

Machine learning

INTERMEDIATE QUESTIONS

Q-1. Imagine you have a dataset where you have different Instagram features like u sername , Caption , Hashtag , Followers , Time_Since_posted , and likes , now your task is to predict the number of likes and Time Since posted and the rest of the features are your input features. Now you have to build a model which can predict the number of likes and Time Since posted.

Dataset This is the Dataset You can use this dataset for this question.

In [3]:

```
1 !pip install wordcloud
```

Collecting wordcloud

Downloading wordcloud-1.9.2-cp39-cp39-win_amd64.whl (153 kB)

Requirement already satisfied: matplotlib in c:\users\abhik\anaconda3\lib\site-packages (from wordcloud) (3.4.3)

Requirement already satisfied: numpy>=1.6.1 in c:\users\abhik\anaconda3\lib\site-packages (from wordcloud) (1.20.3)

Requirement already satisfied: pillow in c:\users\abhik\anaconda3\lib\site-packages (from wordcloud) (8.4.0)

Requirement already satisfied: python-dateutil>=2.7 in c:\users\abhik\anaconda3\lib\site-packages (from matplotlib->wordcloud) (2.8.2)

Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\abhik\anaconda3\lib\site-packages (from matplotlib->wordcloud) (1.3.1)

Requirement already satisfied: pyparsing>=2.2.1 in c:\users\abhik\anaconda3\lib\site-packages (from matplotlib->wordcloud) (3.0.4)

Requirement already satisfied: cycler>=0.10 in c:\users\abhik\anaconda3\lib\site-packages (from matplotlib->wordcloud) (0.10.0)

Requirement already satisfied: six in c:\users\abhik\anaconda3\lib\site-packages (from cycler>=0.10->matplotlib->wordcloud) (1.16.0)

Installing collected packages: wordcloud

Successfully installed wordcloud-1.9.2

In [4]:

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import mean_squared_error
6 from sklearn.preprocessing import OneHotEncoder
7 import re
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 import plotly.express as px
11 from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator
12
```

In [11]:

```
1 ## Load the dataset using pandas:
2 data = pd.read_csv("Downloads/archive/instagram_reach.csv")
```

In [12]:

```
1 ## Checking top 5 rows
2 data.head()
```

Out[12]:

	Unnamed: 0	S.No	USERNAME	Caption	Followers	
0	0	1	mikequindazzi	Who are #DataScientist and what do they do? >>...	1600	#MachineLea
1	1	2	drgorillapaints	We all know where it's going. We just have to ...	880	#deck .#mac #macintosh#
2	2	3	aitrading_official	Alexander Barinov: 4 years as CFO in multinati...	255	#whoiswho #aitrading #a
3	3	4	opensourcedworkplace	sfad	340	#iot #cre#workplace #CDO
4	4	5	crea.vision	Ever missed a call while your phone was chargi...	304	#instamachinelearning #in

In [13]:

```
1 ## Checking Rows & Columns Availabale in Dataset
2 data.shape
```

Out[13]:

(100, 8)

In [14]:

```
1 ## Checking Details Information related with Dataset
2 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 8 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Unnamed: 0            100 non-null   int64
1   S.No                  100 non-null   int64
2   USERNAME              100 non-null   object
3   Caption               94 non-null    object
4   Followers             100 non-null   int64
5   Hashtags              100 non-null   object
6   Time since posted     100 non-null   object
7   Likes                 100 non-null   int64
dtypes: int64(4), object(4)
memory usage: 6.4+ KB
```

In [15]:

```
1 ## Checking ALL Columns name present in dataset
2 data.columns
```

Out[15]:

```
Index(['Unnamed: 0', 'S.No', 'USERNAME', 'Caption', 'Followers', 'Hashtags',
      'Time since posted', 'Likes'],
      dtype='object')
```

In [16]:

```
1 ## checking top 2 rows of dataset
2 data.head(2)
```

Out[16]:

	Unnamed: 0	S.No	USERNAME	Caption	Followers	
0	0	1	mikequindazzi	Who are #DataScientist and what do they do? >>...	1600	#MachineLearning #AI ;
1	1	2	drgorillapaints	We all know where it's going. We just have to ...	880	#deck .#mac #macintosh#sayhello #

In [17]:

```
1 # Remove unnecessary columns
2 data= data.drop(['Unnamed: 0', 'S.No'], axis=1)
```

In [18]:

```
1 ## Checking ALL Columns name present in Dataset
2 data.columns
```

Out[18]:

Index(['USERNAME', 'Caption', 'Followers', 'Hashtags', 'Time since posted',
 'Likes'],
 dtype='object')

In [19]:

```
1 ## Checking top 3 rows of dataset after dropping unnecessary columns.
2 data.head(3)
```

Out[19]:

	USERNAME	Caption	Followers	Hashtags
0	mikequindazzi	Who are #DataScientist and what do they do? >>...	1600	#MachineLearning #AI #DataAnalytics #DataScien...
1	drgorillapaints	We all know where it's going. We just have to ...	880	#deck .#mac #macintosh#sayhello #apple #steve...
2	aitrading_official	Alexander Barinov: 4 years as CFO in multinati...	255	#whoiswho #aitrading #ai #aitradingteam#instat...

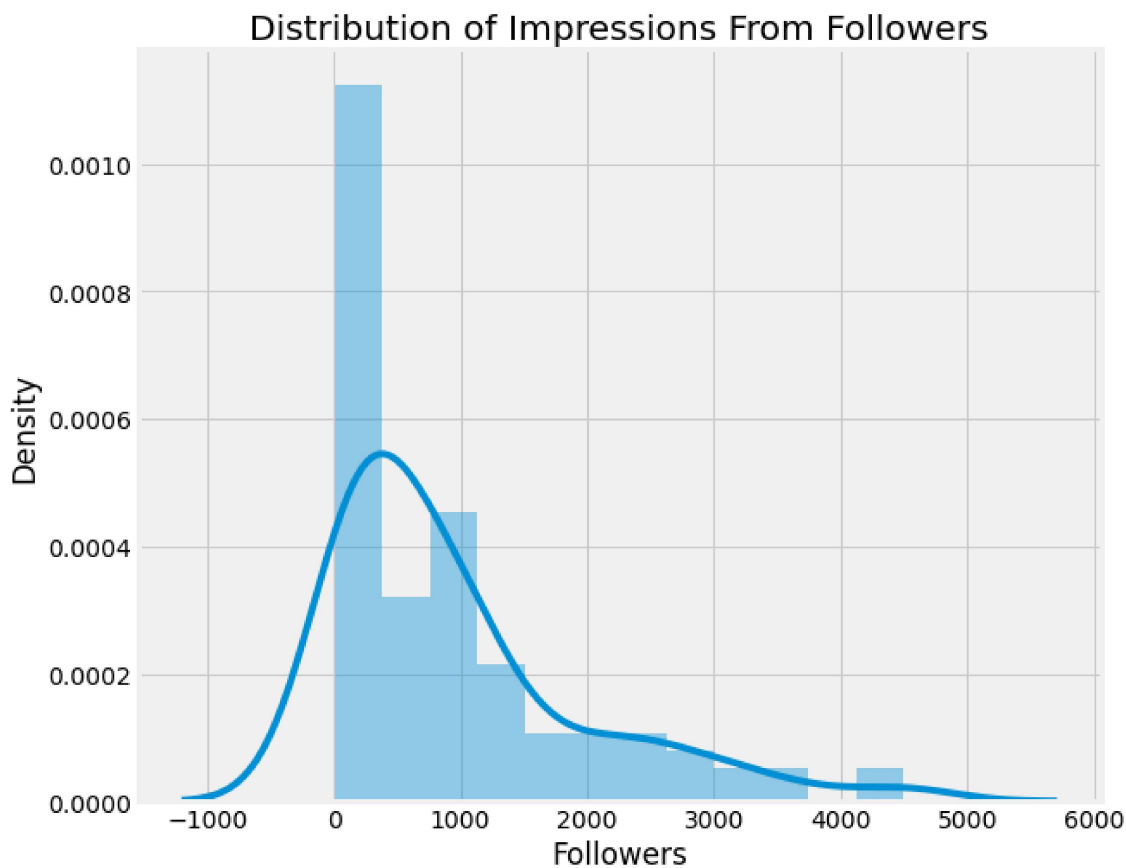
Doing EDA and Analyzing Instagram Reach

In [20]:

```
1  ## Distribution of Impressions From Followers
2  plt.figure(figsize=(10, 8))
3  plt.style.use('fivethirtyeight')
4  plt.title("Distribution of Impressions From Followers")
5  sns.distplot(data['Followers'])
6  plt.show()
7
```

C:\Users\abhik\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)

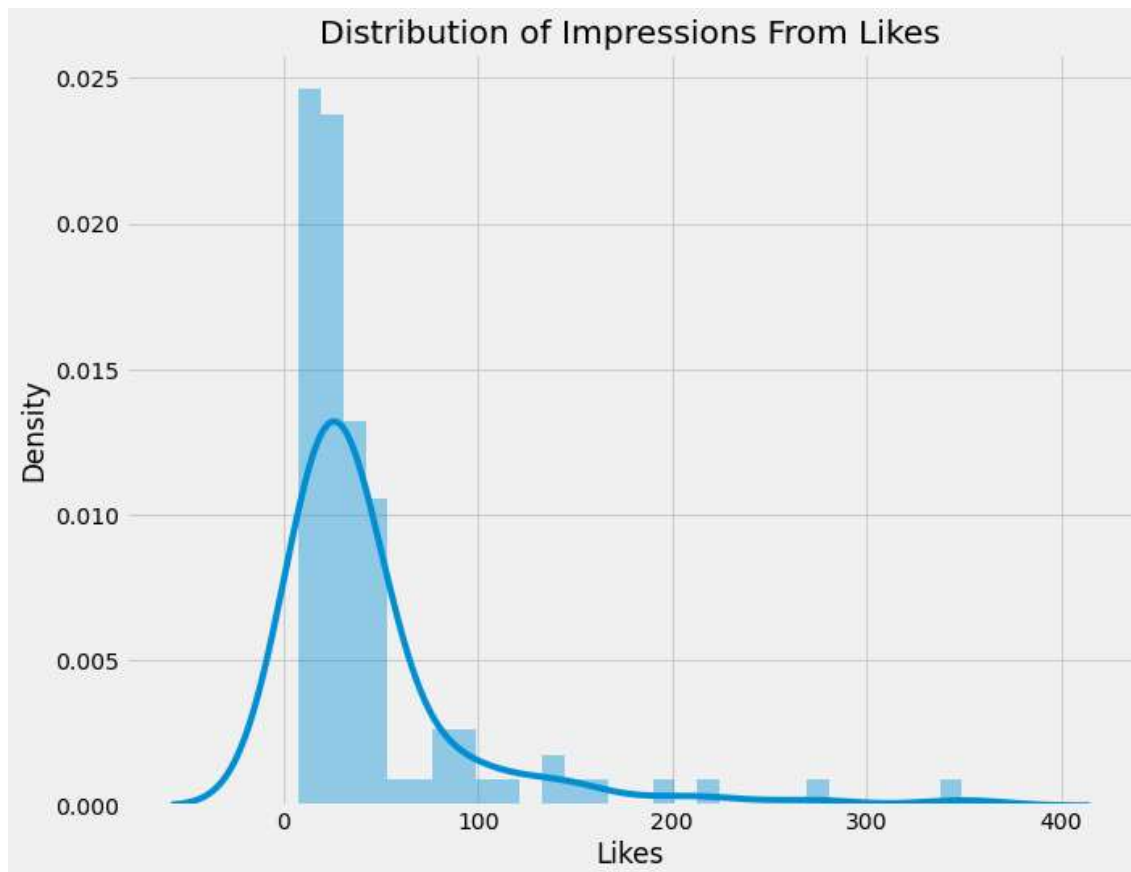


In [21]:

```
1 ## Distribution of Impressions From Likes
2 plt.figure(figsize=(10, 8))
3 plt.title("Distribution of Impressions From Likes")
4 sns.distplot(data['Likes'])
5 plt.show()
```

C:\Users\abhik\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

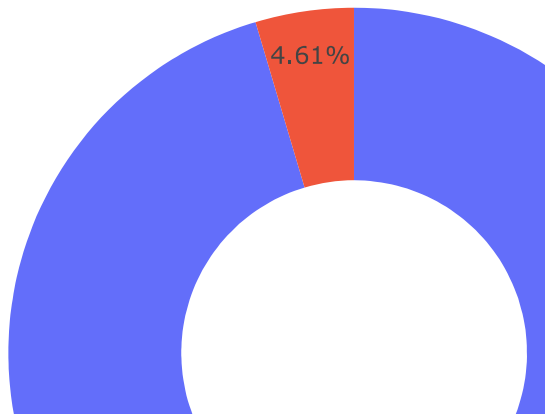
warnings.warn(msg, FutureWarning)



In [25]:

```
1  ## Relation between Likes and Followers
2
3  followers = data["Followers"].sum()
4  likes = data["Likes"].sum()
5
6  labels = ['Followers', 'Likes']
7  values = [followers, likes]
8
9  fig = px.pie(data, values=values, names=labels,
10              title='Impressions on Instagram Posts From Various Sources', hole=0.5)
11  fig.show()
12
13
```

Impressions on Instagram Posts From Various Sources



In [37]:

```
1  from wordcloud import WordCloud
2
3  # Provide the path to the TrueType font file
4  font_path = "/path/to/your/font.ttf"
5
6  # Create a WordCloud object with the specified font
7  wordcloud = WordCloud(font_path=font_path)
```

In []:

```
1  ## Plotting Word-Cloud for Hashtag Related Data
2
3  text = " ".join(i for i in data.Hashtags)
4  stopwords = set(STOPWORDS)
5  wordcloud = WordCloud(stopwords=stopwords, background_color="white").generate(text)
6  plt.style.use('classic')
7  plt.figure( figsize=(12,10))
8  plt.imshow(wordcloud, interpolation='bilinear')
9  plt.axis("off")
10 plt.show()
```

In [27]:

```
1  ## Plotting Scatter-plot for showing Relationship Between Likes and Followers
2
3  figure = px.scatter(data_frame = data, x="Likes",
4                      y="Followers", trendline="ols",
5                      title = "Relationship Between Likes and Followers")
6  figure.show()
```

Relationship Between Likes and Followers



In [28]:

```
1 # Select the relevant features and target variables
2
3 features = ['USERNAME', 'Caption', 'Hashtags', 'Followers']
4 target_likes = 'Likes'
5 target_time_since_posted = 'Time since posted'
6
```

In [29]:

```
1 # Split the data into training and testing sets
2
3 X = data[features]
4 y_likes = data[target_likes]
5 y_time_since_posted = data[target_time_since_posted]
6 X_train, X_test, y_likes_train, y_likes_test, y_time_since_posted_train, y_time_sin
7
```

In [30]:

```
1 # Preprocess the text features using one-hot encoding
2 encoder = OneHotEncoder(sparse=False, handle_unknown='ignore')
3 X_train_encoded = encoder.fit_transform(X_train)
4 X_test_encoded = encoder.transform(X_test)
5
```

Train a model to predict the number of likes:

In [31]:

```
1 # Train a model to predict the number of Likes
2 likes_model = LinearRegression()
3 likes_model.fit(X_train_encoded, y_likes_train)
4 likes_predictions = likes_model.predict(X_test_encoded)
5 likes_mse = mean_squared_error(y_likes_test, likes_predictions)
6 print("Mean Squared Error (Likes):", likes_mse)
```

Mean Squared Error (Likes): 1639.6633407976692

Train a model to predict the time since posted

In [32]:

```
1 # Preprocess the time since posted variable
2 def extract_numerical_value(time_string):
3     numerical_value = re.findall(r'\d+', time_string)[0]
4     return int(numerical_value)
```

In [33]:

```
1 y_time_since_posted_train = y_time_since_posted_train.apply(extract_numerical_value)
2 y_time_since_posted_test = y_time_since_posted_test.apply(extract_numerical_value)
3
```

In [34]:

```
1 # Train a model to predict the time since posted
2 time_since_posted_model = LinearRegression()
3 time_since_posted_model.fit(X_train_encoded, y_time_since_posted_train)
4 time_since_posted_predictions = time_since_posted_model.predict(X_test_encoded)
5 time_since_posted_mse = mean_squared_error(y_time_since_posted_test, time_since_pos
6 print("Mean Squared Error (Time Since Posted):", time_since_posted_mse)
```

Mean Squared Error (Time Since Posted): 12.861392010243321

In []:

```
1
```