

## Outstanding Project-3

By:

Name: Jatin Jindal

E-Mail Id: [jatinjindal1199@gmail.com](mailto:jatinjindal1199@gmail.com)

## 1. Scraping News Headlines

## The Hindu

```
[78] 1 page=20
2 for i in range(1,page+1):
3     count=50
4     ch='https://www.thehindu.com/archive/web/2020/9/'+str(i)+'/'
5     res=requests.get(ch)
6     soup=bs4.BeautifulSoup(res.text,'lxml')
7     for j in soup.select('.archive-list'):
8         if count==0:
9             break
10        else:
11            #print('2020/09/'+str(i))
12            headline['Date'].append('2020/09/'+str(i))
13            headline['News_Channel'].append('The Hindu')
14            #print(j.text)
15            news=j.text.replace('\n',"")
16            headline['Headline'].append(news)
17            count=count-1
18 print(headline)
19 len(headline['Headline'])
```

## DD News

```
[79] 1 page=20
      2 for i in range(1,page+1):
      3     count=50
      4     ch='http://ddnews.gov.in/about/news-archive?page='+str(i)
      5     res=requests.get(ch)
      6     soup=bs4.BeautifulSoup(res.text,'lxml')
      7     for j in soup.select('.archive-title'):
      8         if count==0:
      9             break
     10         else:
     11             #print('2020/09/'+str(i))
     12             headline['Date'].append('2020/09/'+str(i))
     13             headline['News_Channel'].append('DD News')
     14             #print(j.text)
     15             news=j.text.replace('\n',"")
     16             headline['Headline'].append(news)
     17             count=count-1
     18 print(headline)
     19 len(headline['Headline'])
```

Indian Express

[illegible]

## 2. Finding Positive and Negative Sentiments

```
[81] 1 headline['Negative']=[]
      2 headline['Positive']=[]
      3 headline['Neutral']=[]
      4 headline['Total']=[]
      5 headline['Sentiment']=[]
      6 k=0
      7 from nltk.sentiment.vader import SentimentIntensityAnalyzer
      8 sia = SentimentIntensityAnalyzer()
      9 for line in headline['Headline']:
10     pol_score=sia.polarity_scores(line)
11     headline['Negative'].append(pol_score['neg'])
12     headline['Positive'].append(pol_score['pos'])
13     headline['Neutral'].append(pol_score['neu'])
14     headline['Total'].append(pol_score['compound'])
15     if pol_score['compound']> 0.2:
16         headline['Sentiment'].append(1)
17     elif pol_score['compound']> -0.2:
18         headline['Sentiment'].append(-1)
19     else:
20         headline['Sentiment'].append(0)
21 df=pd.DataFrame(headline)
22 df
```

	Date	News_Channel	Headline	Negative	Positive	Neutral	Total	Sentiment
0	2020/09/1	The Hindu	US Open   Fluent start for Djokovic; Osaka str...	0.000	0.000	1.000	0.0000	-1
1	2020/09/1	The Hindu	Facebook wants you to experience its virtual u...	0.043	0.048	0.909	-0.1280	-1
2	2020/09/1	The Hindu	Top news of the day: Former President Pranab M...	0.042	0.261	0.697	0.8020	1
3	2020/09/1	The Hindu	Indian Americans should get involved in U.S. P...	0.128	0.151	0.721	0.3818	

If the overall is less between -0.2 to 0.2 then it is treated as neutral

If the overall is less than -0.2 then it is treated as negative

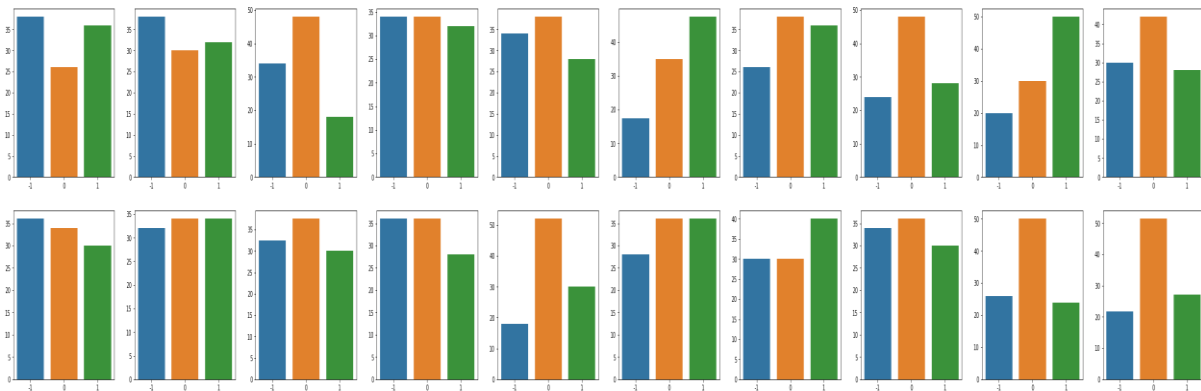
If the overall is more than 0.2 then it is treated as Positive

### 3. Date Wise Sentiment%

-1:Negative(Blue)  
1: Positive(Orange)  
0:Neutral(Green)

#### The Hindu

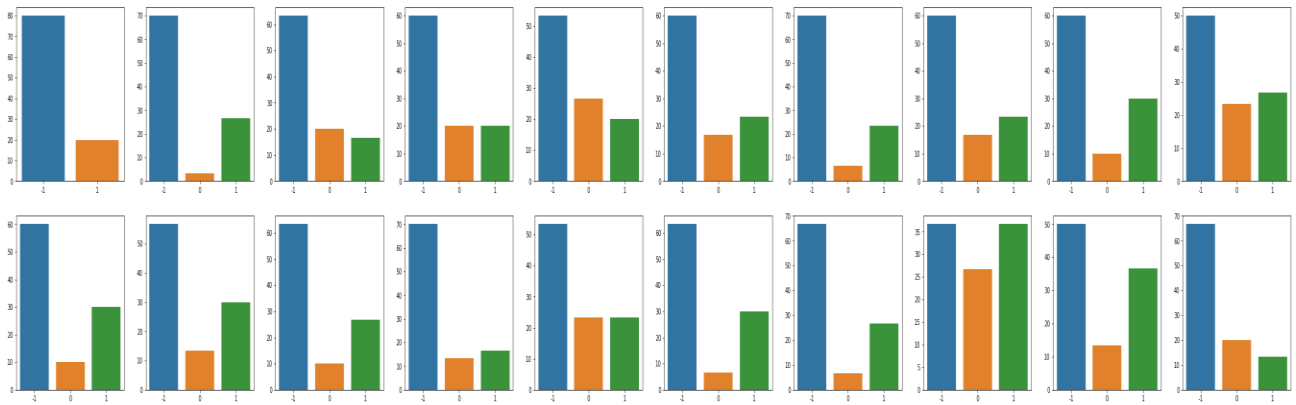
```
[120] 1 print('The Hindu')
      2 plt.figure(figsize=(50,50))
      3 for i in range(1,page+1):
      4     plt.subplot(10,10,i)
      5     df1=df[df['Date']=='2020/09/'+str(i)]
      6     percent=df1[df1['News_Channel']=='The Hindu'].Sentiment.value_counts()/df1[df1['News_Channel']=='The Hindu'].Sentiment.value_counts().sum()*100
      7     sns.barplot(x=percent.index,y=percent.values)
```



We can easily see that there is mostly a good mixture all type sentiments news on each day except on some days.

#### DD News

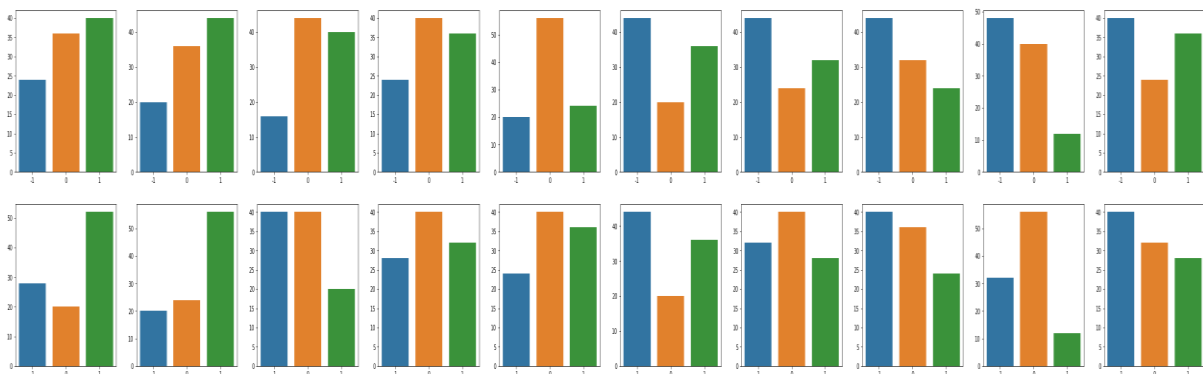
```
1 print('DD News')
2 plt.figure(figsize=(50,50))
3 for i in range(1,page+1):
4     plt.subplot(10,10,i)
5     df1=df[df['Date']=='2020/09/'+str(i)]
6     percent=df1[df1['News_Channel']=='DD News'].Sentiment.value_counts()/df1[df1['News_Channel']=='DD News'].Sentiment.value_counts().sum()*100
7     sns.barplot(x=percent.index,y=percent.values)
```



We can see that the Most Percent of their news are showing negative sentiment.

### Indian Express

```
print('Indian Express')
plt.figure(figsize=(50,50))
for i in range(1,page+1):
    plt.subplot(10,10,i)
    df1=df[df['Date']=='2020/09/'+str(i)]
    percent=df1[df1['News_Channel']=='Indian Express'].Sentiment.value_counts()/df1[df1['News_Channel']=='Indian Express'].Sentiment.value_counts().sum()*100
    sns.barplot(x=percent.index,y=percent.values)
```



It can be seen as some days there it has high percent of negative news and some days it has high negative sentiment and some days high positive sentiment.

## 4. Cleaning News Headlines

```
1 import re
2 nltk.download('stopwords')
3 from nltk.corpus import stopwords
4 from nltk.stem.porter import PorterStemmer
5 clean_headlines=[]
6 for i in range(0,2067):
7     head=re.sub('[^a-zA-Z]', ' ',df['Headline'][i]) # to include only english words
8     head=head.lower()
9     head=head.split() # to get each word seperated from each other
10    ps=PorterStemmer()
11    all_stopwords=stopwords.words('english') # Used to get all stopwords in english language
12    all_stopwords.remove('not')
13    head=[ps.stem(word) for word in head if not word in set(all_stopwords)]
14    # Removing the stopwords and at the same time stem the words
15    head=' '.join(head)
16    clean_headlines.append(head)
```

## 5. Creating Bag of words

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 cv=CountVectorizer()
3 x=cv.fit_transform(df['Clean Headline'])
```

## 6. Splitting

```
1 from sklearn.model_selection import train_test_split
2 x_train,x_test,y_train,y_test=train_test_split(x,df['Sentiment'],random_state=0)
3 print(x_train.shape)
4 print(y_train.shape)
5 print(x_test.shape)
6 print(y_test.shape)
```

```
(1550, 12752)
(1550,)
(517, 12752)
(517,)
```

## 7. Model Training

### a) Naive Bayes

```
9] 1 from sklearn.naive_bayes import MultinomialNB
    2 nb=MultinomialNB()
    3 nb.fit(x_train,y_train)
    4 nb.score(x_test,y_test)
```

```
→ 0.5087040618955513
```

```
1 from sklearn.model_selection import cross_val_score
2 accuracy_nb=cross_val_score(estimator=nb,X=x_train,y=y_train,cv=10)
3 print('Accuracy= {:f} and Standard Deviation= {:f}'.format(accuracy_nb.mean()*100,accuracy_nb.std()*100))
4 accuracy_nb
```

```
Accuracy= 55.161290 and Standard Deviation= 3.675155
array([0.63870968, 0.59354839, 0.52258065, 0.52903226, 0.53548387,
       0.55483871, 0.52903226, 0.52903226, 0.56774194, 0.51612903])
```

### b) K Nearest Neighbors

```
1 from sklearn.neighbors import KNeighborsClassifier
2 kn=KNeighborsClassifier(n_neighbors=5,metric='minkowski',p=2)
3 kn.fit(x_train,y_train)
4 kn.score(x_test,y_test)
```

```
0.38684719535783363
```

```
1 from sklearn.model_selection import cross_val_score
2 accuracy_kn=cross_val_score(estimator=kn,X=x_train,y=y_train,cv=10)
3 print('Accuracy= {:f} and Standard Deviation= {:f}'.format(accuracy_kn.mean()*100,accuracy_kn.std()*100))
4 accuracy_kn
```

```
Accuracy= 41.290323 and Standard Deviation= 0.816072
array([0.41290323, 0.40645161, 0.41290323, 0.41290323, 0.41290323,
       0.41935484, 0.40645161, 0.42580645, 0.42580645, 0.40645161])
```

### C) Random Forest

```
1 from sklearn.ensemble import RandomForestClassifier
2 forest=RandomForestClassifier(n_estimators=10,criterion='entropy',random_state=0)
3 forest.fit(x_train,y_train)
4 forest.score(x_test,y_test)
```

0.5764023210831721

```
1 from sklearn.model_selection import cross_val_score
2 accuracy_forest=cross_val_score(estimator=forest,x=x_train,y=y_train,cv=10)
3 print('Accuracy= {:f} and Standard Deviation= {:f}'.format(accuracy_forest.mean()*100,accuracy_forest.std()*100))
4 accuracy_forest
```

Accuracy= 61.419355 and Standard Deviation= 4.686199  
array([0.64516129, 0.66451613, 0.60645161, 0.61935484, 0.64516129,  
 0.59354839, 0.56774194, 0.56129032, 0.69677419, 0.54193548])

### d) Decision Tree

```
1 from sklearn.tree import DecisionTreeClassifier
2 dtree=DecisionTreeClassifier(criterion='entropy',random_state=0)
3 dtree.fit(x_train,y_train)
4 dtree.score(x_test,y_test)
```

0.5880077369439072

```
1 from sklearn.model_selection import cross_val_score
2 accuracy_dtree=cross_val_score(estimator=dtree,x=x_train,y=y_train,cv=10)
3 print('Accuracy= {:f} and Standard Deviation= {:f}'.format(accuracy_dtree.mean()*100,accuracy_dtree.std()*100))
4 accuracy_dtree
```

Accuracy= 60.451613 and Standard Deviation= 3.054142  
array([0.67741935, 0.58709677, 0.59354839, 0.6 , 0.58064516,  
 0.55483871, 0.6 , 0.61935484, 0.61935484, 0.61290323])

Therefore the best model is Random Forest Classification as the accuracy of this one is highest.

## 8. Finding Sentiment for new news channel

### a) Scrapping News

```
1 count=4
2 lst=[]
3 ch='https://timesofindia.indiatimes.com/home/headlines'
4 res=requests.get(ch)
5 soup=bs4.BeautifulSoup(res.text,'lxml')
6 for j in soup.select('.w_tle'):
7     if count==0:
8         break
9     else:
10        #print(j.text)
11        news=j.text.replace('\n',"")
12        lst.append(news)
13        count=count-1
14 lst
```

```
["Live: Won't contest polls with JD(U) due to ideological reasons, says LJP",
"Live: '500m vaccine doses for 25cr by July 2021'",
"Live: Bhim Army chief meets Hathras victim's kin",
'LAC row: Army, IAF prepare to fight wars jointly']
```

### b) Cleaning News Headlines

```
1 for i in range(0,len(test_lst)):
2     line=lst[i]
3     head=re.sub('[^a-zA-Z]', ' ',line) # to include only english words
4     head=head.lower()
5     head=head.split() # to get each word seperated from each other
6     ps=PorterStemmer()
7     all_stopwords=stopwords.words('english') # Used to get all stopwords in eng
8     all_stopwords.remove('not')
9     head=[ps.stem(word) for word in head if not word in set(all_stopwords)]
10    # Removing the stopwords and at the same time stem the words
11    head=' '.join(head)
12    lst[i]=head
13 lst
```

```
['live contest poll jd u due ideolog reason say ljp',
'live vaccin dose cr juli',
'live bhim armi chief meet hathra victim kin',
'lac row armi iaf prepar fight war jointli']
```



### c) Making Bag of words and Predicting Sentiment

```
[ ] 1 sample=cv.transform(lst).toarray()  
2 sample
```

```
[ ] array([[0, 0, 0, ..., 0, 0, 0],  
          [0, 0, 0, ..., 0, 0, 0],  
          [0, 0, 0, ..., 0, 0, 0],  
          [0, 0, 0, ..., 0, 0, 0]])
```

```
[243] 1 forest.predict(sample)
```

```
[ ] array([-1, -1,  0,  0])
```