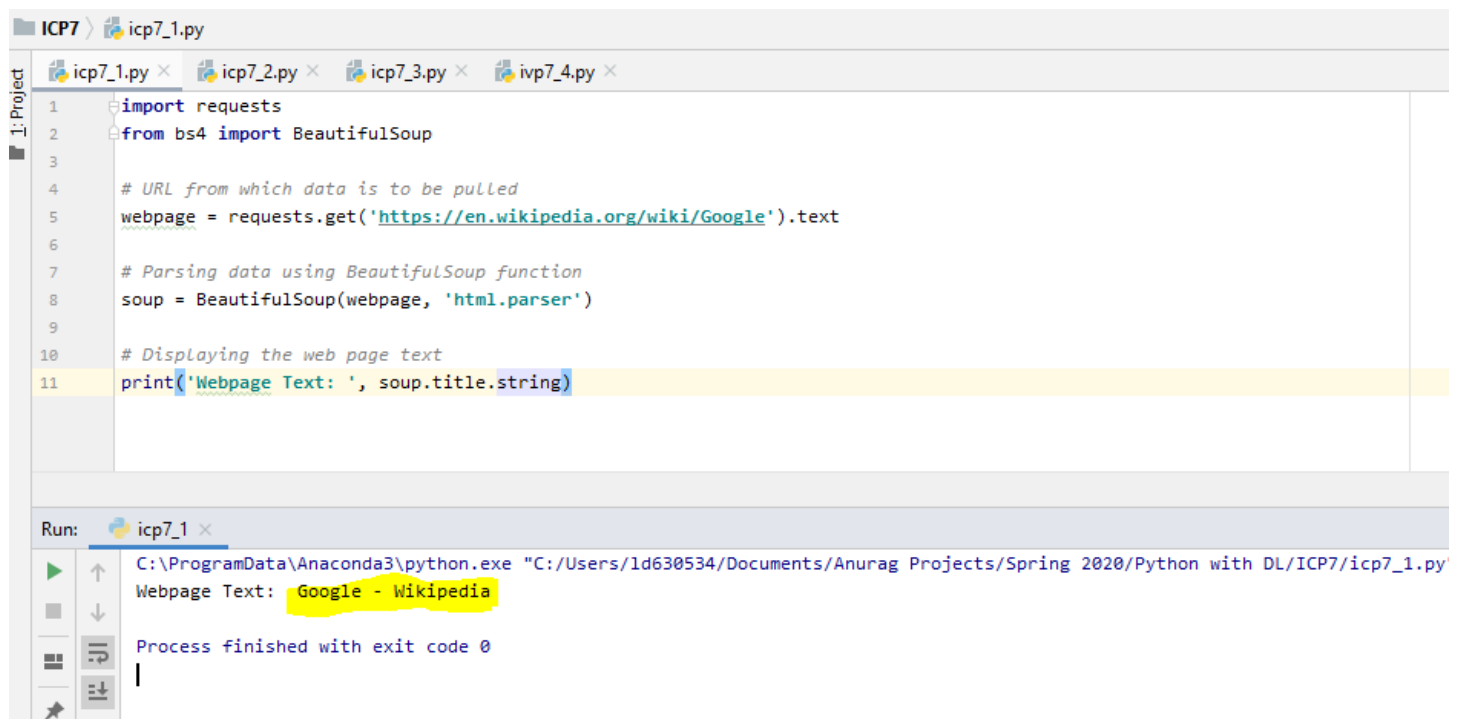


Video Link: <https://www.loom.com/share/6df8d95ff1584195b06569aabe5dac7e>

1. Extract the following web URL text using BeautifulSoup <https://en.wikipedia.org/wiki/Google2>.
2. Save it in input.txt
3. Apply the following on the text and show output:
 - a. Tokenization
 - b. POS
 - c. Stemming
 - d. Lemmatization
 - e. Trigram
 - f. Named Entity Recognition
4. Change the classifier in the given code to:
 - a. KNeighborsClassifier and see how accuracy changes
 - b. change the tfidfvectorizer to use bigram and see how the accuracy changes
TfidfVectorizer(ngram_range=(1,2))
 - c. Put argument stop_words='english' and see how accuracy changes



The screenshot shows a Jupyter Notebook interface with a file explorer on the left displaying a project named 'ICP7'. The notebook has four tabs: 'icp7_1.py', 'icp7_2.py', 'icp7_3.py', and 'icp7_4.py'. The active tab is 'icp7_1.py', which contains the following Python code:

```
1 import requests
2 from bs4 import BeautifulSoup
3
4 # URL from which data is to be pulled
5 webpage = requests.get('https://en.wikipedia.org/wiki/Google').text
6
7 # Parsing data using BeautifulSoup function
8 soup = BeautifulSoup(webpage, 'html.parser')
9
10 # Displaying the web page text
11 print('Webpage Text: ', soup.title.string)
```

Below the code editor, the 'Run' button is clicked, and the terminal output is displayed:

```
Run: icp7_1 x
C:\ProgramData\Anaconda3\python.exe "C:/Users/ld630534/Documents/Anurag Projects/Spring 2020/Python with DL/ICP7/icp7_1.py"
Webpage Text: Google - Wikipedia
Process finished with exit code 0
```

```
ICP7 > icp7_2.py
1 import requests
2 from bs4 import BeautifulSoup
3
4 # URL from which data is to be pulled
5 webpage = requests.get('https://en.wikipedia.org/wiki/Google').text
6
7 # Parsing data using BeautifulSoup function
8 soup = BeautifulSoup(webpage, 'html.parser')
9
10 # Saving the parsed webpage data into the text file titled 'input'
11 text = soup.get_text()
12 f = open('input.txt', 'w', encoding='utf-8')
13 f.write(text)
```

Run: icp7_2 ×

C:\ProgramData\Anaconda3\python.exe "C:/Users/ld630534/Documents/Anurag Projects/Spring 2020/Python with DL/ICP7/icp7_2.py"

Process finished with exit code 0

Input - to.txt
in Github

```
ICP7 > icp7_3.py
1 import nltk
2 from nltk.stem import PorterStemmer
3 from nltk.stem import LancasterStemmer
4 from nltk.stem import SnowballStemmer
5 from nltk.stem import WordNetLemmatizer
6 from nltk import wordpunct_tokenize, pos_tag, ne_chunk
7 from nltk import ngrams
8
9 test_text = open('input.txt', encoding="utf8").read()
10
11 # a) Tokenization
12 token = nltk.word_tokenize(test_text)
13 print('Tokens identified are', token)
14
15
16 # b) Part Of Speech tagging
17 pos = nltk.pos_tag(token)
18 print('Part of Speech associated with the input text', pos)
19
20
21 # c) Stemming - identifying the root or base word of the terms associated
22 pStemmer = PorterStemmer() #Porter Stemming keeps only prefix for each words and leave non English words like tr
23 for x in token:
24     print('Result of Stemming using PorterStemmer for ', x, 'is ', pStemmer.stem(x))
25
26 lStemmer = LancasterStemmer() #Lancaster stemming is a rule-based stemming based on the last letter of the words
27 for y in token:
28     print('Result of Stemming using LancasterStemmer for ', y, 'is ', lStemmer.stem(y))
29
30 sStemmer = SnowballStemmer('english')
31 for z in token:
32     print('Result of Stemming using SnowballStemmer for ', z, 'is ', sStemmer.stem(z))
33
```

```

# d) Lemmatization - normalization of text based on the meaning as part of the speech (converts plurals or adjective to
# their basic, meaningful singular form)
lemmatizer = WordNetLemmatizer()
print('Result of Lemmatization: ', lemmatizer.lemmatize(x))

# e) Trigram
trigram = ngrams(test_text.split(), 3)
for gram in trigram:
    print('Trigram data is ', gram)
print(str(trigram))

# f) Named Entity Recognition
print('Named Entity Recognition is ', ne_chunk(pos_tag(wordpunct_tokenize(test_text))))

```

ICP7 > ivp7_4.py

icp7_1.py x icp7_2.py x icp7_3.py x ivp7_4.py x

```

1  # %%
2  from sklearn.datasets import fetch_20newsgroups
3  from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
4  from sklearn.feature_extraction.text import TfidfTransformer
5  from sklearn import metrics
6  from sklearn.pipeline import Pipeline
7  from sklearn.naive_bayes import MultinomialNB
8  from sklearn.neighbors import KNeighborsClassifier
9  from nltk.corpus import stopwords
10
11  # %%
12  twenty_train = fetch_20newsgroups(subset='train', shuffle=True)
13
14  # %%
15  tfidf_Vect = TfidfVectorizer()
16  X_train_tfidf = tfidf_Vect.fit_transform(twenty_train.data)
17
18  # print(tfidf_Vect.vocabulary_)
19  clf = MultinomialNB()
20  clf.fit(X_train_tfidf, twenty_train.target)
21
22  # %%
23  twenty_test = fetch_20newsgroups(subset='test', shuffle=True)
24  X_test_tfidf = tfidf_Vect.transform(twenty_test.data)
25
26  # %%
27  predicted = clf.predict(X_test_tfidf)
28
29  # 0/0

```

1: Project

Structure

es

```
icp7_1.py × icp7_2.py × icp7_3.py × ivp7_4.py ×
31     print(score)
32
33     # Creating matrix from the data available (provided in Lecture)
34     countVec = CountVectorizer()
35     x_train = countVec.fit_transform(twenty_train.data)
36     x_test = countVec.transform(twenty_test.data)
37
38     # Multinomial NB Model building on training data
39     clfnew = MultinomialNB()
40     clfnew.fit(x_train, twenty_train.target)
41
42     # Evaluating the model fit on test data set
43     clfPredcit = clfnew.predict(x_test)
44
45     # Computing the score of the Multinomial NB model
46     clfScore = metrics.accuracy_score(twenty_test.target, clfPredcit)
47     print('The score for Multinomial NB model is', clfScore)
48
49     # a) Building K-Nearest Neighbours Model on training data set
50     knnModel = KNeighborsClassifier()
51     knnModel.fit(x_train, twenty_train.target)
52
53     # Predicting the model fit on test data set
54     knnPredict = knnModel.predict(x_test)
55
56     # Computing the score of the KNN model
57     knnScore = metrics.accuracy_score(twenty_test.target, knnPredict)
58     print('The score for KNN model is', knnScore)
```

```
59
60 # b) TFIDF vectorization updated to use bigrams as asked in the question
61 bigramnewVec = TfidfVectorizer(ngram_range=(1, 2))
62 x_train_bigram = bigramnewVec.fit_transform(twenty_train.data)
63 x_test_bigram = bigramnewVec.transform(twenty_test.data)
64
65 # Revised models and their accuracy using bigram vectorization of data
66 # Multinomial NB
67 clfBigram = clfnew.fit(x_train_bigram, twenty_train.target)
68 clfPredcitBigram = clfnew.predict(x_test_bigram)
69 clfScoreBigram = metrics.accuracy_score(twenty_test.target, clfPredcitBigram)
70 print('The score for Multinomial NB model after TFIDF vectorization bigram update is', clfScoreBigram)
71
72 # K-Nearest Neighbours
73 knnBigram = knnModel.fit(x_train_bigram, twenty_train.target)
74 knnPredictBigram = knnModel.predict(x_test_bigram)
75 knnScoreBigram = metrics.accuracy_score(twenty_test.target, knnPredictBigram)
76 print('The score for KNN model is', knnScoreBigram)
77
78 # Stop word
79 # TFIDF vectorization updated to accommodate for stop word 'english'
80 stopwordVec = TfidfVectorizer(stop_words='english')
81 x_train_stopWord = stopwordVec.fit_transform(twenty_train.data)
82 x_test_stopWord = stopwordVec.transform(twenty_test.data)
83
84 # Revised models and their accuracy when stopword english is applied on data
85 # Multinomial NB
86 clfStopWord = clfnew.fit(x_train_stopWord, twenty_train.target)
87 clfPredcitStopWord = clfnew.predict(x_test_stopWord)
```

```

# Stop word
# TFIDF vectorization updated to accommodate for stop word 'english'
stopwordVec = TfidfVectorizer(stop_words='english')
x_train_stopWord = stopwordVec.fit_transform(twenty_train.data)
x_test_stopWord = stopwordVec.transform(twenty_test.data)

# Revised models and their accuracy when stopword english is applied on data
# Multinomial NB
clfStopWord = clfnew.fit(x_train_stopWord, twenty_train.target)
clfPredictStopWord = clfnew.predict(x_test_stopWord)
clfScoreStopWord = metrics.accuracy_score(twenty_test.target, clfPredictStopWord)
print('The score for Multinomial NB model after stopword english update is', clfScoreStopWord)

# K-Nearest Neighbours
knnStopWord = knnModel.fit(x_train_stopWord, twenty_train.target)
knnPredictStopWord = knnModel.predict(x_test_stopWord)
knnScoreStopWord = metrics.accuracy_score(twenty_test.target, knnPredictStopWord)
print('The score for KNN model is after stop word english is updated is', knnScoreStopWord)

```

ivp7_4 x

```

C:\ProgramData\Anaconda3\python.exe C:/Users/10050534/Documents/Anurag-Projects/Spring-2020/Python-with-DL/ICF7/ivp7_4.py
0.7738980350504514
The score for Multinomial NB model is 0.7728359001593202
The score for KNN model is 0.3524960169941583
The score for Multinomial NB model after TFIDF vectorization bigram update is 0.765400955921402
The score for KNN model is 0.6196229421136484
The score for Multinomial NB model after stopword english update is 0.8169144981412639
The score for KNN model is after stop word english is updated is 0.6757833244822092

```