain is the student of text and web analytics
 Enter second string: shruti jain is the student of northcap university

```
In [36]: 1 documents = [I,II]
2 from sklearn.feature_extraction.text import CountVectorizer
3 import pandas as pd
4 count_vectorizer = CountVectorizer(stop_words='english')
5 count_vectorizer = CountVectorizer()
6 sparse_matrix = count_vectorizer.fit_transform(documents)
7 doc_term_matrix = sparse_matrix.todense()
8 df = pd.DataFrame(doc_term_matrix,columns=count_vectorizer.get_feature_names(),index=['I', 'II'])
9 df
10
```

Out[36]:

analytics and is jain northcap of shruti student text the university web

I	1	1	1	1	0	1	1	1	1	0	1
II	0	0	1	1	1	1	1	1	0	1	0

2. Calculate the TF-IDF vectorizer on 2 documents.

CODE AND OUTPUT:

```
In [37]: 1 #2.Calculate the TF-IDF vectorizer on 2 documents.  
2 #TF-IDF  
3 import pandas as pd  
4 from sklearn.feature_extraction.text import CountVectorizer  
5 from sklearn.feature_extraction.text import TfidfVectorizer  
6 text = ["Do not limit your challenges,challenge your limits",  
7         "your challenges",  
8         "their limits"]  
9 vect = TfidfVectorizer()  
10 X = vect.fit_transform(text)  
11 print(vect.get_feature_names())  
12 vect.idf_  
  
['challenge', 'challenges', 'do', 'limit', 'limits', 'not', 'their', 'your']  
  
Out[37]: array([1.69314718, 1.28768207, 1.69314718, 1.69314718, 1.28768207,  
1.69314718, 1.69314718, 1.28768207])
```

3. Apply the max-df, min-df param in the TF-IDF function.

CODE AND OUTPUT:

```
In [42]: 1 #3.Apply the max-df, min-df param in the TF-IDF function.  
2 data=[I,II]  
3 from sklearn.feature_extraction.text import CountVectorizer  
4 import pandas as pd  
5 count_vec = CountVectorizer(stop_words='english', analyzer="word", ngram_range=(1,1), max_df=0.50, min_df=1, max_features=None)  
6 count_train= count_vec.fit(data)  
7 bag_of_words=count_vec.transform(data)  
8 print("features:",count_vec.get_feature_names())  
9  
features: ['analytics', 'northcap', 'text', 'university', 'web']
```

4. Compute Cosine Similarity using both TF-IDF and Count vectorizer CODE

AND OUTPUT:

```
In [50]: 1 #4.Compute Cosine Similarity using both TF-IDF and Count vectorizer  
2 from sklearn.feature_extraction.text import CountVectorizer  
3 from sklearn.feature_extraction.text import TfidfTransformer  
4 from nltk.corpus import stopwords  
5 import numpy as np  
6 import numpy.linalg as LA  
7 import nltk  
8 nltk.download("stopwords")  
  
[nltk_data] Downloading package stopwords to  
[nltk_data]     C:\Users\Shruti\AppData\Roaming\nltk_data...  
[nltk_data]   Package stopwords is already up-to-date!  
  
Out[50]: True
```

```
In [51]: 1 train_dataset=["shruti is a good girl","shruti studies at northcap university"]#document  
2 test_dataset=["a good girl named shruti studies at nothcap university"]#query  
3 stopWords=stopwords.words("English")
```

```
In [56]: 1 vectorizer = CountVectorizer(stop_words = stopWords)
2 print('Vectorizer : ',vectorizer)
3 transformer = TfidfTransformer()
4 print('\n\nTF-IDF Transformer : ',transformer)
5 #Using Count Vectorizer
6 trainVectorizerArray = vectorizer.fit_transform(train_dataset).toarray()
7 testVectorizerArray =vectorizer.transform(test_dataset).toarray()
8 print ('Fit Vectorizer to train set \n', trainVectorizerArray)
9 print ('\n\nTransform Vectorizer to test set\n', testVectorizerArray)
10
```

Vectorizer : CountVectorizer(analyzer='word', binary=False, decode_error='strict',
 dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
 lowercase=True, max_df=1.0, max_features=None, min_df=1,
 ngram_range=(1, 1), preprocessor=None,
 stop_words=['i', 'me', 'my', 'myself', 'we', 'our', 'ours',
 'ourselves', 'you', "you're", "you've", "you'll",
 "you'd", 'your', 'yours', 'yourself', 'yourselves',
 'he', 'him', 'his', 'himself', 'she', "she's",
 'her', 'hers', 'herself', 'it', "it's", 'its',
 'itself', ...],
 strip_accents=None, token_pattern='(?u)\\b\\w+\\b',
 tokenizer=None, vocabulary=None)

TF-IDF Transformer : TfidfTransformer(norm='l2', smooth_idf=True, sublinear_tf=False, use_idf=True)
Fit Vectorizer to train set
TF-IDF Transformer : TfidfTransformer(norm='l2', smooth_idf=True, sublinear_tf=False, use_idf=True)
Fit Vectorizer to train set
[[1 1 0 1 0 0]
 [0 0 1 1 1 1]]

Transform Vectorizer to test set
[[1 1 0 1 1 1]]

```
In [57]: 1 #Using
2 transformer.fit(trainVectorizerArray)
3 print('Fit transformer to train set \n',transformer.transform(trainVectorizerArray).toarray())
4 transformer.fit(testVectorizerArray)
5 tfidf = transformer.transform(testVectorizerArray)
6 print('\n\nFit transformer to test set\n',tfidf.todense())
7
```

Fit transformer to train set
[[0.6316672 0.6316672 0. 0.44943642 0. 0.]
 [0. 0. 0.53404633 0.37997836 0.53404633 0.53404633]]

Fit transformer to test set
[[0.4472136 0.4472136 0. 0.4472136 0.4472136 0.4472136]]

EXPERIMENT NO. 5

Student Name and Roll Number: SHRUTI JAIN (17CSU186)
Semester /Section: VIII-D
Date: 6 March 2021
Faculty Signature:
Link to Code: https://github.com/shruti17csu186/TWIA-LAB-EXPERIMENTS/blob/main/Lesk%20Implementation.ipynb
Grade:

Objective: Implementation of the Lesk algorithm for Word Sense Disambiguation

Background Study:

Lesk Algorithm

The Lesk algorithm assumes that words in each "neighbourhood" (section of text) will tend to share a common topic. A simplified version of the Lesk algorithm is to compare the dictionary definition of an ambiguous word with the terms contained in its neighbourhood.

An implementation might look like this:

1. for every sense of the word being disambiguated one should count the number of words that are in both neighbourhood of that word and in the dictionary definition of that sense

the sense that is to be chosen is the sense which has the biggest number of this count.

Outcome: Students will be able to demonstrate how to Lesk

Problem Statement: Implement Lesk algorithm for Word Sense Disambiguation

CODE AND OUTPUT:

```
▶ import nltk
nltk.download('averaged_perceptron_tagger')
nltk.download('wordnet')
nltk.download('punkt')

▷ [nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]      /root/nltk_data...
[nltk_data]      Package averaged_perceptron_tagger is already up-to-
[nltk_data]      date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]      Package wordnet is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]      Unzipping tokenizers/punkt.zip.
True
```

```
[12] from pywsd.lesk import simple_lesk
```

```
[16] sentences = ['I went to the bank to deposit my money','The river bank had a lot of fishes and crocodiles.']

#ambiguous word - Bank
```

· LESK WORKS CORRECTLY

```
[13] # Context 1 - Financial institution

print ("Context-1:", sentences[0])
answer = simple_lesk(sentences[0], 'bank')
print ("Sense:", answer)
print ("Definition : ", answer.definition())

# Correct Output - Financial Institution printed
# No disambiguity
```

```
Context-1: I went to the bank to deposit my money
Sense: Synset('repository_financial_institution.n.01')
Definition : a financial institution that accepts deposits and channels the money into lending activities
```

```
▶ # Context 2 - River Bank

print ("Context-2:", sentences[1])
answer = simple_lesk(sentences[1], 'bank')
print ("Sense:", answer)
print ("Definition : ", answer.definition())

# Correct Output - River Bank or sloping land printed
# No disambiguity
```

```
Context-2: The river bank had a lot of fishes and crocodiles.
Sense: Synset('bank.n.01')
Definition : sloping land (especially the slope beside a body of water)
```

LESK WORKS INCORRECTLY

```
[17] new_sentences = ['The workers at the plant were overworked', 'The plant was no longer bearing flowers', 'The workers at the industrial plant were overworked']

#ambiguous word - Plant
```

```
[18] # Context 1 - Industrial plant

print ("Context-1:", new_sentences[0])
answer = simple_lesk(new_sentences[0], 'plant')
print ("Sense:", answer)
print ("Definition : ", answer.definition())

# Incorrect output - Industrial plant not printed
# Disambiguity occurred
```

Context-1: The workers at the plant were overworked
Sense: Synset('plant.v.06')
Definition : put firmly in the mind

```
[19] # Context 2 - Tree/Seedling/Sapling
```

```
print ("Context-2:", new_sentences[1])
answer = simple_lesk(new_sentences[1], 'plant')
print ("Sense:", answer)
print ("Definition : ", answer.definition())

# Correct output - Plant bearing flower sense printed
# No disambiguation
```

Context-2: The plant was no longer bearing flowers
Sense: Synset('plant.v.01')
Definition : put or set (seeds, seedlings, or plants) into the ground

▶ # Context 3 - Industrial plant (Added the word industrial before plant in
Context 1)

print ("Context-3:", new_sentences[2])
answer = simple_lesk(new_sentences[2], 'plant')
print ("Sense:", answer)
print ("Definition : ", answer.definition())

Correct output - Industrial plant printed
(Disambiguation resolved by adding the word "industrial" in the sentence.)

⇒ Context-3: The workers at the industrial plant were overworked
Sense: Synset('plant.n.01')
Definition : buildings for carrying on industrial labor

EXPERIMENT NO. 6

Student Name and Roll Number: SHRUTI JAIN (17CSU186)
Semester /Section: VIII-D
Date: 1 April 2021
Faculty Signature:
Link to Code: https://github.com/shruti17csu186/TWIA-LAB-EXPERIMENTS/blob/main/LDA%20implementation%20and%20comparision.ipynb
Grade:

Objectives:

1. Implement LDA with Bag of Words
2. Implement LDA with TF-IDF
3. Compare the accuracy using Score values make the analysis.

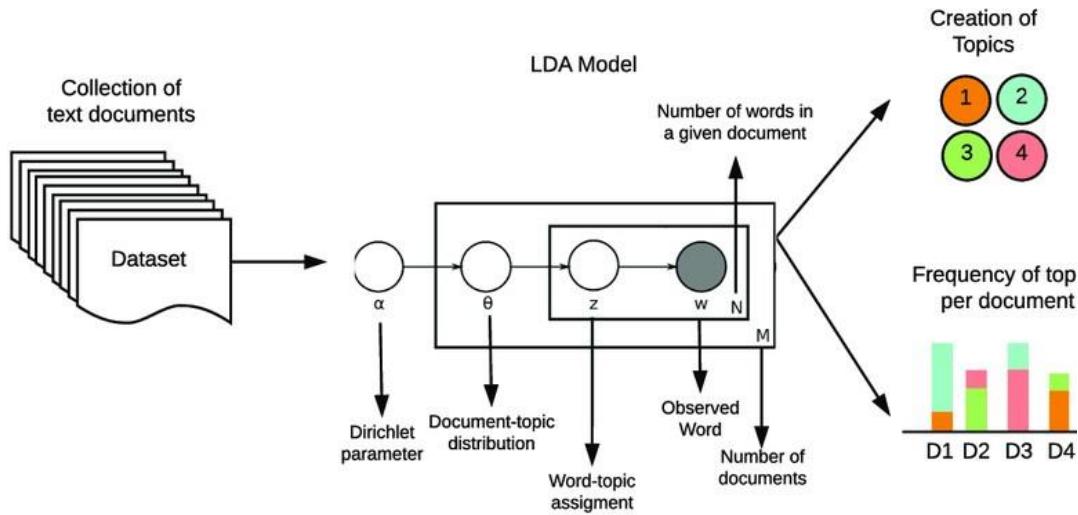
Dataset: <https://www.kaggle.com/therohk/million-headlines>

Background Study:

Latent Dirichlet Allocation (LDA)

It is an example of topic model and is used to classify text in a document to a particular topic. It builds a topic per document model and words per topic model, modeled as Dirichlet distributions.

For example, if observations are words collected into documents, it posits that each document is a mixture of a small number of topics and that each word's presence is attributable to one of the document's topics. LDA is an example of a topic model and belongs to the machine learning toolbox and in wider sense to the artificial intelligence toolbox.



Outcome: Students will be able to demonstrate how LDA works using Bag of Words and using TF-IDF.

Problem Statement: Implement LDA using BOW and TF-IDF on Million News Dataset

CODE AND OUTPUT:

Dataset

```
In [1]: 1 import pandas as pd
2
3 data = pd.read_csv('abcnews-date-text.csv', error_bad_lines=False);
4 data_text = data[['headline_text']]
5 data_text['index'] = data_text.index
6 documents = data_text
```

Pre-Processing

```
In [13]: 1 import gensim
2 from gensim.utils import simple_preprocess
3 from gensim.parsing.preprocessing import STOPWORDS
4 from nltk.stem import WordNetLemmatizer
5 from nltk.stem.porter import *
6 import numpy as np
7 np.random.seed(2018)
```

```
In [14]: 1 import nltk
2 nltk.download('wordnet')

[nltk_data] Downloading package wordnet to
[nltk_data]     C:\Users\Shruti\AppData\Roaming\nltk_data...
[nltk_data]     Package wordnet is already up-to-date!
```

```
Out[14]: True
```

```
In [15]: 1 # not stemming as it will not provide valid results
2 #Lemmatizing
3 #removing stopwords and words with Len<3
4 def lemmatize_stemming(text):
5     return WordNetLemmatizer().lemmatize(text, pos='v')
6
7 def preprocess(text):
8     result = []
9     for token in gensim.utils.simple_preprocess(text):
10         if token not in gensim.parsing.preprocessing.STOPWORDS and len(token) > 3:
11             result.append(lemmatize_stemming(token))
12
13 return result
```

```
In [16]: 1 processed_docs = documents['headline_text'].map(preprocess)
2 processed_docs[:10]
```

```
Out[16]: 0      [decide, community, broadcast, licence]
1          [witness, aware, defamation]
2      [call, infrastructure, protection, summit]
3          [staff, aust, strike, rise]
4      [strike, affect, australian, travellers]
5          [ambitious, olsson, win, triple, jump]
6          [antic, delight, record, break, barca]
7      [aussie, qualifier, stosur, waste, memphis, ma...]
8          [aust, address, security, council, iraq]
9          [australia, lock, timetable]
Name: headline_text, dtype: object
```

```
In [22]: 1 dictionary = gensim.corpora.Dictionary(processed_docs)
```

```
In [49]: 1 dictionary.filter_extremes(no_below=15, no_above=0.5, keep_n=100000)
2 bow_corpus = [dictionary.doc2bow(doc) for doc in processed_docs]
```

```
In [35]: 1 from gensim import corpora, models
2
3 tfidf = models.TfidfModel(bow_corpus)
4 corpus_tfidf = tfidf[bow_corpus]
5 from pprint import pprint
6 for doc in corpus_tfidf:
7     pprint(doc)
8     break
```

```
[(0, 0.5918674193999763),
(1, 0.3937180767686992),
(2, 0.5009876624450964),
(3, 0.49365007440105513)]
```

```
In [36]: 1 # training Lda model
2 lda_model = gensim.models.LdaMulticore(bow_corpus, num_topics=10, id2word=dictionary, passes=2, workers=2)
```

```
In [40]: 1 for idx, topic in lda_model.print_topics(-1):
2     print('Topic: {} \nWords: {}'.format(idx, topic))
```

```
Topic: 0
Words: 0.057*"australia" + 0.041*"trump" + 0.024*"australian" + 0.022*"china" + 0.019*"world" + 0.019*"sydney" + 0.017*"open"
+ 0.017*"coronavirus" + 0.015*"border" + 0.012*"win"
Topic: 1
Words: 0.023*"market" + 0.019*"year" + 0.016*"record" + 0.012*"care" + 0.012*"price" + 0.012*"years" + 0.012*"australian" +
0.011*"business" + 0.011*"country" + 0.010*"age"
Topic: 2
Words: 0.065*"coronavirus" + 0.032*"covid" + 0.029*"government" + 0.015*"rise" + 0.015*"restrictions" + 0.014*"water" + 0.01
2*"royal" + 0.012*"scott" + 0.011*"tasmanian" + 0.010*"commission"
Topic: 3
Words: 0.027*"kill" + 0.022*"die" + 0.019*"coast" + 0.018*"shoot" + 0.017*"miss" + 0.016*"crash" + 0.015*"attack" + 0.015*
"old" + 0.015*"dead" + 0.014*"island"
Topic: 4
Words: 0.040*"police" + 0.026*"charge" + 0.026*"case" + 0.025*"court" + 0.020*"death" + 0.020*"murder" + 0.017*"face" + 0.01
3*"jail" + 0.013*"people" + 0.012*"arrest"
Topic: 5
Words: 0.027*"test" + 0.020*"tasmania" + 0.015*"morrison" + 0.014*"drum" + 0.012*"premier" + 0.011*"save" + 0.010*"race" +
0.009*"show" + 0.009*"park" + 0.009*"alan"
Topic: 6
Words: 0.050*"say" + 0.024*"victoria" + 0.022*"health" + 0.019*"change" + 0.019*"adelaide" + 0.014*"indigenous" + 0.011*"rur
al" + 0.011*"service" + 0.010*"national" + 0.010*"climate"
Topic: 7
Words: 0.030*"queensland" + 0.028*"election" + 0.022*"live" + 0.013*"federal" + 0.013*"work" + 0.012*"school" + 0.011*"coun
cils" + 0.010*"farm" + 0.009*"labour" + 0.009*"fund"
```

```
In [42]: 1 for idx, topic in lda_model_tfidf.print_topics(-1):
2     print('Topic: {} Word: {}'.format(idx, topic))
```

```
Topic: 0 Word: 0.007*"coronavirus" + 0.007*"border" + 0.007*"kill" + 0.006*"september" + 0.005*"footage" + 0.005*"china" +
0.005*"brexit" + 0.005*"protesters" + 0.005*"biden" + 0.005*"attack"
Topic: 1 Word: 0.013*"crash" + 0.011*"queensland" + 0.010*"bushfire" + 0.010*"coast" + 0.007*"weather" + 0.007*"die" + 0.007
*"michael" + 0.007*"police" + 0.006*"beach" + 0.006*"death"
Topic: 2 Word: 0.006*"plan" + 0.005*"coal" + 0.005*"water" + 0.005*"action" + 0.005*"government" + 0.005*"legal" + 0.005*"co
uncil" + 0.004*"spring" + 0.004*"closure" + 0.004*"coronavirus"
Topic: 3 Word: 0.021*"coronavirus" + 0.016*"covid" + 0.010*"health" + 0.009*"government" + 0.007*"royal" + 0.007*"commis
sion" + 0.007*"federal" + 0.006*"care" + 0.006*"update" + 0.006*"sport"
Topic: 4 Word: 0.012*"scott" + 0.011*"morrison" + 0.010*"restrictions" + 0.010*"coronavirus" + 0.009*"christmas" + 0.008*"tu
rnbull" + 0.007*"island" + 0.007*"stories" + 0.006*"morning" + 0.005*"malcolm"
Topic: 5 Word: 0.025*"news" + 0.016*"rural" + 0.010*"climate" + 0.010*"business" + 0.009*"thursday" + 0.009*"john" + 0.009
*"victorian" + 0.008*"david" + 0.008*"national" + 0.006*"farm"
Topic: 6 Word: 0.013*"market" + 0.009*"australian" + 0.009*"friday" + 0.008*"price" + 0.008*"share" + 0.007*"coronavirus" +
0.007*"rise" + 0.006*"grandstand" + 0.006*"finance" + 0.005*"festival"
Topic: 7 Word: 0.030*"trump" + 0.018*"donald" + 0.017*"interview" + 0.015*"country" + 0.010*"speak" + 0.010*"hour" + 0.007
*"australians" + 0.007*"cricket" + 0.007*"peter" + 0.007*"alan"
Topic: 8 Word: 0.014*"drum" + 0.011*"australia" + 0.010*"world" + 0.009*"tuesday" + 0.009*"league" + 0.008*"wednesday" + 0.0
08*"street" + 0.008*"wall" + 0.007*"final" + 0.006*"brief"
Topic: 9 Word: 0.018*"charge" + 0.016*"murder" + 0.016*"police" + 0.012*"court" + 0.011*"woman" + 0.010*"jail" + 0.010*"arre
st" + 0.009*"sentence" + 0.009*"assault" + 0.009*"guilty"
```

Testing both the models

```
In [51]: # Bag Of Words
# Compute Perplexity
from gensim.models.coherencemodel import CoherenceModel
print('\nPerplexity: ', lda_model.log_perplexity(bow_corpus)) # a measure of how good the model is. Lower the better.
#
# Compute Coherence Score
coherence_model_lda = CoherenceModel(model=lda_model, texts=processed_docs, dictionary=dictionary, coherence='c_v')
coherence_lda = coherence_model_lda.get_coherence()
print('\nCoherence Score: ', coherence_lda)
```

Perplexity: -9.118709657723862

Coherence Score: 0.24535051024041796

```
In [54]: # TFIDF
# Compute Perplexity
from gensim.models.coherencemodel import CoherenceModel
print('\nPerplexity: ', lda_model_tfidf.log_perplexity(bow_corpus)) # a measure of how good the model is. Lower the bet:
#
# Compute Coherence Score
coherence_model_lda = CoherenceModel(model=lda_model_tfidf, texts=processed_docs, dictionary=dictionary, coherence='c_v')
coherence_lda = coherence_model_lda.get_coherence()
print('\nCoherence Score: ', coherence_lda)
```

Perplexity: -8.997436855159867

Coherence Score: 0.30983791499722557

```
In [ ]: #since tfidif has more coherence score therefore it is more effective than bow
```

Conclusion:

Since Tf-Idf approach has more coherence score than Bag of Words approach therefore it more effective in text classification

EXPERIMENT NO. 7

Student Name and Roll Number: SHRUTI JAIN (17CSU186)
Semester /Section: VIII-D
Date: 8 April 2021
Faculty Signature:
Link to Code: https://github.com/shruti17csu186/TWIA-LAB-EXPERIMENTS/blob/main/knn%20na%C3%A9ve%20Bayes(Gaussian%20and%20multinomial).ipynb
Grade:

Objectives:

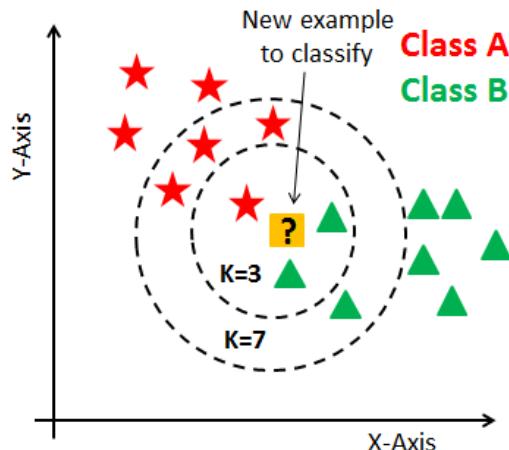
1. Perform KNN, Naïve Bayes and Multinomial Naïve Bayes on 20 Newsgroup Dataset.

Dataset: <https://www.kaggle.com/crawford/20-newsgroups>

Background Study:

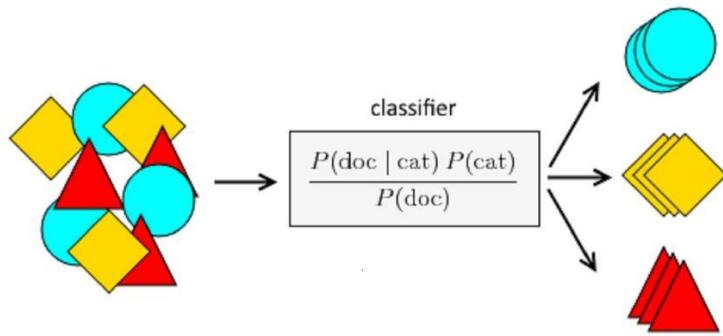
K- Nearest Neighbors (KNN)

The k-nearest neighbors (KNN) algorithm is a simple, easy-to-implement supervised machine learning algorithm that can be used to solve both classification and regression problems. The KNN algorithm assumes that similar things exist in close proximity. In other words, similar things are near to each other.



Naïve Bayes

Naïve Bayes classifiers are a collection of classification algorithms based on Bayes Theorem. It is not a single algorithm but a family of algorithms where all of them share a common principle, i.e., every pair of features being classified is independent of each other.



Multinomial Naïve Bayes

The multinomial Naive Bayes classifier is suitable for classification with discrete features (e.g., word counts for text classification). The multinomial distribution normally requires integer feature counts. However, in practice, fractional counts such as tf-idf may also work.

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

$$P(c | X) = P(x_1 | c) \times P(x_2 | c) \times \dots \times P(x_n | c) \times P(c)$$

Outcome: Students will be able to understand the difference between the implementations of KNN, Naïve Bayes and Multinomial Naïve Bayes.

Problem Statement: Implement KNN, Naïve Bayes and Multinomial Naïve Bayes on 20 Newsgroup dataset and compare their accuracy scores.

CODE AND OUTPUT:

Importing libraries

```
In [1]: 1 import numpy as np
2 from sklearn.datasets import fetch_20newsgroups
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.feature_extraction.text import TfidfTransformer
5 from sklearn.neighbors import KNeighborsClassifier
6 from sklearn.pipeline import Pipeline
7 from sklearn.naive_bayes import GaussianNB
8 from sklearn.naive_bayes import MultinomialNB
9 import pandas as pd
10 from sklearn import metrics
11 import seaborn as sns; sns.set()
```

importing Data

```
In [2]: categories = ['rec.motorcycles', 'sci.electronics',
                   'comp.graphics', 'sci.med']

# sklearn provides us with subset data for training and testing
train_data = fetch_20newsgroups(subset='train',
                                categories=categories, shuffle=True, random_state=42)

test_data = fetch_20newsgroups(subset='test',
                               categories=categories, shuffle=True, random_state=42)
```

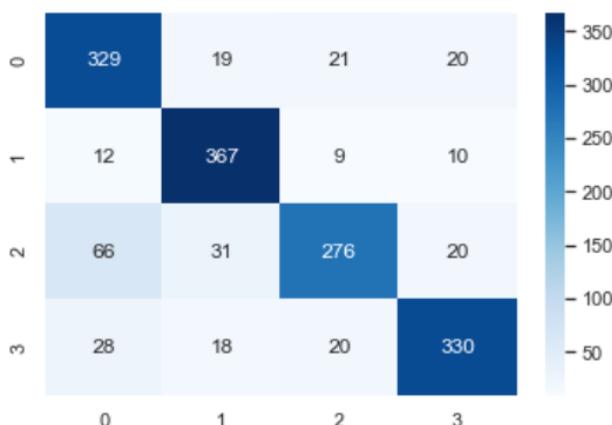
KNN Classifier

```
In [3]: knn = KNeighborsClassifier(n_neighbors=3)
#pipelining the fuctions
pipeline = Pipeline([
    ('vect', CountVectorizer()),
    ('tfidf', TfidfTransformer()),
    ('clf', knn)
])
pipeline.fit(train_data.data, train_data.target)
predicted = pipeline.predict(test_data.data)
```

```
In [4]: # accuracy of the above
cm=metrics.confusion_matrix(test_data.target,predicted)
ax = sns.heatmap(cm, annot=True, fmt="d", cmap="Blues")
print(metrics.classification_report(test_data.target,predicted))
print(metrics.accuracy_score(test_data.target,predicted))
```

	precision	recall	f1-score	support
0	0.76	0.85	0.80	389
1	0.84	0.92	0.88	398
2	0.85	0.70	0.77	393
3	0.87	0.83	0.85	396
accuracy			0.83	1576
macro avg	0.83	0.83	0.82	1576
weighted avg	0.83	0.83	0.82	1576

0.8261421319796954



MultiNomial Naive Bayes Classifier

```
In [5]: 1 mnb = MultinomialNB()
2 #pipelinnig the fuctions
3 pipeline = Pipeline([
4     ('vect', CountVectorizer()),
5     ('tfidf', TfidfTransformer()),
6     ('clf', mnb)
7 ])
8 pipeline.fit(train_data.data, train_data.target)
9 predicted = pipeline.predict(test_data.data)
```

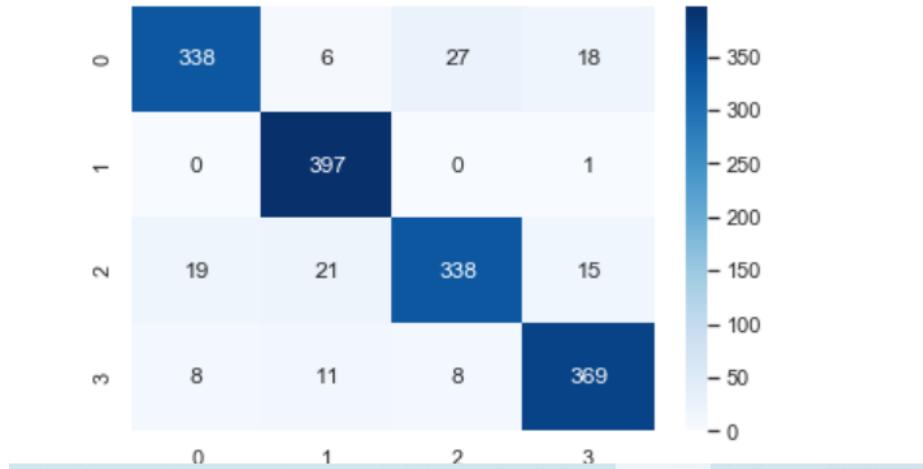
```
In [6]: 1 # accuracy of the above
2 cm=metrics.confusion_matrix(test_data.target,predicted)
3 ax = sns.heatmap(cm, annot=True, fmt="d", cmap="Blues")
4 print(metrics.classification_report(test_data.target,predicted))
5 print(metrics.accuracy_score(test_data.target,predicted))
```

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.93	0.87	0.90	389
1	0.91	1.00	0.95	398
2	0.91	0.86	0.88	393
3	0.92	0.93	0.92	396

accuracy				0.91	1576
macro avg	0.92	0.91	0.91	1576	
weighted avg	0.92	0.91	0.91	1576	

0.9149746192893401



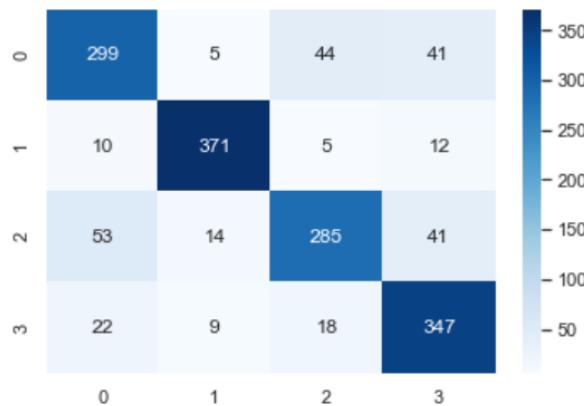
Gaussian Naive bayes Classifier

```
In [7]: 1 gnb = GaussianNB()
2
3 #countvectorizer
4 count_vect = CountVectorizer()
5 X_train_counts = count_vect.fit_transform(train_data.data)
6
7 #tfidf vectorizer
8 tfidf_transformer = TfidfTransformer()
9 X_train_tfidf = tfidf_transformer.fit_transform(X_train_counts)
10 X_train_tfidf=X_train_tfidf.todense()
11
12 #fitting data into gaussian NB classifier
13 clf = gnb.fit(X_train_tfidf, train_data.target)
14
15 #transforming the test data
16 X_new_counts = count_vect.transform(test_data.data)
17
18 X_new_tfidf = tfidf_transformer.transform(X_new_counts)
19 X_new_tfidf=X_new_tfidf.todense()
20
21 # making prediction
22 predicted = gnb.predict(X_new_tfidf)
23
```

```
In [8]: 1 #accuracy of the above
2 cm=metrics.confusion_matrix(test_data.target,predicted)
3 ax = sns.heatmap(cm, annot=True, fmt="d", cmap="Blues")
4 print(metrics.classification_report(test_data.target,predicted))
5 print(metrics.accuracy_score(test_data.target,predicted))
```

	precision	recall	f1-score	support
0	0.78	0.77	0.77	389
1	0.93	0.93	0.93	398
2	0.81	0.73	0.77	393
3	0.79	0.88	0.83	396
accuracy			0.83	1576
macro avg	0.83	0.83	0.82	1576
weighted avg	0.83	0.83	0.83	1576

0.8261421319796954



```
In [ ]: 1 # as we can see above the best result is given by Multinomial Naïve bayes classifier
```

Conclusion:

After seeing the accuracy of the above algos we see that multinomial naïve bayes has the highest accuracy of 92%.

EXPERIMENT NO. 8

Student Name and Roll Number: SHRUTI JAIN (17CSU186)
Semester /Section: VIII-D
Date: 2 May 2021
Faculty Signature:
Link to Code: https://github.com/shruti17csu186/TWIA-LAB-EXPERIMENTS/blob/main/k-means%20hierarchical%20and%20k-medoids%20clustering%20.ipynb
Grade:

Objectives:

1. Implementation of K Means, K Medoids and Hierarchical Clustering algorithms on text data.

Background Study:

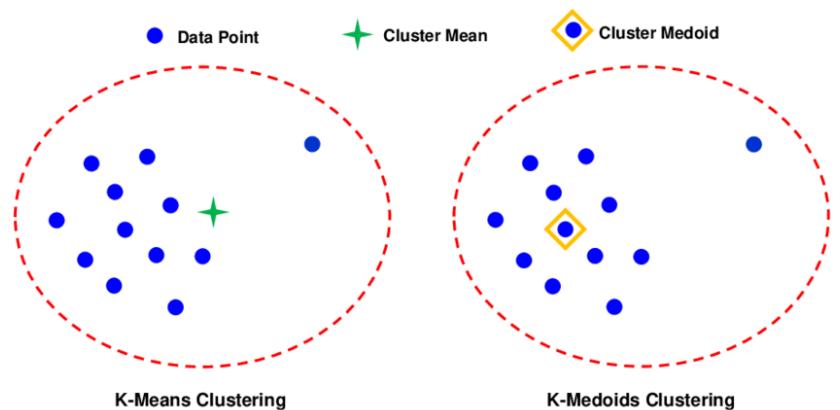
K-Means

K-Means Clustering is an unsupervised Learning algorithm, which groups the unlabeled dataset into different clusters in such a way that each dataset belongs to the group which has similar properties. Here K defines the number of pre-defined clusters that need to be created in the process.

It is a centroid-based algorithm, where each cluster is associated with a centroid. The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters.

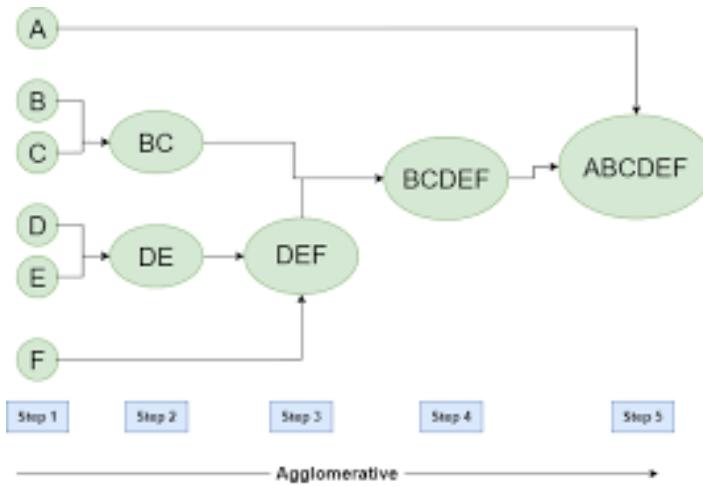
K-Medoids

K-medoids chooses actual data points as centers (medoids or exemplars), and thereby allows for greater interpretability of the cluster centers than in k -means, where the center of a cluster is not necessarily one of the input data points (it is the average between the points in the cluster). Furthermore, k -medoids can be used with arbitrary dissimilarity measures, whereas k -means generally requires Euclidean distance for efficient solutions.



Hierarchical Clustering

Hierarchical clustering, also known as hierarchical cluster analysis, is an algorithm that groups similar objects into groups called clusters. The endpoint is a set of clusters, where each cluster is distinct from each other cluster, and the objects within each cluster are broadly similar to each other.



Outcome: Student will be able to differentiate between the three clustering algorithms and compare their accuracies.

Problem Statement: Implement K-Means, K-Medoids and Hierarchical Clustering Algorithms on text data and plot visualizations for the same.

CODE AND OUTPUT:

importing libraries

```
In [2]: 1 import numpy as np
2 from sklearn.datasets import fetch_20newsgroups
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.feature_extraction.text import TfidfTransformer
5 from sklearn_extra.cluster import KMedoids
6 from sklearn.pipeline import Pipeline
7 from sklearn.cluster import KMeans
8 import scipy.cluster.hierarchy as shc
9 from sklearn.cluster import AgglomerativeClustering
10 import pandas as pd
11 from sklearn import metrics
12 import matplotlib.pyplot as plt
13 import seaborn as sns; sns.set()
14 from sklearn.decomposition import TruncatedSVD
```

importing Data

```
In [3]: 1 categories = ['rec.motorcycles', 'sci.electronics',
2                  'comp.graphics', 'sci.med']
3
4 # sklearn provides us with subset data for training and testing
5 train_data = fetch_20newsgroups(subset='train',
6                                 categories=categories, shuffle=True, random_state=42)
7
8 test_data = fetch_20newsgroups(subset='test',
9                               categories=categories, shuffle=True, random_state=42)
```

K-means Clustering

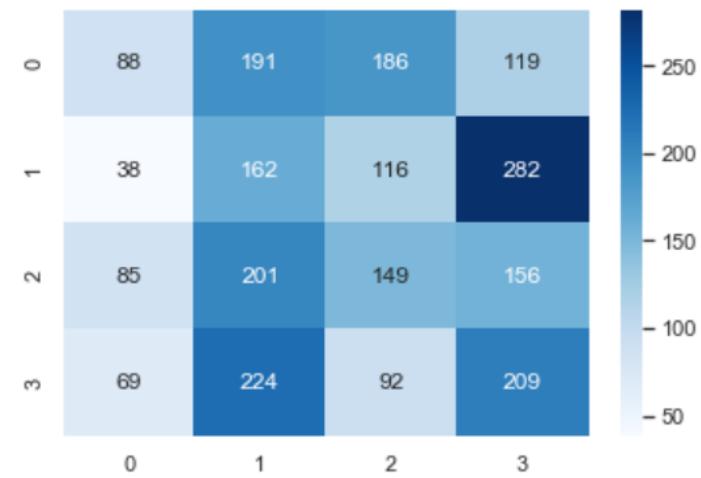
```
In [4]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 tfidf_transformer = TfidfVectorizer(max_df=0.5, min_df=0.2, stop_words='english')
3 X_train_tfidf = tfidf_transformer.fit_transform(train_data.data)
4 km = KMeans(n_clusters=4, init='k-means++', max_iter=100, n_init=1)
5 km.fit(X_train_tfidf)
6 predicted=km.predict(X_train_tfidf)
7 print("Homogeneity: %0.3f" % metrics.homogeneity_score(train_data.target, predicted))
8 print("Completeness: %0.3f" % metrics.completeness_score(train_data.target, predicted))
9 print("V-measure: %0.3f" % metrics.v_measure_score(train_data.target, predicted))
10 print("Silhouette Coefficien: %0.3f" % metrics.silhouette_score(X_train_tfidf, predicted , sample_size=1000))
11
```

Homogeneity: 0.023
Completeness: 0.024
V-measure: 0.023
Silhouette Coefficien: 0.158

```
In [5]: 1 cm=metrics.confusion_matrix(train_data.target,predicted)
2 ax = sns.heatmap(cm, annot=True, fmt="d", cmap="Blues")
3 print(metrics.classification_report(train_data.target,predicted))
4 print(metrics.accuracy_score(train_data.target,predicted))
```

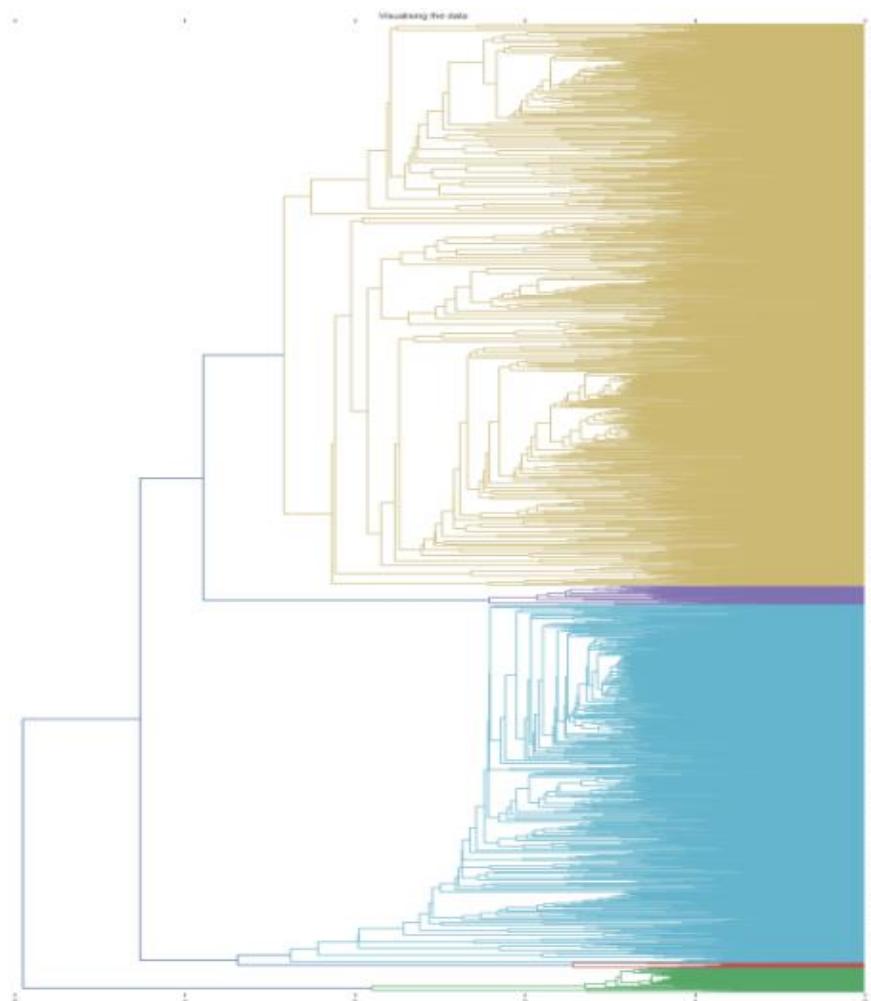
	precision	recall	f1-score	support
0	0.31	0.15	0.20	584
1	0.21	0.27	0.24	598
2	0.27	0.25	0.26	591
3	0.27	0.35	0.31	594
accuracy			0.26	2367
macro avg	0.27	0.26	0.25	2367
weighted avg	0.27	0.26	0.25	2367

0.25686523024926067



Hierachical Clustering

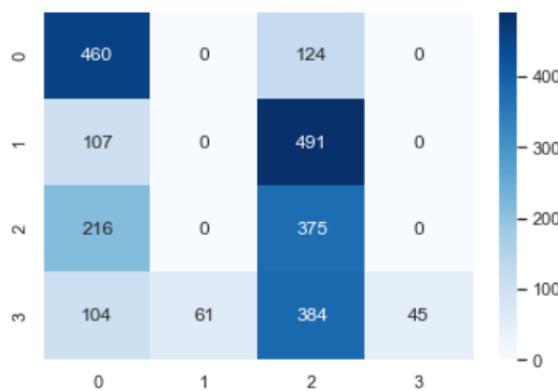
```
In [*]: 1 hc = AgglomerativeClustering(n_clusters=4, affinity='euclidean', linkage='ward')
2 #countvectorizer
3 count_vect = CountVectorizer()
4 X_train_counts = count_vect.fit_transform(train_data.data)
5
6 #tfidf vectorizer
7 tfidf_transformer = TfidfTransformer()
8 X_train_tfidf = tfidf_transformer.fit_transform(X_train_counts)
9 X_train_tfidf=X_train_tfidf.todense()
10
11 #plotting of dendrogram
12 fig, ax = plt.subplots(figsize=(15, 20)) # set size
13 plt.title('Visualising the data')
14 ax= shc.dendrogram(shc.linkage(X_train_tfidf, method = 'ward')),orientation="left")
15
16
17 plt.tick_params(\n    axis= 'x',
18     which='both',
19     bottom='off',
20     top='off',
21     labelbottom='off')
22
23 plt.tight_layout()
24
25
26 #fitting data into Heirarchical classifier
27 hc = hc.fit_predict(X_train_tfidf)
28
29
```



```
In [11]: 1 cm=metrics.confusion_matrix(train_data.target, hc)
2 ax = sns.heatmap(cm, annot=True, fmt="d", cmap="Blues")
3 print(metrics.classification_report(train_data.target, hc))
4 print(metrics.accuracy_score(train_data.target, hc))
5
```

	precision	recall	f1-score	support
0	0.52	0.79	0.63	584
1	0.00	0.00	0.00	598
2	0.27	0.63	0.38	591
3	1.00	0.08	0.14	594
accuracy			0.37	2367
macro avg	0.45	0.37	0.29	2367
weighted avg	0.45	0.37	0.28	2367

0.3717786227291931



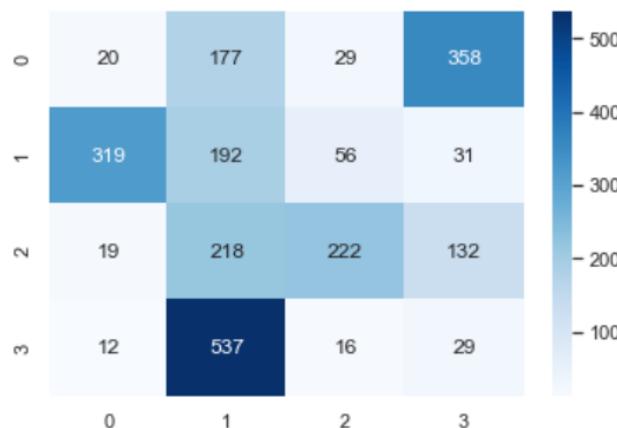
K-medoids Clustering

```
In [6]: 1 kMedoids = KMedoids(n_clusters = 4, random_state = 0)
2 #pipelining the fuctions
3 pipeline = Pipeline([
4     ('vect', CountVectorizer()),
5     ('tfidf', TfidfTransformer()),
6     ('clf', kMedoids)
7 ])
8 pipeline.fit(train_data.data)
9 predicted = pipeline.predict(train_data.data)
10
```

```
In [7]: 1 cm=metrics.confusion_matrix(train_data.target,predicted)
2 ax = sns.heatmap(cm, annot=True, fmt="d", cmap="Blues")
3 print(metrics.classification_report(train_data.target,predicted))
4 print(metrics.accuracy_score(train_data.target,predicted))
5
```

	precision	recall	f1-score	support
0	0.05	0.03	0.04	584
1	0.17	0.32	0.22	598
2	0.69	0.38	0.49	591
3	0.05	0.05	0.05	594
accuracy			0.20	2367
macro avg	0.24	0.19	0.20	2367
weighted avg	0.24	0.20	0.20	2367

0.19560625264047318



```
In [9]: ┌─ 1 ── from sklearn.feature_extraction.text import TfidfVectorizer
  2 ── tfidf_transformer = TfidfVectorizer(max_df=0.5, min_df=0.2, stop_words='english')
  3 ── X_train_tfidf = tfidf_transformer.fit_transform(train_data.data)
  4 ── kMedoids = KMedoids(n_clusters = 3, random_state = 0)
  5 ── kMedoids.fit(X_train_tfidf)
  6 ── predicted=km.predict(X_train_tfidf)
  7 ── print("Homogeneity: %0.3f" % metrics.homogeneity_score(train_data.target, predicted))
  8 ── print("Completeness: %0.3f" % metrics.completeness_score(train_data.target, predicted))
  9 ── print("V-measure: %0.3f" % metrics.v_measure_score(train_data.target,predicted))
 10 ── print("Silhouette Coefficien: %0.3f" % metrics.silhouette_score(X_train_tfidf, predicted , sample_size=1000))
 11
```

Homogeneity: 0.023
Completeness: 0.024
V-measure: 0.023
Silhouette Coefficien: 0.151

EXPERIMENT NO. 9

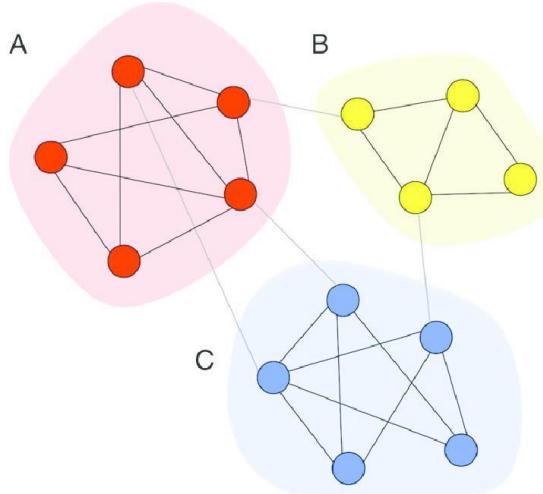
Student Name and Roll Number: SHRUTI JAIN (17CSU186)
Semester /Section: VIII-D
Date: 9 May 2021
Faculty Signature:
Link to Code: https://github.com/shruti17csu186/TWIA-LAB-EXPERIMENTS/blob/main/homophily.docx
Grade:

Objectives:

1. Implement homophily for social media analysis.

Background Study:
Homophily

Network homophily refers to the theory in network science which states that, based on node attributes, similar nodes may be more likely to attach to each other than dissimilar ones. The hypothesis is linked to the model of preferential attachment and it draws from the phenomenon of homophily in social sciences and much of the scientific analysis of the creation of social ties based on similarity comes from network science. In fact, empirical research seems to indicate the frequent occurrence of homophily in real networks.



Outcome:

Students will be able to perform a case study on Social Network Analysis for a population.

Problem Statement:

In this case study, you will investigate homophily of several characteristics of individuals connected in social networks in rural India.

You will calculate the chance homophily for an arbitrary characteristic. Homophily is the proportion of edges in the network whose constituent nodes share that characteristic. How much homophily do we expect by chance? If characteristics are distributed completely randomly, the probability that two nodes x and y share characteristic a is the probability both nodes have characteristic a , which is the frequency of a squared. The total probability that

nodes x and y share their characteristic is therefore the sum of the frequency of each characteristic in the network.

DATA CAMP CASE STUDY QUESTIONS AND ANSWERS

Q1. Network homophily occurs when nodes that share an edge share a characteristic more often than nodes that do not share an edge. In this case study, we will investigate homophily of several characteristics of individuals connected in social networks in rural India.

In this exercise, we will calculate the chance homophily for an arbitrary characteristic. Homophily is the proportion of edges in the network whose constituent nodes share that characteristic. How much homophily do we expect by chance? If characteristics are distributed completely randomly, the probability that two nodes x and y share characteristic a is the probability both nodes have characteristic a , which is the frequency of a squared. The total probability that nodes x and y share their characteristic is therefore the sum of the frequency of each characteristic in the network. For example, in the dictionary `favorite_colors` provided, the frequency of `red` and `blue` is $1/3$ and $2/3$ respectively, so the chance homophily is $(1/3)^2 + (2/3)^2 = 5/9$.

INSTRUCTIONS

- Create a function that takes a dictionary `chars` with personal IDs as keys and characteristics as values, and returns a dictionary with characteristics as keys, and the frequency of their occurrence as values.
- Create a function `chance_homophily(chars)` that takes a dictionary `chars` defined as above and computes the chance homophily for that characteristic.
- A sample of three peoples' favorite colors is given in `favorite_colors`. Use your function to compute the chance homophily in this group, and store as `color_homophily`.
- Print `color_homophily`

CODE:

```
script.py | Light Mode
1  from collections import Counter
2
3
4  def frequency(chars):
5      # Enter code here!
6      chars_counts_dict = Counter(chars.values())
7      print(chars_counts_dict)
8      return chars_counts_dict
9
10 def chance_homophily(chars):
11     # Enter code here!
12     chars_counts_dict=frequency(chars)
13     chars_counts = np.array(list(chars_counts_dict.values()))
14     print(chars_counts)
15     chars_props  = chars_counts / sum(chars_counts)
16     return sum(chars_props**2)
17
18
19 favorite_colors = {
20     "ankit": "red",
21     "xiaoyu": "blue",
22     "mary": "blue"
23 }
24
25 color_homophily = chance_homophily(favorite_colors)
26 print(color_homophily)
27
```

5 Run Code Submit Answer

IPython Shell

OUTPUT:

```
IPython Shell
def chance_homophily(chars):
    # Enter code here!
    chars_counts_dict=frequency(chars)
    chars_counts = np.array(list(chars_counts_dict.values()))
    print(chars_counts)
    chars_props  = chars_counts / sum(chars_counts)
    return sum(chars_props**2)

favorite_colors = {
    "ankit": "red",
    "xiaoyu": "blue",
    "mary": "blue"
}

color_homophily = chance_homophily(favorite_colors)
print(color_homophily)

Counter({'blue': 2, 'red': 1})
[2 1]
0.5555555555555556
```

VALUE ADDED

EXPERIMENT NO. 1

Student Name and Roll Number: SHRUTI JAIN (17CSU186)
Semester /Section: VIII-D
Date: 6 March 2021
Faculty Signature:
Link to Code: https://github.com/shruti17csu186/TWIA-LAB-EXPERIMENTS/blob/main/chunking%20and%20chinkin.ipynb
Grade:

Objective:

1. Apply chunking on a sentence and explore text corpora.
2. Apply chinking on a piece of text.
3. Implementation of Named Entity recognition.

Background Study:

Chunking

Chunking refers to the process of taking individual pieces of information and grouping them into larger units. By grouping each data point into a larger whole, we can improve the amount of information we can remember.

Chinking

Chinking is the process of removing a sequence of tokens from a chunk. If the matching sequence of tokens spans an entire chunk, then the whole chunk is removed; if the sequence of tokens appears in the middle of the chunk, these tokens are removed, leaving two chunks where there was only one before. If the sequence is at the periphery of the chunk, these tokens are removed, and a smaller chunk remains.

Named Entity

Named entities are definite noun phrases that refer to specific types of individuals, such as organizations, persons, dates, and so on.

Named Entity Recognition

Named entity recognition (NER) is probably the first step towards information extraction that seeks to locate and classify named entities in text into pre-defined categories such as the names of persons, organizations, locations, expressions of times, quantities, monetary values, percentages, etc.

Outcome:

Students will be able to understand the concept such as chunking, regex parser, chinking and tagging. They will also be able to explore the corpus brown and apply named entity recognition using ne_chunk() functions

Problem Statement:

CODE AND OUTPUT:

Chunking

```
In [2]: #Chunking
# Loading Libraries
from nltk.chunk.regexp import ChunkString, ChunkRule
from nltk.tree import Tree
# ChunkString() starts with the flat tree
tree = Tree('S', [('the', 'DT'), ('book', 'NN'),
                  ('has', 'VBZ'), ('many', 'JJ'), ('chapters', 'NNS')])
tree.draw()
chunk_string = ChunkString(tree)
print ("Chunk String : ", chunk_string)
```

Chunk String : <DT> <NN> <VBZ> <JJ> <NNS>

```
In [3]: grammar = r"""
NP: {<DT>|<JJ>|<NN>.*+}          # Chunk sequences of DT, JJ, NN
PP: {<IN><NP>}                 # Chunk prepositions followed by NP
VP: {<VB.*>} # Chunk verbs and their arguments
"""
In [5]: import nltk
chunker = nltk.RegexpParser(grammar)
sent=chunker.parse(tree)
sent.draw()
```

```
In [6]: chunker1 = nltk.RegexpParser(grammar)
sent = [('the', 'DT'), ('sushi', 'NN'), ('roll', 'NN'), ('was', 'VBD'),
        ('filled', 'VBN'), ('with', 'IN'), ('the', 'DT'), ('fish', 'NN')]
result2=chunker.parse(sent)
result2.draw()
```

```
In [10]: #Name Entity Recognition using NLTK
import nltk
nltk.download('maxent_ne_chunker')
nltk.download('words')
my_sent = "WASHINGTON -- In the wake of a string of abuses by New York police officers in the 1990s, Loretta E. Lynch, th
parse_tree = nltk.ne_chunk(nltk.tag.pos_tag(nltk.word_tokenize(my_sent)), binary=True)
parse_tree.draw()
print(parse_tree)
12
```

```
In [11]: ┶ 1 from nltk.tree import Tree
2 nltk.download('maxent_ne_chunker')
3 txt="WASHINGTON -- In the wake of a string of abuses by New York police officers in the 1990s, Loretta E. Lynch, the top
4
5 pos_tag = nltk.pos_tag(txt.split())
6 parse_tree = nltk.ne_chunk(pos_tag )
7 # print(chunk)
8 parse_tree.draw()
9 NE=[]
10 for chunk in parse_tree:
11     if hasattr(chunk, 'label'):
12         NE=(chunk.label(), ' '.join(c[0] for c in chunk))
13         print(NE)
14

[nltk_data] Downloading package maxent_ne_chunker to
[nltk_data]      C:\Users\Shruti\AppData\Roaming\nltk_data...
[nltk_data]  Package maxent_ne_chunker is already up-to-date!

('GPE', 'WASHINGTON')
('GPE', 'New York')
('PERSON', 'Loretta E.')
```

VALUE ADDED

EXPERIMENT NO. 2

Student Name and Roll Number: SHRUTI JAIN (17CSU186)
Semester /Section: VIII-D
Date: 9 May 2021
Faculty Signature:
Link to Code: https://github.com/shruti17csu186/TWIA-LAB-EXPERIMENTS/blob/main/FATMAN%20EVOLUTIONARY%20NETWORK%20IMPLEMENTATION.ipynb
Grade:

Objective:

To implement Fatman Evolutionary Model in Social Networks on python jupyter notebook.

Background Study:

In the Fatman evolutionary model, there are 3 main concepts to make an evolutionary model i.e Homophily, Social influence, and closure.

1. **Homophily-** People who are similar to each other, they tend to make friends with each other.
 2. **Social Influence-** People change their behavior and properties because of social influence.
 3. **Closure-** There are 3 main kinds of closures i.e triadic closure, membership closure, and foci closure.
-
- **Triadic closure-** In triadic closure, if A is a friend of B and B is a friend of C then C will eventually become friends of A.
 - **Membership closure-** In membership closure, if A and B are in the same club then there is a tendency that they will become friends.
 - **Foci closure-** It is the probability of becoming friends when they have the same foci.
So, According to Fatman Hypothesis, beware of fat friends, if your friend is fat then the probability of you gaining the weight becomes high.

Outcome:

Students will be able to understand the concept of Fatman evolutionary model and able to implement it on python.

Problem Statement:

Assume there is a city with several people, and we have everyone's BMI. Then we will see as evolution continues people having similar BMI become friends to each other which is Homophily. Then we will create some foci like Gym etc. and see how eventually people with similar foci become a friend and then we will implement social foci

CODE AND OUTPUT:

```
In [2]: M
1 import networkx as nx
2 import matplotlib.pyplot as plt
3 import random
4 import math
5
6 # 1- Create a graph with Lets say 100 nodes
7 # 2- Add edges and Labels.
8 def create():
9     G = nx.Graph()
10    G.add_nodes_from(range(1, 101))
11    return G
12
13 def visualize(G):
14     labeledict = get_labels(G)
15     nodesize = get_size(G)
16     color = get_colors(G)
17     nx.draw(G, labels=labeledict, node_size=nodesize, node_color=color)
18     plt.show()
19
20
21 def assign_bmi(G):
22
23     for each in G.nodes():
24         G.nodes[each]['name'] = random.randint(15, 40)
25         G.nodes[each]['type'] = 'person'
26
27
28 def get_labels(G):
29     dict1 = {}
30
31     for each in G.nodes():
32         dict1[each] = G.nodes[each]['name']
33
34     return dict1
35
36 def get_size(G):
37     array1 = []
38     for each in G.nodes():
39         if (G.nodes[each]['type'] == 'person'):
40             array1.append(G.nodes[each]['name'] * 20)
41         else:
42             array1.append(500)
43
44     return array1
45
46 # 3. Add foci nodes
47 def add_foci_nodes(G):
48     n = G.number_of_nodes()
49     i = n + 1
50     foci_nodes = ['gym', 'eatout', 'movie_club',
51                  'karate_club', 'yoga_club']
52
53     for j in range(5):
54         G.add_node(i)
55         G.nodes[i]['name'] = foci_nodes[j]
56         G.nodes[i]['type'] = 'foci'
57         i += 1
```

```

59 def get_colors(G):
60     c = []
61
62     for i in G.nodes():
63         if (G.nodes[i]['type'] == 'person'):
64             if (G.nodes[i]['name'] == 15):
65                 c.append('yellow')
66             elif (G.nodes[i]['name'] == 40):
67                 c.append('green')
68             else:
69                 c.append('blue')
70         else:
71             c.append('red')
72     return c
73
74 def get_person_nodes(G):
75     p = []
76
77     for i in G.nodes():
78         if (G.nodes[i]['type'] == 'person'):
79             p.append(i)
80     return p
81
82 def get_foci_nodes(G):
83     f = []
84
85     for i in G.nodes():
86         if (G.nodes[i]['type'] == 'foci'):
87             f.append(i)
88     return f
89
90 def add_foci_edges(G):
91     foci_nodes = get_foci_nodes(G)
92     person_nodes = get_person_nodes(G)
93
94     for i in person_nodes:
95         r = random.choice(foci_nodes)
96         G.add_edge(i, r)
97
98 def homophily(G):
99     pnodes = get_person_nodes(G)
100
101    for u in pnodes:
102        for v in pnodes:
103            if (u != v):
104                diff = abs(G.nodes[u]['name'] - G.nodes[v]['name'])
105                p = 1 / (diff + 1000)
106                r = random.uniform(0, 1)
107                if (r < p):
108                    G.add_edge(u, v)
109
110
111 def common(u, v, G):
112     nu = set(G.neighbors(u))
113     nv = set(G.neighbors(v))
114     return (len(nu & nv))
115
116 def closure(G):
117     array1 = []
118     for u in G.nodes():
119         for v in G.nodes():
120             if (u != v and G.nodes[u]['type'] == 'person' or G.nodes[v]['type'] == 'person'):
121
122                 # common function will return start node,
123                 # end node and Graph itself.
124                 k = common(u, v, G)
125                 p = 1 - math.pow(1 - 0.01, k)
126                 tmp = []
127                 tmp.append(u)
128                 tmp.append(v)
129                 tmp.append(p)
130                 array1.append(tmp)
131
132     for i in array1:
133         u = i[0]
134         v = i[1]
135         p = i[2]
136         r = random.uniform(0, 1)
137         if r < p:
138             G.add_edge(u, v)
139
140 def social_influence(G):
141     fnodes = get_foci_nodes(G)
142
143     for i in fnodes:
144         if (G.nodes[i]['name'] == 'eatout'):
145             for j in G.neighbors(i):
146                 if (G.nodes[j]['name'] != 40):
147                     G.nodes[j]['name'] = G.nodes[j]['name'] + 1
148
149         if (G.nodes[i]['name'] == 'gym'):
150             for j in G.neighbors(i):
151                 if (G.nodes[j]['name'] != 15):
152                     G.nodes[j]['name'] = G.nodes[j]['name'] - 1
153
154 G = create()
155 assign_bmi(G)
156
157 # to add foci field in the nodes
158 add_foci_nodes(G)
159 add_foci_edges(G)
160 visualize(G)
161
162 for t in range(1, 3):
163     homophily(G)
164     visualize(G)
165     closure(G)
166     visualize(G)
167     social_influence(G)
168     visualize(G)

```

