

Task 4- GRIP at Sparks Foundation

Stock Market Prediction using Numerical and Textual Analysis*

Create a hybrid model for stock price/performance prediction using numerical analysis of historical stock prices, and sentimental analysis of news headlines - Stock to analyze and predict - SENSEX (S&P BSE SENSEX)

Sentimental Analysis of News headlines

```
In [1]: # pip install TextBlob
```

```
In [2]: # Importing libraries in Python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [3]: df = pd.read_csv('D://DK//dddddd//Task_grip//TASK _3//task 4//india-new
s-headlines.csv')
```

```
In [4]: df.head()
```

Out[4]:

	publish_date	headline_category	headline_text
0	20010101	sports.wwe	win over cena satisfying but defeating underta...

	publish_date	headline_category	headline_text
1	20010102	unknown	Status quo will not be disturbed at Ayodhya; s...
2	20010102	unknown	Fissures in Hurriyat over Pak visit
3	20010102	unknown	America's unwanted heading for India?
4	20010102	unknown	For bigwigs; it is destination Goa

In [5]: `df.shape`

Out[5]: (3297172, 3)

As we can see in headline_category column there are some cities name, let's separate this cities name

In [6]: `df['headline_category'].value_counts()`

```
Out[6]: india                285619
unknown                207732
city.mumbai            132649
city.delhi              124658
business.india-business 115246
city.chandigarh         107464
city.hyderabad           99014
city.bengaluru           91857
entertainment.hindi.bollywood 90374
city.ahmedabad           85813
city.pune                84620
city.lucknow             83591
city.kolkata             81298
city.nagpur              77300
city.goa                 76409
city.chennai             72264
city.patna               70386
city.jaipur              53695
sports.icc-world-cup-2015 40537
business.international-business 36356
```

city.gurgaon	31312
citizen.stories	29543
city.bhubaneswar	29098
city.kochi	27861
tv.news.hindi	27791
home.campaigns	26834
entertainment.english.hollywood	25748
city.bhopal	24484
city.thiruvananthapuram	24047
city.guwahati	24017
...	...
times-fact-check.news	8
pms-us-visit	8
business.faqs.gst-faqs	8
best-products.beauty.grooming	8
elections.assembly-elections.goa	8
lifespan-news	8
life-style.parenting.getting-pregnant	8
himachal-pradesh	8
scorecard-and-statistics	8
indias-vision	8
best-products.home-decor-and-garden.living-room-decor	8
advanis-us-visit	8
year-ender-2015.march	8
sports.hockey.hockey-india-league.interviews	8
iim-fee-row	8
preeti-shenoy	8
editorialt	8
sports.cricket.bangladesh-in-west-indies	8
brandwire.media-entertainment.newspapers-magazines-books	8
entertainment.hindi.music.singer-of-the-week	8
profiles.india-profiles	8
sports.asian-games-2018	8
nepal-india-earthquake.opinion	8
ballot-talk	7
delhi-ncr	7
sports.headline3	7
actresses	7
did-you-know	6

```
sports.headline6      3
party-manifestos      2
Name: headline_category, Length: 1016, dtype: int64
```

Data Cleaning

```
In [7]: city_df = df[df['headline_category'].str.contains('^city\[a-z]+\$', reg
ex=True)]
```

```
In [8]: city_df
```

Out[8]:

	publish_date	headline_category	headline_text
274	20010104	city.bengaluru	Three in race for chief secy's post
275	20010104	city.patna	Druggists' stir leads to shortage of medicines
278	20010104	city.bengaluru	He's not so inscrutable
279	20010104	city.delhi	DPCC stages Nyay rally
643	20010110	city.patna	Fend for yourselves; Pande tells doctors
644	20010110	city.patna	Bureaucracy undermining legislature's 'existence'
645	20010110	city.patna	State police collapses under pressure
646	20010110	city.patna	Court declares Pappu Yadav an absconder
2195	20010309	city.delhi	Maneka asks govt to take back land for cow she...
3559	20010423	city.ahmedabad	Killer was promised Rs 25,000 for killing Nair...
4425	20010518	city.thiruvananthapuram	Antony is new CM

	publish_date	headline_category	headline_text
4428	20010518	city.bengaluru	Don't take that biscuit; you dope
4429	20010518	city.bengaluru	'I've done my bit when it comes to preparing'
4436	20010518	city.bengaluru	He is etched in the chapters of Bangalore
4441	20010518	city.bengaluru	Old homes live again as pubs; shops; restaurants
4444	20010518	city.bengaluru	CET students see a kind Karnataka
4453	20010518	city.thiruvananthapuram	Cong 'groups' oppose Muraleedharan's leadership
4461	20010518	city.thiruvananthapuram	I have a Herculean task ahead: Antony
4463	20010518	city.bengaluru	CET results to be announced on June 10
4464	20010518	city.bengaluru	Two terrific cuisines
4465	20010518	city.thiruvananthapuram	Strict financial discipline need of the hour'
4466	20010518	city.bengaluru	2 arrested for selling fake CET question papers
4468	20010518	city.bengaluru	Task forces to check pollution
4472	20010518	city.bengaluru	Petition challenging budget proposal adjourned
4474	20010518	city.bengaluru	Krishna waters: State plans contempt plea agai...
4476	20010518	city.bengaluru	He believes in making customers happy
4478	20010518	city.bengaluru	Cottonpet: The bustle of the bazaar
4479	20010518	city.bengaluru	Second day hassle-free
4480	20010518	city.bengaluru	Will for sale at just Rs 15,000

	publish_date	headline_category	headline_text
4484	20010518	city.bengaluru	Krishna will not meet M'rasht CM over border...
...
3297106	20200630	city.hyderabad	Hyderabad: Two die of Covid-19; their videos p...
3297107	20200630	city.hyderabad	First in India: Telangana HC takes court to la...
3297108	20200630	city.goa	Cooperative banks get govt approval to offer g...
3297109	20200630	city.hyderabad	Hyderabad: Wife's murder accused out on bail k...
3297110	20200630	city.goa	Goa: Police oppose anticipatory bail plea by Z...
3297111	20200630	city.chennai	tamil nadu custodial deaths cctv footage shows...
3297112	20200630	city.goa	Dongrim-Neura locals begin community farming
3297113	20200630	city.goa	Congress wants white paper on Goa's economic s...
3297114	20200630	city.ranchi	BCCL confident of meeting production target: CMD
3297116	20200630	city.mumbai	Maharashtra: Use Marathi or risk pay hike free...
3297117	20200630	city.jamshedpur	Tent owners; caterers want gatherings curbs ea...

	publish_date	headline_category	headline_text
3297118	20200630	city.mumbai	Mumbai cops target 'illegal' travellers; seize...
3297119	20200630	city.ranchi	Plantation in Dalma to reduce man-elephant con...
3297120	20200630	city.ranchi	BJP meets Guv with CBI probe plea into death o...
3297121	20200630	city.ranchi	When admn is appealing to follow lockdown norm...
3297127	20200630	city.chennai	Internal policing a must for checks and balances
3297128	20200630	city.chennai	How the Madras Magician wowed SW19 60 years ago
3297131	20200630	city.amaravati	Andhra Pradesh: Rise in caseload; Covid care c...
3297132	20200630	city.vadodara	Hotspots changing: Surat sees rise in Covid-19...
3297133	20200630	city.surat	Hotspots changing: Surat sees rise in Covid-19...
3297134	20200630	city.jodhpur	Record wheat procurement at MSP this year desp...
3297135	20200630	city.udaipur	Record wheat procurement at MSP this year desp...
3297136	20200630	city.nagpur	14 new cases in Chandrapur
3297137	20200630	city.nagpur	Aurangabad returnee tests positive
3297138	20200630	city.nagpur	One tola gold & 1kg silver create price ratio ...

	publish_date	headline_category	headline_text
3297139	20200630	city.delhi	Remove illegal signage or face penalty: NDMC
3297158	20200630	city.bengaluru	what bengaluru can do to tackle covid surge
3297161	20200630	city.bengaluru	karnataka may adopt keralas triple lockdown plan
3297162	20200630	city.kanpur	vehicle of up stf team bringing gangster vikas...
3297166	20200630	city.lucknow	up govt imposes weekend restrictions from tonight

1842288 rows × 3 columns

*create a new dataframe which have only one column city_name

```
In [9]: city_split = pd.DataFrame(columns=['city_name'])
```

```
In [10]: city_split['city_name'] = city_df.headline_category.str.split('.', expand=True)[1]
```

```
In [11]: # concat the data set city_data and city_split
city_df= pd.concat([city_df,city_split],axis=1)
```

```
In [12]: # removing unnecessary column from dataset
city_df.drop(columns=['headline_category'],axis=1, inplace=True)
```

```
In [13]: # re-setting the index
city_df.reset_index(inplace=True)
```

```
In [14]: # see the final data which contain city_name column
city_df.head()
```


Out[14]:

	index	publish_date	headline_text	city_name
0	274	20010104	Three in race for chief secy's post	bengaluru
1	275	20010104	Druggists' stir leads to shortage of medicines	patna
2	278	20010104	He's not so inscrutable	bengaluru
3	279	20010104	DPCC stages Nyay rally	delhi
4	643	20010110	Fend for yourselves; Pande tells doctors	patna

In [15]: `city_df.city_name.unique()`

Out[15]: array(['bengaluru', 'patna', 'delhi', 'ahmedabad', 'thiruvananthapura
m',
'pune', 'mumbai', 'chandigarh', 'lucknow', 'kolkata', 'hyderaba
d',
'chennai', 'bareilly', 'aurangabad', 'nagpur', 'bhopal',
'vadodara', 'jaipur', 'goa', 'thane', 'hubballi', 'mangaluru',
'mysuru', 'rajkot', 'surat', 'kanpur', 'varanasi', 'allahabad',
'ludhiana', 'guwahati', 'bhubaneswar', 'ranchi', 'gurgaon',
'coimbatore', 'noida', 'madurai', 'indore', 'kochi', 'kozhikod
e',
'nashik', 'raipur', 'visakhapatnam', 'kolhapur', 'navimumbai',
'trichy', 'puducherry', 'dehradun', 'meerut', 'agra', 'vijayawad
a',
'jamshedpur', 'imphal', 'shillong', 'amritsar', 'shimla',
'cuttack', 'jind', 'agartala', 'jammu', 'faridabad', 'srinagar',
'salem', 'rajahmundry', 'erode', 'ghaziabad', 'itanagar', 'ajme
r',
'kohima', 'gaya', 'jodhpur', 'udaipur', 'amaravati'], dtype=obje
ct)

In [18]: `# create a group of cities and let's see which city have, how many no.
of headline text
city_headline = city_df.groupby(['city_name']).agg({'headline_text': 'co
unt'})`

```
In [19]: # for better understanding remane the column
city_headline.rename(columns={'headline_text':'No._headline_text'},inplace=True)
```

```
In [20]: # finding the top 15 cities which have high number of headlines text
city_headline = city_headline.sort_values(by='No._headline_text',ascending=False)
top_15_headline_city = city_headline.head(15)
top_15_headline_city
```

Out[20]:

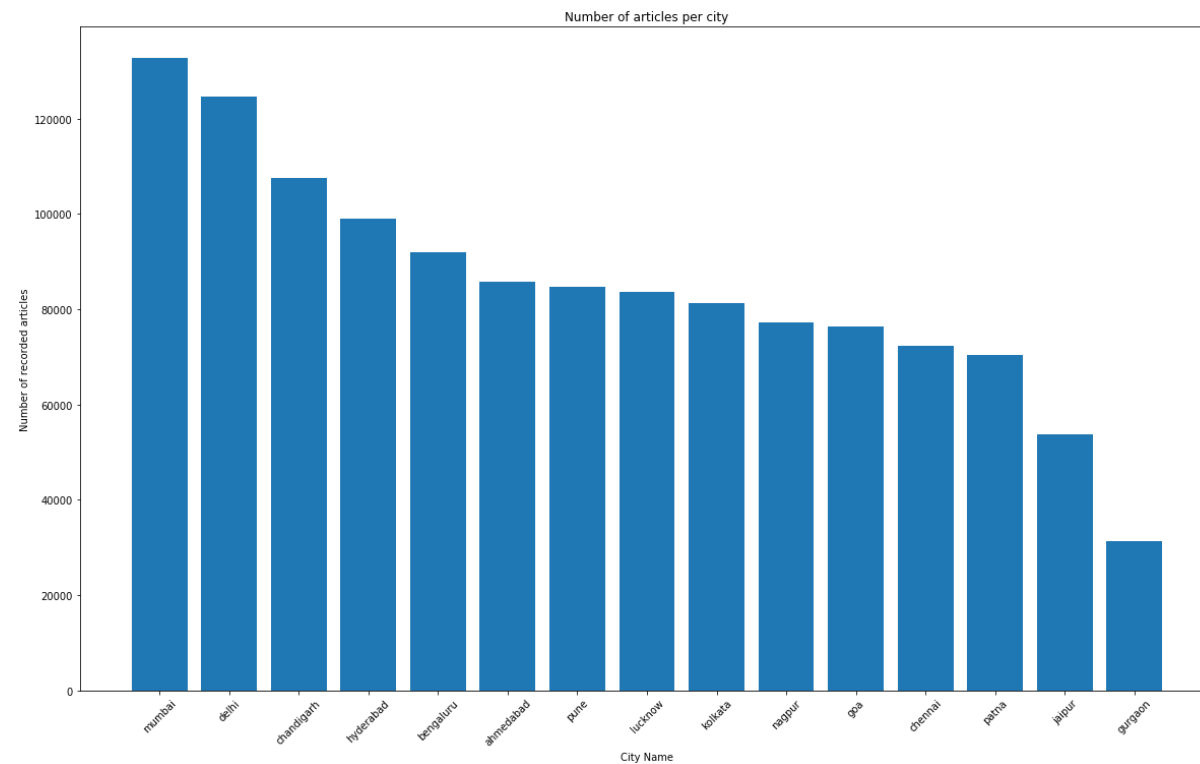
	No._headline_text
city_name	
mumbai	132649
delhi	124658
chandigarh	107464
hyderabad	99014
bengaluru	91857
ahmedabad	85813
pune	84620
lucknow	83591
kolkata	81298
nagpur	77300
goa	76409
chennai	72264
patna	70386
jaipur	53695
aurangabad	31312

gurgaon

0.012

```
In [21]: # plotting a graph for better understanding for top 15 cities which have high no. of headline text

plt.figure(figsize=(20,12))
plt.bar(top_15_headline_city.index,height=top_15_headline_city['No._headline_text'])
plt.xticks(rotation = 45)
plt.xlabel('City Name')
plt.ylabel('Number of recorded articles')
plt.title('Number of articles per city')
plt.show()
```



```
In [22]: city_df.dtypes
# publish data have "int" datatype we have to change that into "datetime"
```

```
e"
```

```
Out[22]: index          int64
publish_date      int64
headline_text     object
city_name         object
dtype: object
```

```
In [23]: city_df['publish_date'] = pd.to_datetime(city_df['publish_date'])
city_df.head()
```

```
Out[23]:
```

	index	publish_date	headline_text	city_name
0	274	1970-01-01 00:00:00.020010104	Three in race for chief secy's post	bengaluru
1	275	1970-01-01 00:00:00.020010104	Druggists' stir leads to shortage of medicines	patna
2	278	1970-01-01 00:00:00.020010104	He's not so inscrutable	bengaluru
3	279	1970-01-01 00:00:00.020010104	DPCC stages Nyay rally	delhi
4	643	1970-01-01 00:00:00.020010110	Fend for yourselves; Pande tells doctors	patna

```
In [24]: city_df['Year'] = city_df['publish_date'].apply(lambda x: (x.microsecond)//10)
```

```
In [25]: # dropping unnecessary columns
city_df.drop(columns=['publish_date', 'index'], axis=1, inplace=True)
```

```
In [26]: city_df.head()
```

```
Out[26]:
```

	headline_text	city_name	Year
0	Three in race for chief secy's post	bengaluru	2001
1	Druggists' stir leads to shortage of medicines	patna	2001
2	He's not so inscrutable	bengaluru	2001
3	DPCC stages Nyay rally	delhi	2001
4	Fend for yourselves; Pande tells doctors	patna	2001

```
In [27]: city_df_2 = city_df.copy()
city_df_2['No_headline_text'] = 1
```

```
In [28]: # dropping unwanted column
city_df_2.drop(columns=['headline_text'],axis=1, inplace=True)
```

```
In [29]: city_df_2.head()
```

Out[29]:

	city_name	Year	No_headline_text
0	bengaluru	2001	1
1	patna	2001	1
2	bengaluru	2001	1
3	delhi	2001	1
4	patna	2001	1

```
In [30]: top_15_headline_city.head(7)
```

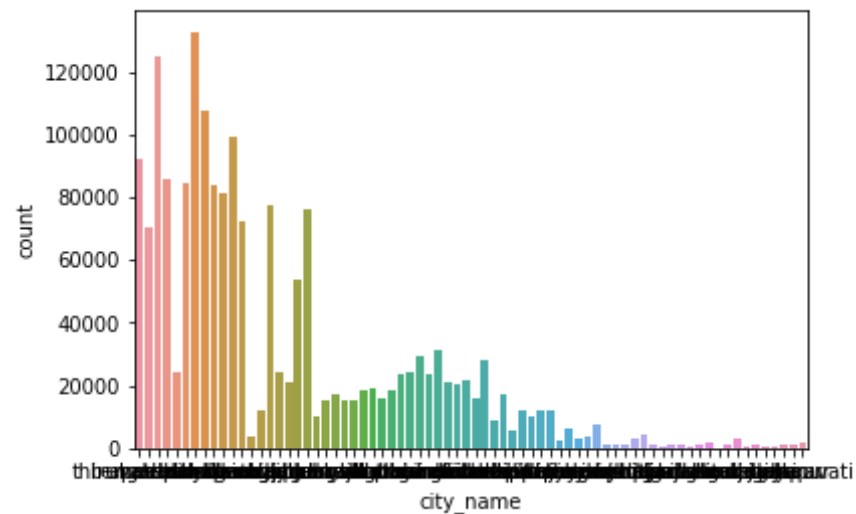
Out[30]:

	No._headline_text
city_name	

	No._headline_text
city_name	
mumbai	132649
delhi	124658
chandigarh	107464
hyderabad	99014
bengaluru	91857
ahmedabad	85813
pune	84620

```
In [31]: sns.countplot(city_df_2['city_name'])
```

```
Out[31]: <matplotlib.axes._subplots.AxesSubplot at 0x54212208>
```



```
In [32]: # creating new data_set for top 7 cities with year
```

```
city_del = city_df_2[city_df_2.city_name == 'mumbai']
city_mum = city_df_2[city_df_2.city_name == 'delhi']
city_chan = city_df_2[city_df_2.city_name == 'chandigarh']
city_hyd = city_df_2[city_df_2.city_name == 'hyderabad']
city_bang = city_df_2[city_df_2.city_name == 'bengaluru']
city_ahme = city_df_2[city_df_2.city_name == 'ahmedabad']
city_pune = city_df_2[city_df_2.city_name == 'pune']
```

```
In [33]: famous_cities = [city_del, city_mum, city_chan, city_hyd, city_bang, ci
ty_ahme, city_pune]
data_famous_cities = pd.concat(famous_cities)
```

```
In [34]: # reset index numbers
data_famous_cities.reset_index(inplace=True)
```

```
In [35]: # top 7 cities
print(data_famous_cities.shape)
data_famous_cities.head()
```

```
(726075, 4)
```

Out[35]:

	index	city_name	Year	No_headline_text
0	37	mumbai	2001	1
1	41	mumbai	2001	1
2	51	mumbai	2001	1
3	55	mumbai	2001	1
4	57	mumbai	2001	1

```
In [36]: # year-wise no. of headlines text
city_year_count = pd.Series(city_df_2.groupby(['Year'])['No_headline_t
xt'].count())

# year-wise data which shows which city have how many no. of headline t
```

```
ext
city_count = data_famous_cities.groupby(['Year', 'city_name']).sum()['No_headline_text'].unstack()
```

*year-wise data which shows which city have how many no. of headline text

```
In [37]: city_count.head()
```

```
Out[37]:
```

city_name	ahmedabad	bengaluru	chandigarh	delhi	hyderabad	mumbai	pune
Year							
2001	3131	5739	2283	2462	4493	2465	1534
2002	5251	7170	5047	4845	6547	4366	3357
2003	4446	4382	5476	4850	5475	3762	3790
2004	2657	2759	3714	5926	2523	2822	2904
2005	1736	1714	1106	2880	1539	1556	974

```
In [38]: city_year_count.head()
```

```
Out[38]: Year
2001      31365
2002      53721
2003      48409
2004      34030
2005      15788
Name: No_headline_text, dtype: int64
```

```
In [39]: fig=plt.figure()

a = fig.add_subplot(111,label="1")
b = fig.add_subplot(111,label="2", frame_on = False)
```

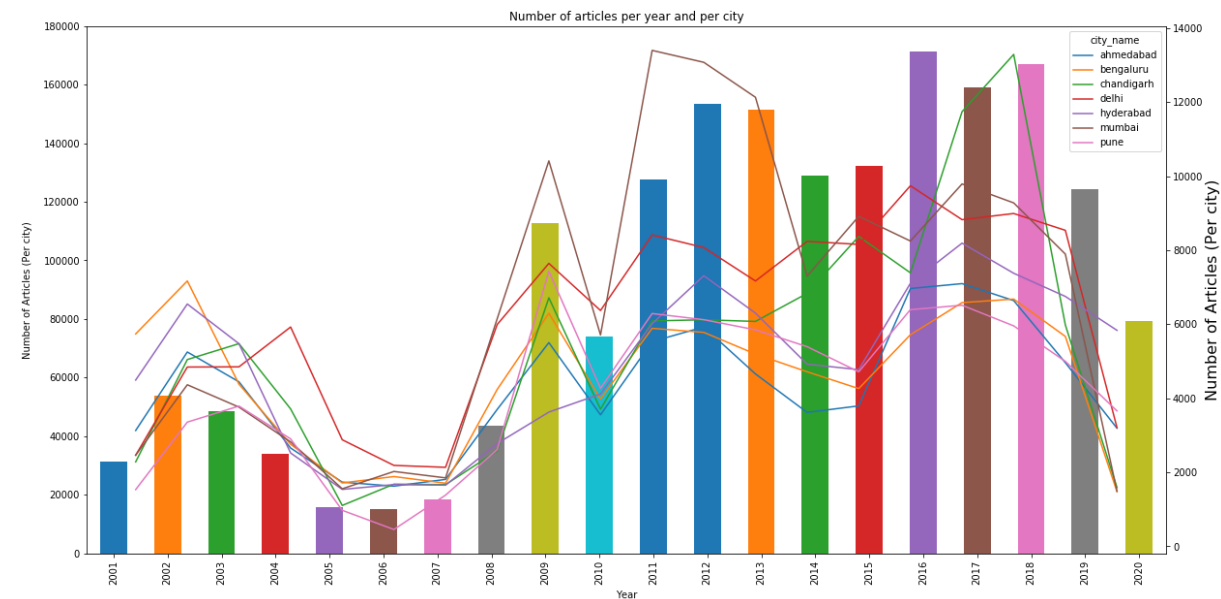


```
plt.figure(figsize=(20,10))
city_year_count.plot(kind='bar',figsize=(20,10), ax=a,title="Number of
articles per year and per city")
a.set_ylabel('Number of Articles (Per city)')
a.set_xlabel("")

city_count.plot(ax=b)
b.set_xticks([])
b.set_ylabel('Number of Articles (Per city)', size =16)
b.yaxis.tick_right()
b.set_xlabel('Year')

a.get_yaxis().set_label_coords(-.05,0.5)
b.get_yaxis().set_label_coords(1.05,0.5)
b.get_xaxis().set_label_coords(0.5, -0.07)

plt.show()
```



<Figure size 1440x720 with 0 Axes>

The graph for the cities was not uniform and hence we overlap on it a graph for the number of city-based articles by year. This way, we can disregard the apparent non-uniformity in city-based reporting over the years.

```
In [40]: df.head()
```

```
Out[40]:
```

	publish_date	headline_category	headline_text
0	20010101	sports.wwe	win over cena satisfying but defeating underta...
1	20010102	unknown	Status quo will not be disturbed at Ayodhya; s...
2	20010102	unknown	Fissures in Hurriyat over Pak visit
3	20010102	unknown	America's unwanted heading for India?
4	20010102	unknown	For bigwigs; it is destination Goa

```
In [41]: # for safety always take copy of original data
new_df = df.copy()
```

```
In [42]: new_df['category'] = new_df['headline_category'].str.split('.').map(lambda x : x[0])
```

```
In [43]: new_df.head()
```

```
Out[43]:
```

	publish_date	headline_category	headline_text	category
0	20010101	sports.wwe	win over cena satisfying but defeating underta...	sports
1	20010102	unknown	Status quo will not be disturbed at Ayodhya; s...	unknown
2	20010102	unknown	Fissures in Hurriyat over Pak visit	unknown

	publish_date	headline_category	headline_text	category
3	20010102	unknown	America's unwanted heading for India?	unknown
4	20010102	unknown	For bigwigs; it is destination Goa	unknown

```
In [44]: top_categories = new_df.groupby(['category']).agg({'headline_text': 'count'}).sort_values(by='headline_text', ascending = False)
```

```
In [45]: top_10_cat = top_categories.drop(['unknown', 'tv', 'top-stories', 'city', 'citizen', 'edit-page', 'top-headlines', 'world']).head(10)
top_10_cat
```

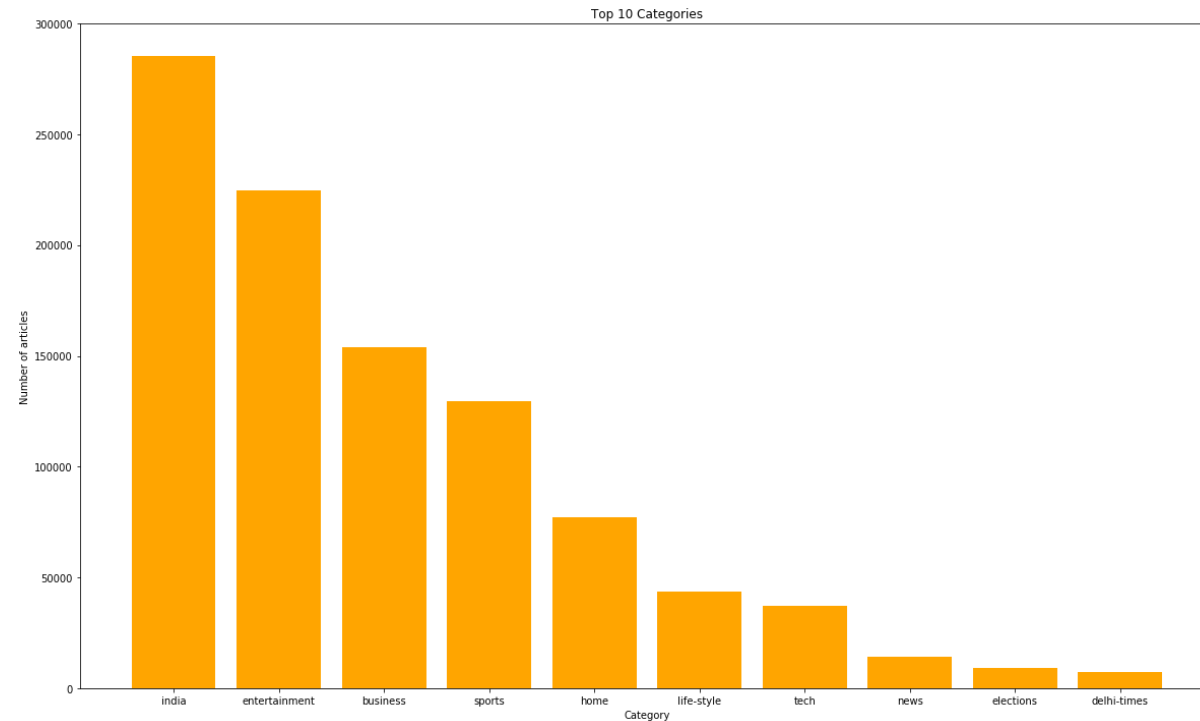
Out[45]:

	headline_text
category	
india	285619
entertainment	224877
business	153876
sports	129494
home	77208
life-style	43370
tech	37220
news	13987
elections	9003
delhi-times	7356

```
In [46]: # reset index
top_10_cat.reset_index(inplace=True)
```

```
In [47]: # visualizing top 10 categories

plt.figure(figsize=(20,12))
plt.bar(top_10_cat.category, height=top_10_cat.headline_text,color = 'orange')
plt.xlabel('Category')
plt.ylabel('Number of articles')
plt.title('Top 10 Categories')
plt.show()
```



India is the most published since the publication is based in India. Interestingly, Entertainment is the second-most covered topic by TOI. In fact, it's eye-opening how Bollywood is given more coverage than every other topic except for Indian Business.

Sentiment Analysis

```
In [12]: def getAnalysis_polarity(headline):  
         result = TextBlob(headline)  
  
         if result.sentiment.polarity < 0:  
             return 'negative'  
         elif result.sentiment.polarity == 0:  
             return 'neutral'  
         else:  
             return 'positive'
```

```
In [13]: from textblob import TextBlob
```

```
In [ ]: df.head()
```

```
In [ ]: df['Result'] = np.array([getAnalysis_polarity(headline) for headline in  
                                df['headline_text']])
```

```
In [ ]: final_result = df.groupby(['Result']).agg({'headline_text': 'count'})
```

```
In [ ]: final_result
```

```
In [ ]: negative = (final_result.loc['negative'] / len(data))*100  
         neutral = (final_result.loc['neutral'] / len(data))*100  
         positive = (final_result.loc['positive'] / len(data))*100
```

```
In [ ]: print('Positive Headlines: ', positive )  
         print('\n\nNegative Headlines: ', negative )  
         print('\n\nUnbiased Headlines: ', str(neutral))
```

Stock Price Prediction

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from pandas.plotting import lag_plot
from pandas import datetime
from statsmodels.tsa.arima_model import ARIMA
from sklearn.metrics import mean_squared_error, mean_absolute_error
from statsmodels.tsa.seasonal import seasonal_decompose
```

```
In [2]: df = pd.read_csv('D://DK//dddddd//Task_grip//TASK_3//task 4//History_S
tock.csv')
```

```
In [3]: df.head()
```

Out[3]:

	Date	Open	High	Low	Close	Adj Close	Volume
0	2019-10-21	192.960007	196.550003	192.199997	196.009995	195.600586	6546700
1	2019-10-22	198.300003	202.529999	195.300003	195.610001	195.201431	8218000
2	2019-10-23	192.309998	195.660004	191.029999	195.089996	194.682495	6997900
3	2019-10-24	196.750000	198.210007	195.380005	196.860001	196.448822	5800500
4	2019-10-25	200.100006	205.380005	199.789993	204.539993	204.112762	10577300

```
In [4]: df.shape
```

Out[4]: (252, 7)

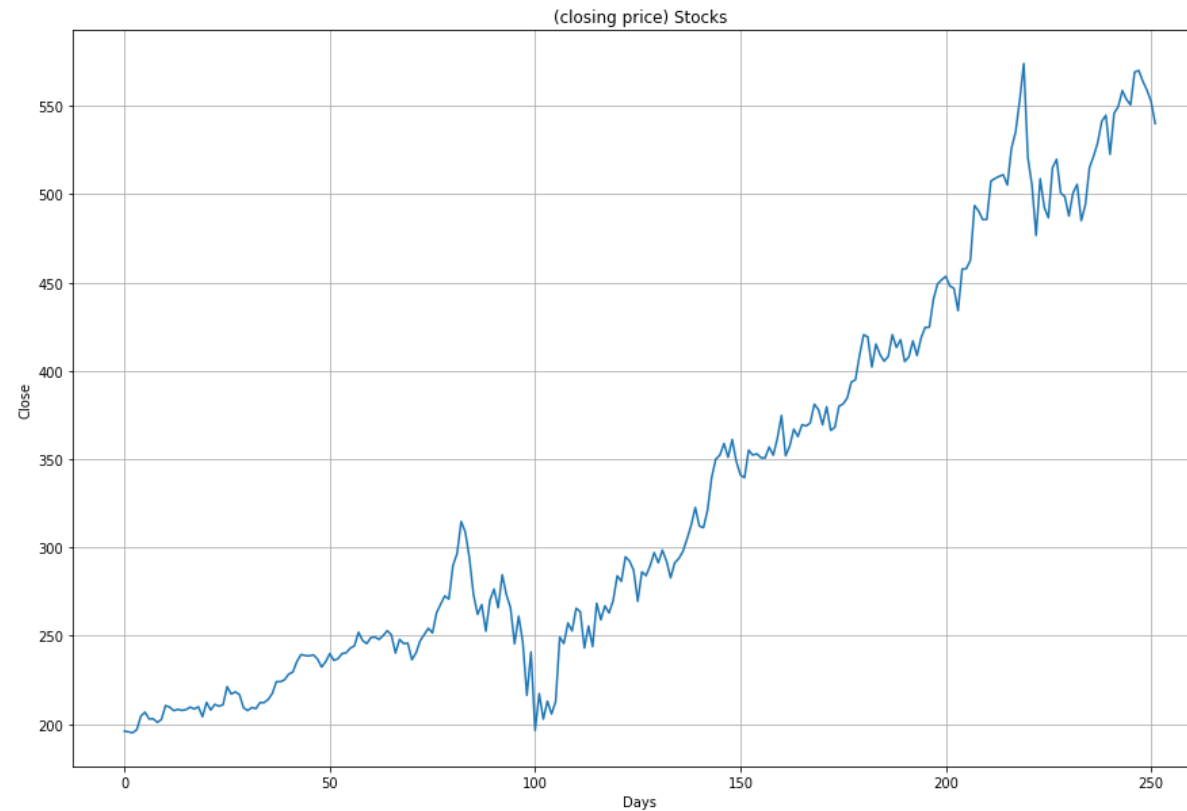
```
In [5]: df.describe()
```

Out[5]:

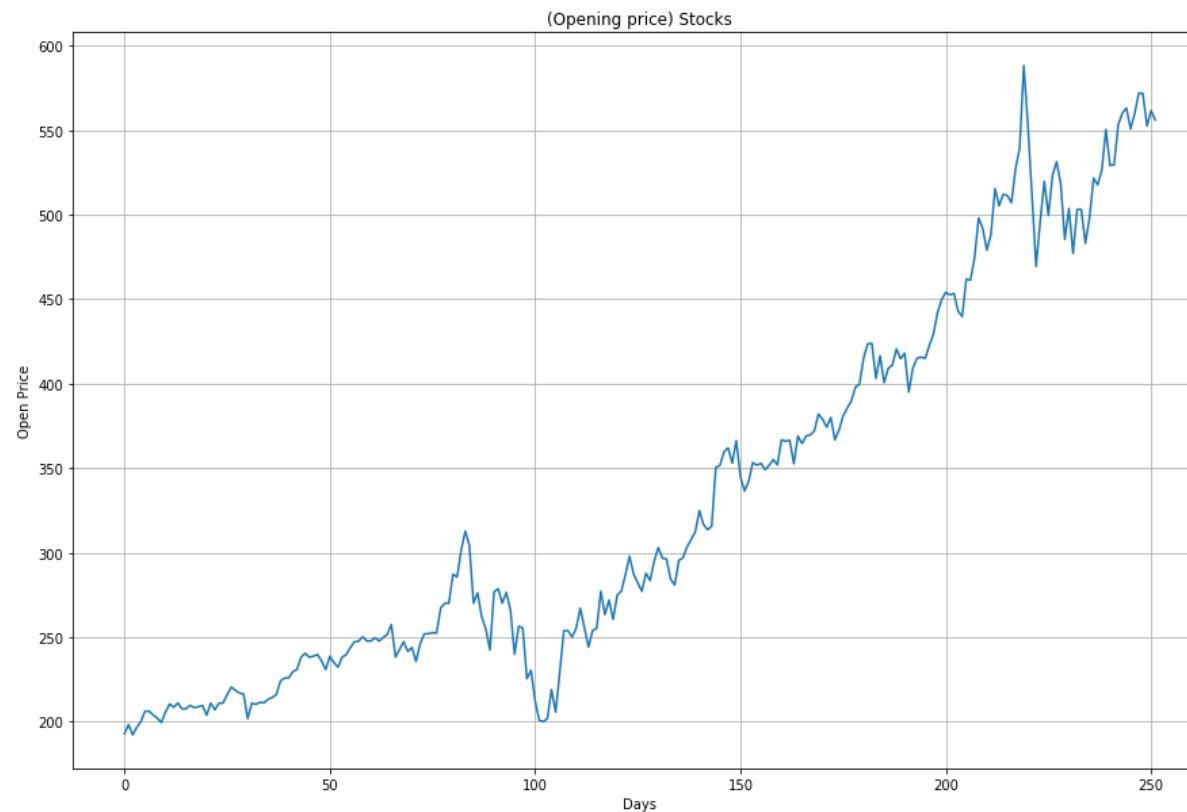
	Open	High	Low	Close	Adj Close	Volume
count	252.000000	252.000000	252.000000	252.000000	252.000000	2.520000e+02

	Open	High	Low	Close	Adj Close	Volume
mean	331.695040	338.014763	325.542262	332.267104	332.052070	1.216431e+07
std	112.555459	114.099520	110.064080	112.066886	112.176138	5.775373e+06
min	192.309998	195.660004	180.679993	195.089996	194.682495	3.471600e+06
25%	239.509995	241.982494	237.225002	240.297501	239.972519	7.770400e+06
50%	287.550003	297.884995	286.274994	291.819992	291.599625	1.059065e+07
75%	415.052498	421.407501	408.454994	415.524987	415.400704	1.558578e+07
max	588.150024	589.070007	560.750000	573.859985	573.859985	3.659210e+07

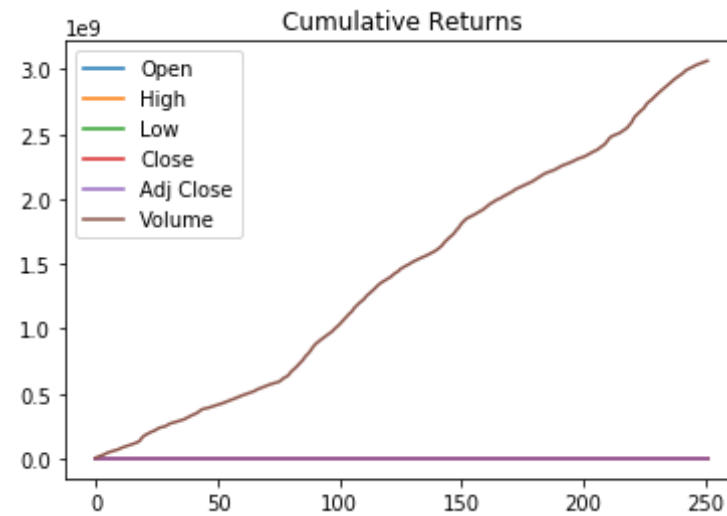
```
In [6]: plt.figure(figsize=(15,10))
plt.grid(True)
plt.plot(df['Close'])
plt.xlabel('Days')
plt.ylabel('Close')
plt.title('(closing price) Stocks')
plt.show()
```



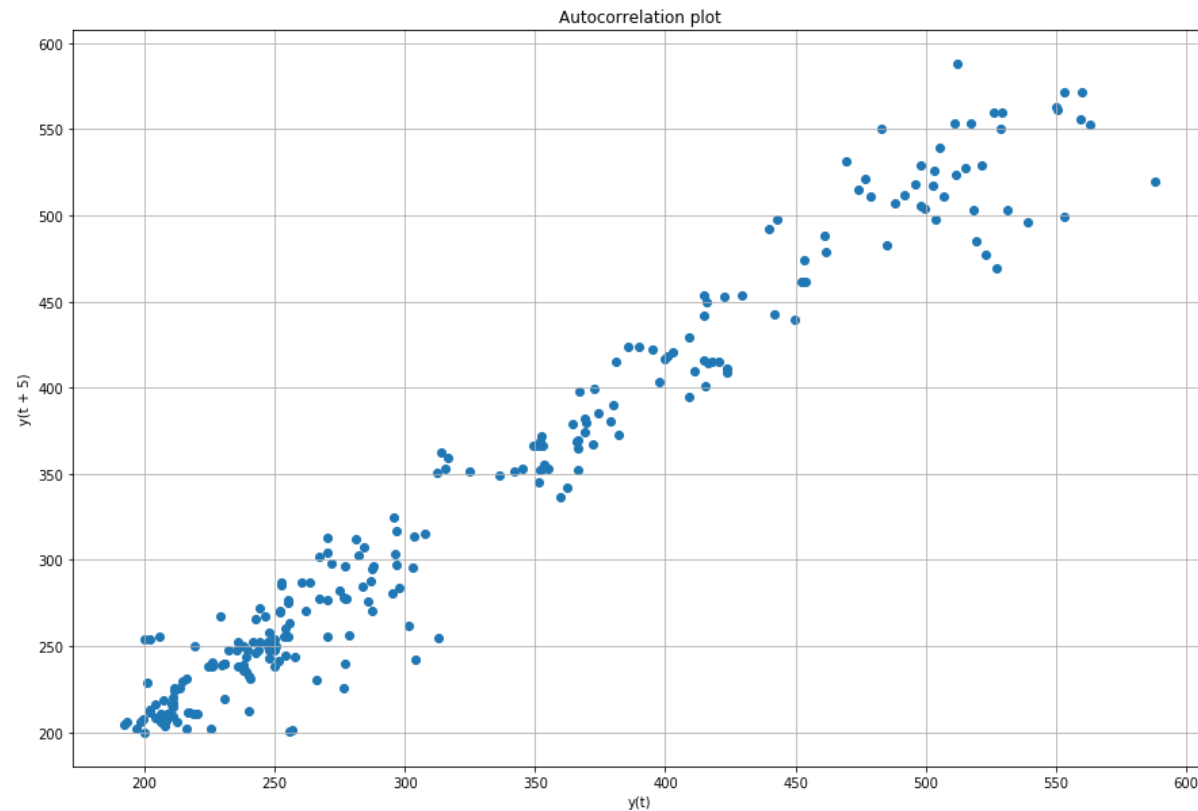
```
In [7]: plt.figure(figsize=(15,10))
plt.grid(True)
plt.plot(df['Open'])
plt.xlabel('Days')
plt.ylabel('Open Price')
plt.title('(Opening price) Stocks')
plt.show()
```

```
In [8]: # Cumulative Return
cumsum_data = df.cumsum()
cumsum_data.plot()
plt.title('Cumulative Returns')
plt.show()
```

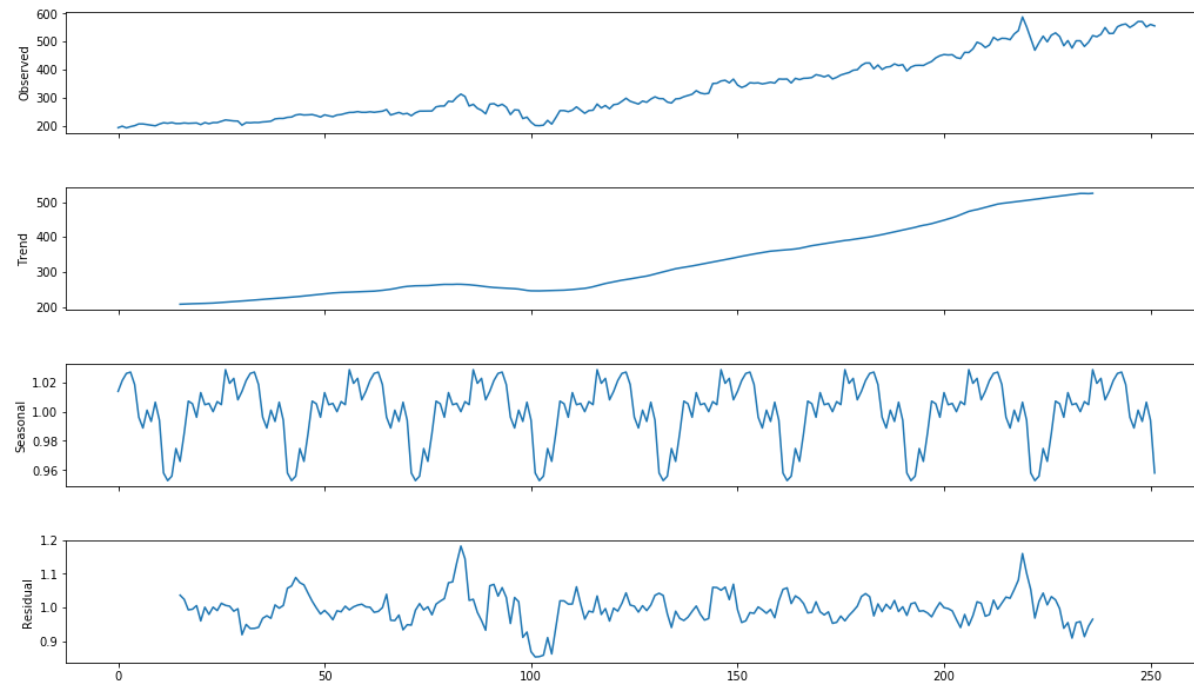


```
In [9]: plt.figure(figsize=(15,10))
plt.grid(True)
lag_plot(df['Open'], lag=5)
plt.title('Autocorrelation plot')
plt.show()
```



```
In [10]: sea_dec = seasonal_decompose(df['Open'], model='multiplicative', freq =
        30)
fig = plt.figure()
fig = sea_dec.plot()
fig.set_size_inches(16, 9)
plt.show()
```

<Figure size 432x288 with 0 Axes>



```
In [11]: # Extract date frame and plot closing stock price w.r.t time
df['Date'] = pd.to_datetime(df.Date, format='%Y-%m-%d')
df.index = df['Date']
df.dropna(inplace=True)
#plot
plt.figure(figsize=(16,8))
plt.plot(df['Close'], label='Close Price history')
```

```
Out[11]: [<matplotlib.lines.Line2D at 0x1844db00>]
```

