

CS-313 : TEXT MINING (LAB)

Shubham Dey (M.Sc. CS)

Roll No. : 23419CMP026

[GitHub Repo](#)

Program 01

Read from a '.txt' file and count the number of sentences, words, and characters in the file.

```
1  # count no. of sentences
2  def countLines(file):
3      f = open(file, 'r')
4      dat = f.read().split('.')
5      f.close()
6      dat.pop()
7      count = len(dat)
8      return count
9
10 # count no. of words
11 def countWords(file):
12     f = open(file, 'r')
13     count = 0
14     line = f.readline()
15     while line:
16         count += len(line.split())
17         line = f.readline()
18     f.close()
19     return count
20
21 # count no. of characters
22 def countChars(file):
23     f = open(file, 'r')
24     count = 0
25     line = f.readline()
26     while line:
27         for word in line.split():
28             for i in word:
29                 if i.isalnum():
30                     count += 1
31         line = f.readline()
32     f.close()
33     return count
34
35 print('No. of Sentences = ', countLines(r'Files/data.txt'))
36 print('No. of Words = ', countWords(r'Files/data.txt'))
37 print('No. of Characters = ', countChars(r'Files/data.txt'))
```

Program 02

Read from a '.pdf' file and count the number of sentences, words, and characters in the file.

```
1 from pypdf import PdfReader
2
3 # count no. of sentences
4 def countLines(file):
5     f = PdfReader(file)
6     count = 0
7     for page in f.pages:
8         dat = page.extract_text().split('.')
9         dat.pop()
10        count += len(dat)
11    return count
12
13 # count no. of words
14 def countWords(file):
15     f = PdfReader(file)
16     count = 0
17     for page in f.pages:
18         dat = page.extract_text().split(' ')
19         for i in dat:
20             if i != '':
21                 count += 1
22    return count
23
24 # count no. of characters
25 def countChars(file):
26     f = PdfReader(file)
27     count = 0
28     for page in f.pages:
29         dat = page.extract_text()
30         for i in dat:
31             if i.isalnum():
32                 count += 1
33    return count
34
35 print('No. of Sentences = ', countLines(r'Files/data.pdf'))
36 print('No. of Words = ', countWords(r'Files/data.pdf'))
37 print('No. of Characters = ', countChars(r'Files/data.pdf'))
```

Program 03

Read from a '.docx' file and count the number of sentences, words, and characters in the file.

```
1 from docx import Document
2
3 # count no. of sentences
4 def countLines(file):
5     f = Document(file)
6     count = 0
7     for para in f.paragraphs:
8         dat = para.text.split('.')
9         dat.pop()
10        count += len(dat)
11    return count
12
13 # count no. of words
14 def countWords(file):
15     f = Document(file)
16     count = 0
17     for para in f.paragraphs:
18         dat = para.text.split(' ')
19         for i in dat:
20             if i != '':
21                 count += 1
22    return count
23
24 # count no. of characters
25 def countChars(file):
26     f = Document(file)
27     count = 0
28     for para in f.paragraphs:
29         dat = para.text
30         for i in dat:
31             if i.isalnum():
32                 count += 1
33    return count
34
35 print('No. of Sentences = ', countLines(r'Files/data.docx'))
36 print('No. of Words = ', countWords(r'Files/data.docx'))
37 print('No. of Characters = ', countChars(r'Files/data.docx'))
```

Program 04

Implement Term-Document incidence matrix for boolean retrieval.

```
1 import os
2 import pandas as pd
3
4 def getFiles():
5     lst = []
6     for f in os.listdir():      # gives list of files in current
        ↳ directory
7         if f.endswith('.txt'):  # taking only .txt files from current
            ↳ directory
8             lst.append(f)
9     return lst
10
11 # create term-document matrix
12 def createMatrix(files):
13     fp = []      # index of this list work as documentID
14     for file in files:
15         f = open(file, 'r')
16         fp.append(f)
17     words = []  # to store unique words which acts as rows
18     matrix = []
19     for i, f in enumerate(fp):
20         dat = f.read().split()
21         for word in dat:
22             if word not in words:
23                 words.append(word)
24                 matrix.append([0 for _ in range(len(fp))])
25                 matrix[words.index(word)][i] = 1
26     return matrix, words
27
28 # display term-document matrix
29 def printMatrix(matrix, words, files):
30     df = pd.DataFrame(matrix, columns=files, index=words)
31     print(df)
32
33 # perform AND Query
34 def andQuery(w1, w2, matrix, words):
35     try:      # gets row of word-1 if available, else error
36         l1 = matrix[words.index(w1)]
37         l1 = int(''.join(map(str, l1)), 2)
38     except ValueError:
39         return -1
40     try:      # gets row of word-2 if available, else error
41         l2 = matrix[words.index(w2)]
```

```

42         l2 = int(''.join(map(str, l2)), 2)
43     except ValueError:
44         return -2
45     return bin(l1 & l2)[2:]      # bitwise and
46
47 # perform OR Query
48 def orQuery(w1, w2, matrix, words):
49     try:      # gets row of word-1 if available, else works with word-2
50         ↪ only
51         l1 = matrix[words.index(w1)]
52         l1 = int(''.join(map(str, l1)), 2)
53     except ValueError:
54         print(f'{w1} Not Found')
55         l1 = 0
56     try:      # gets row of word-2 if available, else works with word-1
57         ↪ only
58         l2 = matrix[words.index(w2)]
59         l2 = int(''.join(map(str, l2)), 2)
60     except ValueError:
61         print(f'{w2} Not Found')
62         l2 = 0
63     return bin(l1 | l2)[2:]      # bitwise or
64
65 # perform NOT Query
66 def notQuery(w, matrix, words):
67     try:      # gets row of word if available, else gives all 0's
68         l = matrix[words.index(w)]
69         l = ''.join(map(str, l))
70     except:
71         l = '0' * len(matrix[0])
72     return ''.join(['0' if i=='1' else '1' for i in l])      #
73     ↪ complementing
74
75 if __name__ == '__main__':
76     files = getFiles()
77     matrix, words = createMatrix(files)
78
79     print('\nTerm-Document Incidence Matrix: ')
80     printMatrix(matrix, words, files)
81
82     print('\nVocabulary:')
83     print(words)
84
85     print('\nTerms in lowercase & Operator in uppercase')
86     query = input('Enter query: ').split()

```

```

85
86 if len(query) == 3:
87     if query[1] == 'AND':
88         ans = andQuery(query[0], query[2], matrix, words)
89         if type(ans) == 'str':
90             ans = ans.zfill(len(matrix[0]))
91         if ans == -1:
92             print(f'\n{query[0]} Not Found.')
93         elif ans == -2:
94             print(f'\n{query[2]} Not Found.')
95         elif ans == '0'*len(matrix[0]):
96             print(f'\n{query[0]} AND {query[2]} Not Available.')
97         else:
98             print(f'\n{query[0]} AND {query[2]} are Available in
99                 ↪ :')
100             for i in range(len(ans)):
101                 if ans[i] == '1':
102                     print(files[i])
103
104     elif query[1] == 'OR':
105         ans = orQuery(query[0], query[2], matrix, words)
106         ans = ans.zfill(len(matrix[0]))
107         if ans != '0'*len(matrix[0]):
108             print(f'\n{query[0]} OR {query[2]} are Available in
109                 ↪ :')
110             for i in range(len(ans)):
111                 if ans[i] == '1':
112                     print(files[i])
113
114     else:
115         print("Invalid Query")
116
117 elif len(query) == 2 and query[0] == 'NOT':
118     ans = notQuery(query[1], matrix, words)
119     if ans == '0'*len(matrix[0]):
120         print(f'\n{query[1]} Available in ALL Docs.')
121     else:
122         print(f'\n{query[1]} is NOT Available in :')
123         for i in range(len(ans)):
124             if ans[i] == '1':
125                 print(files[i])
126
127 else:
128     print("Invalid Query")

```

Program 05

Implement Inverted-Index for boolean retrieval.

```
1 import os
2
3 def getFiles():
4     lst = []
5     for f in os.listdir():      # gives list of files in current
        ↳ directory
6         if f.endswith('.txt'): # taking only .txt files from current
            ↳ directory
7             lst.append(f)
8     return lst
9
10 # create inverted index
11 def index(files):
12     fp = []      # index of this list work as documentID
13     for file in files:
14         f = open(file, 'r')
15         fp.append(f)
16     dix = {}     # dictionary: key= terms & value= postings
17     for i, f in enumerate(fp):
18         dat = f.read().split()
19         for word in dat:
20             try:
21                 flag = True # doc not listed for the term
22                 for j in range(len(dix[word])):
23                     if dix[word][j] == i:
24                         flag = False # doc already listed for the term
25                         break
26                     elif dix[word][j] > i:
27                         dix[word].insert(j, i)
28                         flag = False
29                         break
30                 if flag:
31                     dix[word].append(i)
32             except KeyError:
33                 dix[word] = [i] # inserting new term
34     return dix
35
36 # perform AND Query
37 def andQuery(dix, w1, w2):
38     try:      # get posting for word-1, else error
39         l1 = dix[w1]
40     except KeyError:
41         return -1
```

```

42     try:      # get posting for word-2, else error
43         l2 = dix[w2]
44     except KeyError:
45         return -2
46
47     # intersection of the two postings
48     i, j = 0, 0
49     lst = []
50     while i < len(l1) and j < len(l2):
51         if l1[i] == l2[j]:
52             lst.append(l1[i])
53             i += 1
54             j += 1
55         elif l1[i] < l2[j]:
56             i += 1
57         else:
58             j += 1
59     return lst
60
61 # perform OR Query
62 def orQuery(dix, w1, w2):
63     try:      # get posting for word-1, else works with word-2 only
64         l1 = dix[w1]
65     except KeyError:
66         print(f'{w1} Not Found')
67         l1 = []
68     try:      # get posting for word-2, else works with word-1 only
69         l2 = dix[w2]
70     except KeyError:
71         print(f'{w2} Not Found')
72         l2 = []
73
74     # union of the two postings
75     i, j = 0, 0
76     lst = []
77     while i < len(l1) and j < len(l2):
78         if l1[i] == l2[j]:
79             lst.append(l1[i])
80             i += 1
81             j += 1
82         elif l1[i] < l2[j]:
83             lst.append(l1[i])
84             i += 1
85         else:
86             lst.append(l2[j])
87             j += 1

```



```

88     while i < len(l1):
89         lst.append(l1[i])
90         i += 1
91     while j < len(l2):
92         lst.append(l2[j])
93         j += 1
94     return lst
95
96     # perform NOT Query
97     def notQuery(dix, w, N):
98         try:     # get posting for word
99             l = dix[w]
100         except KeyError:
101             l = []
102         lst = [i for i in range(N) if i not in l]
103         return lst
104
105 if __name__ == '__main__':
106     files = getFiles()
107     idx = index(files)
108     print('\nInverted Index:')
109     for i in idx:
110         print(i, ':', idx[i])
111     print('\nVocabulary')
112     for i in idx:
113         print(i, end=', ')
114     print()
115
116     print('\nTerms in lowercase & Operator in uppercase')
117     query = input('Enter query: ').split()
118
119     if len(query) == 3:
120         if query[1] == 'AND':
121             ans = andQuery(idx, query[0], query[2])
122             if ans == -1:
123                 print(f'{query[0]} Not Found.')
124             elif ans == -2:
125                 print(f'{query[2]} Not Found.')
126             elif not ans:
127                 print(f'{query[0]} AND {query[2]} Not Available.')
128             else:
129                 print(f'{query[0]} AND {query[2]} are Available in :')
130                 for i in ans:
131                     print(files[i])
132
133         elif query[1] == 'OR':

```

```
134         ans = orQuery(idx, query[0], query[2])
135         if ans:
136             print(f'{query[0]} OR {query[2]} are Available in :')
137             for i in ans:
138                 print(files[i])
139         else:
140             print("Invalid Query")
141
142     elif len(query) == 2 and query[0] == 'NOT':
143         ans = notQuery(idx, query[1], len(files))
144         if ans:
145             print(f'{query[1]} is NOT Available in :')
146             for i in ans:
147                 print(files[i])
148         else:
149             print(f'\n{query[1]}Available in ALL Docs.')
150     else:
151         print("Invalid Query")
```

Program 06

Using NLTK perform Tokenization, Normalization, Stemming & Lemmetization.

```
1 import nltk
2 nltk.download('punkt')      # download resources for tokenization &
   ↪  punctuations
3 nltk.download('stopwords')  # download resources for stopwords
4 nltk.download('wordnet')    # download resources for lemmetization
5
6 from nltk.tokenize import sent_tokenize
7 from nltk.tokenize import word_tokenize
8 from nltk.probability import FreqDist
9 from nltk.corpus import stopwords
10 from nltk.stem import PorterStemmer
11 from nltk.stem.wordnet import WordNetLemmatizer
12 from truecase import get_true_case
13 import string
14
15 # document
16 text = "Mr. Smith is feeling Relaxed today, as The weather in USA is
   ↪  awesome. Did something troubled Him in U.S.A.? The birds are
   ↪  flying."
17 print('\nDocument:\n', text)
18
19 # Sentence Tokenization
20 tokenized_sent = sent_tokenize(text)
21 print('\nSentence Tokenization:\n', tokenized_sent)
22
23 # Word Tokenization
24 tokenized_word = word_tokenize(text)
25 print('\nWord Tokenization\n', tokenized_word)
26
27 # Frequency Distribution
28 fdist = FreqDist(tokenized_word)
29 print('\nMost frequent 5 words:\n', fdist.most_common(5))
30
31 # Lowercasing
32 lower_token = []
33 for token in tokenized_word:
34     lower_token.append(token.lower())
35 print('\nLowercasing:\n', lower_token)
36
37 # Truecasing Sentences
38 true_text = []
39 for text in tokenized_sent:
40     true_text.append(get_true_case(text))
```

```

41 print('\nTruecasing Sentences:\n', true_text)
42
43 # Tokenize Truecase Words
44 true_token = []
45 for token in true_text:
46     true_token.extend(word_tokenize(token))
47 print('\nTruecase Words\n', true_token)
48
49 # Removing Punctuations
50 punctuations = list(string.punctuation)
51 tokens = []
52 for i in lower_token:
53     if i not in punctuations:
54         tokens.append(i)
55 print('\nAfter removing Punctuations:\n', tokens)
56
57 # Removing Stopwords
58 stopwords_english = stopwords.words("english")
59 filtered_tokens=[]
60 for w in tokens:
61     if w not in stopwords_english:
62         filtered_tokens.append(w)
63 print('\nAfter removing Stopwords:\n', filtered_tokens)
64
65 # Stemming
66 ps = PorterStemmer()
67 stem_words=[]
68 for w in filtered_tokens:
69     stem_words.append(ps.stem(w))
70 print('\nAfter Stemming:\n', stem_words)
71
72 # Lemmetization
73 lem = WordNetLemmatizer()
74 lem_words=[]
75 for w in filtered_tokens:
76     lem_words.append(lem.lemmatize(w))
77 print('\nAfter Lemmetization:\n', lem_words)

```

Program 07

Naive Bayes classification using scikit-learn.

```
1 from sklearn.naive_bayes import GaussianNB, MultinomialNB, BernoulliNB
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 import pandas as pd
4
5 # documents
6 doc1 = 'Chinese Beijing Chinese'
7 doc2 = 'Chinese Chinese Shanghai'
8 doc3 = 'Chinese Macao'
9 doc4 = 'Tokyo Japan Chinese'
10
11 docs = [doc1, doc2, doc3, doc4]    # doc list
12 y = [1, 1, 1, 0]                  # labels: 1=> chinese, 0=> not chinese
13
14 # displaying the training dataset
15 d = {'Documents': docs, 'Labels': y}
16 df = pd.DataFrame(d)
17 print('\n1=> chinese, 0=> not chinese')
18 print('\nTraining Dataset:\n', df)
19
20 # converting text to numeric using tf-idf
21 tfidf = TfidfVectorizer()
22 X = tfidf.fit_transform(docs).toarray()
23 print('\nAfter tf-idf :\n', X)
24 print('\nFeatures:\n', tfidf.get_feature_names_out())
25
26 # training Gaussian Naive Bayes
27 gnb = GaussianNB()
28 gnb.fit(X, y)
29 # training Multinomial Naive Bayes
30 mnb = MultinomialNB()
31 mnb.fit(X, y)
32 # training Bernoulli Naive Bayes
33 bnb = BernoulliNB()
34 bnb.fit(X, y)
35
36 # classifying new document
37 doc5 = ['Chinese Chinese Chinese Tokyo Japan']
38 a = tfidf.transform(doc5).toarray()
39
40 print('\nNew document: ', doc5[0])
41 print('\nGaussian NB: ', gnb.predict(a)[0])
42 print('\nMultinomial NB: ', mnb.predict(a)[0])
43 print('\nBernoulli NB: ', bnb.predict(a)[0])
```

Program 08

Rocchio classification using scikit-learn.

```
1 import pandas as pd
2 from sklearn.neighbors import NearestCentroid
3 from sklearn.feature_extraction.text import TfidfVectorizer
4
5 # documents
6 doc1 = 'Chinese Beijing Chinese'
7 doc2 = 'Chinese Chinese Shanghai'
8 doc3 = 'Chinese Macao'
9 doc4 = 'Tokyo Japan Chinese'
10
11 docs = [doc1, doc2, doc3, doc4]    # doc list
12 y = [1, 1, 1, 0]                  # labels: 1=> chinese, 0=> not chinese
13
14 # displaying the training dataset
15 d = {'Documents': docs, 'Labels': y}
16 df = pd.DataFrame(d)
17 print('\n1=> chinese, 0=> not chinese')
18 print('\nTraining Dataset:\n', df)
19
20 # converting text to numeric using tf-idf
21 tfidf = TfidfVectorizer()
22 X = tfidf.fit_transform(docs).toarray()
23 print('\nAfter tf-idf :\n', X)
24 print('\nFeatures:\n', tfidf.get_feature_names_out())
25
26 # training rocchio classifier
27 rocchio = NearestCentroid()
28 rocchio.fit(X, y)
29
30 # classifying new document
31 doc5 = ['Chinese Chinese Chinese Tokyo Japan']
32 a = tfidf.transform(doc5).toarray()
33
34 print('\nNew document: ', doc5[0])
35 print('Prediction: ', rocchio.predict(a)[0])
```

Program 09

k-means Clustering using scikit-learn.

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.cluster import KMeans
3
4 docs = ["This is the most beautiful place in the world.",
5         "This man has more skills to show in cricket than any other
6         ↪ game.",
7         "Hi there! how was your ladakh trip last month?",
8         "There was a player who had scored 200+ runs in single cricket
9         ↪ innings in his career.",
10        "I have got the opportunity to travel to Paris next year for my
11        ↪ internship.",
12        "May be he is better than you in batting but you are much
13        ↪ better than him in bowling.",
14        "That was really a great day for me when I was there at Lavasa
15        ↪ for the whole night.",
16        "That's exactly I wanted to become, a highest ratting batsmen
17        ↪ ever with top scores.",
18        "Does it really matter wether you go to Thailand or Goa, its
19        ↪ just you have spend your holidays.",
20        "Why don't you go to Switzerland next year for your 25th
21        ↪ Wedding anniversary?",
22        "Travel is fatal to prejudice, bigotry, and narrow mindedness.,
23        ↪ and many of our people need it sorely on these accounts.",
24        "Stop worrying about the potholes in the road and enjoy the
25        ↪ journey.",
26        "No cricket team in the world depends on one or two players.
27        ↪ The team always plays to win.",
28        "Cricket is a team game. If you want fame for yourself, go play
29        ↪ an individual game.",
30        "Because in the end, you won't remember the time you spent
31        ↪ working in the office or mowing your lawn. Climb that
32        ↪ goddamn mountain.",
33        "Isn't cricket supposed to be a team sport? I feel people
34        ↪ should decide first whether cricket is a team game or an
35        ↪ individual sport."]
```

```
21 label = ['travel', 'cricket', 'travel', 'cricket', 'travel', 'cricket',
22         ↪ 'travel', 'cricket', 'travel', 'travel', 'travel', 'travel',
23         ↪ 'cricket', 'cricket', 'travel', 'cricket']
24 true_k = 2      # no. of clusters
25
26 # converting text to numeric using tf-idf
27 tfidf = TfidfVectorizer(stop_words='english')
```

```

26 X = tfidf.fit_transform(docs)
27
28 # training k-means
29 model = KMeans(n_clusters=true_k, max_iter=100)
30 model.fit(X)
31
32 # display top 10 terms in the clusters
33 order_centroids = model.cluster_centers_.argsort()[:, :-1]
34 terms = tfidf.get_feature_names_out()
35 for i in range(true_k):
36     print("\nCluster-%d (top 10 terms):" % i)
37     for ind in order_centroids[i, :10]:
38         print(terms[ind], end=" ")
39     print()
40
41 new = ["Nothing is easy in cricket. Maybe when you watch it on TV, it
↪ looks easy. But it is not. You have to use your brain and time the
↪ ball."]
42 X = tfidf.transform(new)
43 predicted = model.predict(X)[0]
44 print('\nNew Document:\n', new[0])
45 print("\nPrediction = Cluster-", predicted)

```


Program 10

Perform Parts-Of-Speech (POS) Tagging using NLTK

```
1 from nltk import word_tokenize, pos_tag
2
3 import nltk
4 nltk.download('averaged_perceptron_tagger')    # download resources
5         ↪ for POS tagger
6
7 # document
8 sent = "Albert Einstein was born in Ulm Germany in 1879"
9
10 # tokenizing the document
11 tokens = word_tokenize(sent)
12
13 print('Tokens:\n', tokens)
14
15 # POS tagging
16 pos = pos_tag(tokens)
17
18 print('\nPOS tagging of the tokens:')
19 for i in pos:
20     print(i)
```