

Problem statement

Data collection

Importing libraries

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

Importing dataset

```
In [2]: data=pd.read_csv(r"C:\Users\user\Downloads\5_Instagram data - 5_Instagram data.csv")
data
```

```
Out[2]:
```

	Impressions	From Home	From Hashtags	From Explore	From Other	Saves	Comments	Shares	Likes	Profile Visits	Follows
0	3920	2586	1028	619	56	98	9	5	162	35	2
1	5394	2727	1838	1174	78	194	7	14	224	48	10
2	4021	2085	1188	0	533	41	11	1	131	62	12
3	4528	2700	621	932	73	172	10	7	213	23	8
4	2518	1704	255	279	37	96	5	4	123	8	0

	Impressions	From Home	From Hashtags	From Explore	From Other	Saves	Comments	Shares	Likes	Profile Visits	Follows
...
114	13700	5185	3041	5352	77	573	2	38	373	73	80
115	5731	1923	1368	2266	65	135	4	1	148	20	18
116	4139	1133	1538	1367	33	36	0	1	92	34	10
117	32695	11815	3147	17414	170	1095	2	75	549	148	214
118	36919	13473	4176	16444	2547	653	5	26	443	611	228

119 rows × 13 columns

head

In [3]:

```
# to display first 8 dataset values
da=data.head(8)
da
```

Out[3]:

	Impressions	From Home	From Hashtags	From Explore	From Other	Saves	Comments	Shares	Likes	Profile Visits	Follows
0	3920	2586	1028	619	56	98	9	5	162	35	2
1	5394	2727	1838	1174	78	194	7	14	224	48	10

	Impressions	From Home	From Hashtags	From Explore	From Other	Saves	Comments	Shares	Likes	Profile Visits	Follows	
												L
2	4021	2085	1188	0	533	41	11	1	131	62	12	m
3	4528	2700	621	932	73	172	10	7	213	23	8	c pr
4	2518	1704	255	279	37	96	5	4	123	8	0	an v
5	3884	2046	1214	329	43	74	7	10	144	9	2	i
6	2621	1543	599	333	25	22	5	1	76	26	0	L t c ch
7	3541	2071	628	500	60	135	4	9	124	12	6	b y

info

```
In [4]: # to identify missing values
data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 119 entries, 0 to 118
Data columns (total 13 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Impressions     119 non-null    int64
1   From Home       119 non-null    int64
2   From Hashtags   119 non-null    int64
3   From Explore    119 non-null    int64
4   From Other      119 non-null    int64
```

```
5   Saves          119 non-null    int64
6   Comments       119 non-null    int64
7   Shares         119 non-null    int64
8   Likes          119 non-null    int64
9   Profile Visits 119 non-null    int64
10  Follows        119 non-null    int64
11  Caption        119 non-null    object
12  Hashtags       119 non-null    object
dtypes: int64(11), object(2)
memory usage: 12.2+ KB
```

describe

```
In [5]: # to display summary of the dataset
data.describe()
```

Out[5]:

	Impressions	From Home	From Hashtags	From Explore	From Other	Saves	Comments
count	119.000000	119.000000	119.000000	119.000000	119.000000	119.000000	119.000000
mean	5703.991597	2475.789916	1887.512605	1078.100840	171.092437	153.310924	6.663866
std	4843.780105	1489.386348	1884.361443	2613.026132	289.431031	156.317731	3.544576
min	1941.000000	1133.000000	116.000000	0.000000	9.000000	22.000000	0.000000
25%	3467.000000	1945.000000	726.000000	157.500000	38.000000	65.000000	4.000000
50%	4289.000000	2207.000000	1278.000000	326.000000	74.000000	109.000000	6.000000
75%	6138.000000	2602.500000	2363.500000	689.500000	196.000000	169.000000	8.000000
max	36919.000000	13473.000000	11817.000000	17414.000000	2547.000000	1095.000000	19.000000

columns

```
In [6]: # to display headings of the dataset
data.columns
```

Out[6]: Index(['Impressions', 'From Home', 'From Hashtags', 'From Explore', 'From Other', 'Saves', 'Comments', 'Shares', 'Likes', 'Profile Visits', 'Follows', 'Caption', 'Hashtags'], dtype='object')

```
In [7]: a=data.dropna(axis=1)
a
```

Out[7]:

	Impressions	From Home	From Hashtags	From Explore	From Other	Saves	Comments	Shares	Likes	Profile Visits	Follows
0	3920	2586	1028	619	56	98	9	5	162	35	2

	Impressions	From Home	From Hashtags	From Explore	From Other	Saves	Comments	Shares	Likes	Profile Visits	Follows
1	5394	2727	1838	1174	78	194	7	14	224	48	10
2	4021	2085	1188	0	533	41	11	1	131	62	12
3	4528	2700	621	932	73	172	10	7	213	23	8
4	2518	1704	255	279	37	96	5	4	123	8	0
...
114	13700	5185	3041	5352	77	573	2	38	373	73	80
115	5731	1923	1368	2266	65	135	4	1	148	20	18
116	4139	1133	1538	1367	33	36	0	1	92	34	10
117	32695	11815	3147	17414	170	1095	2	75	549	148	214
118	36919	13473	4176	16444	2547	653	5	26	443	611	228

Impressions	From Home	From Hashtags	From Explore	From Other	Saves	Comments	Shares	Likes	Profile Visits	Follows
-------------	-----------	---------------	--------------	------------	-------	----------	--------	-------	----------------	---------

119 rows × 13 columns

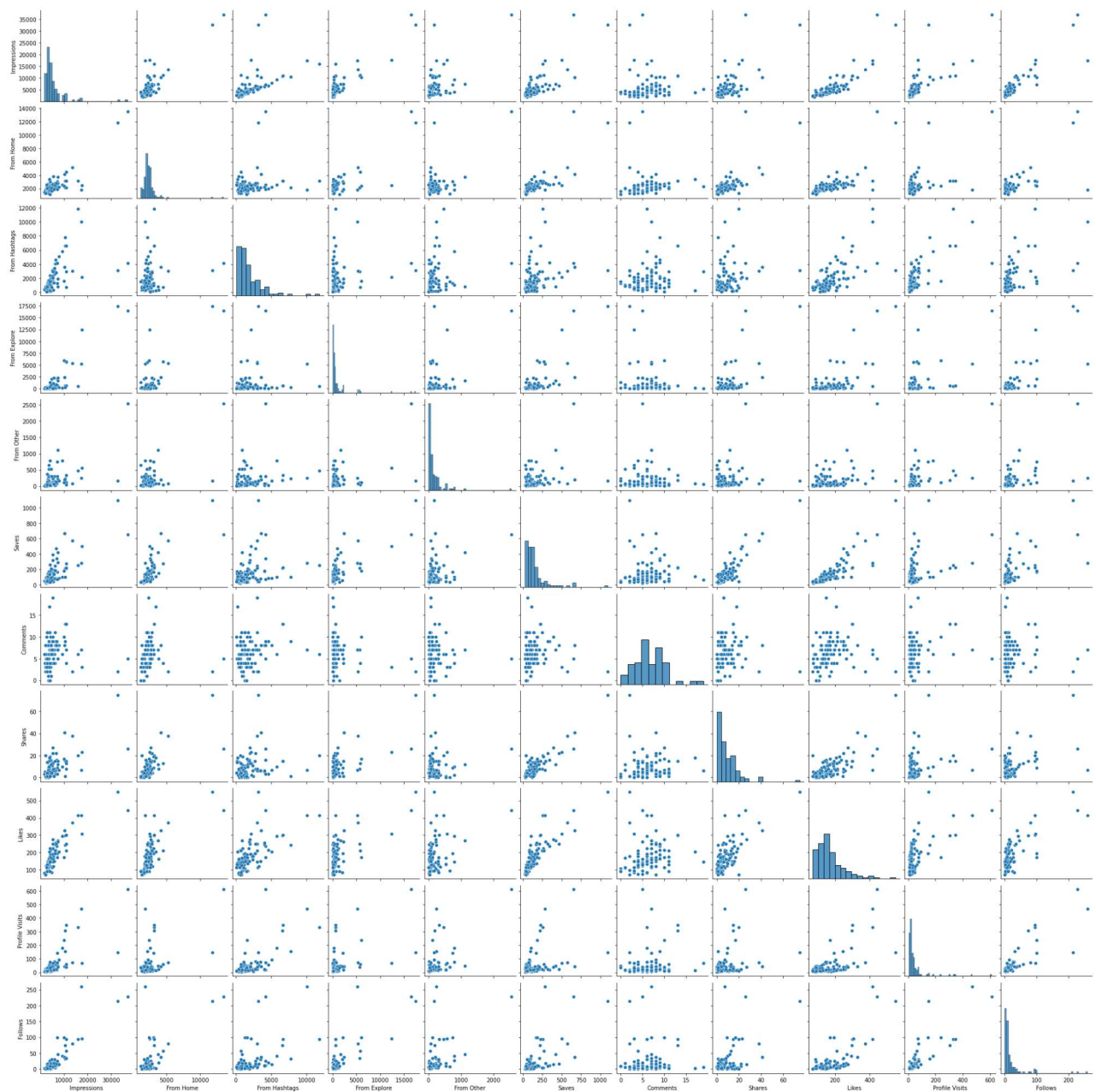
```
In [8]: a.columns
```

```
Out[8]: Index(['Impressions', 'From Home', 'From Hashtags', 'From Explore',  
             'From Other', 'Saves', 'Comments', 'Shares', 'Likes', 'Profile Visits',  
             'Follows', 'Caption', 'Hashtags'],  
            dtype='object')
```

EDA and Visualization

```
In [9]: sns.pairplot(a)
```

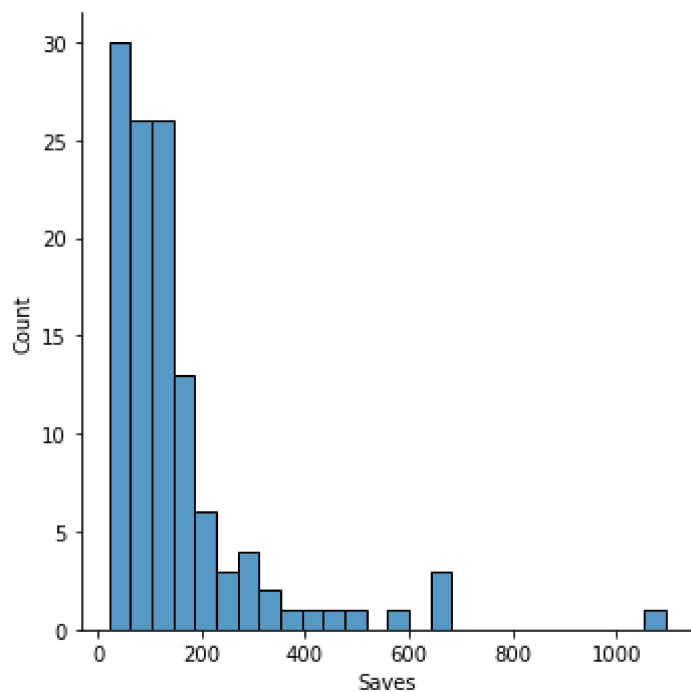
```
Out[9]: <seaborn.axisgrid.PairGrid at 0x2f7a5ce9b80>
```



distribution plot

```
In [10]: sns.displot(a["Saves"])
```

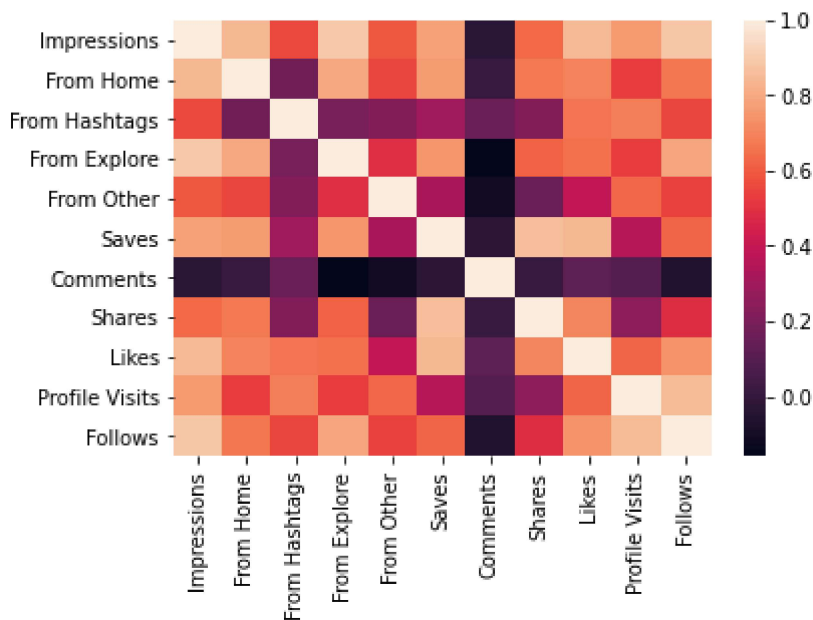
```
Out[10]: <seaborn.axisgrid.FacetGrid at 0x2f7a8f20f70>
```



correlation

```
In [11]: dat=data[['Impressions', 'From Home', 'From Hashtags', 'From Explore',
                  'From Other', 'Saves', 'Comments', 'Shares', 'Likes', 'Profile Visits',
                  'Follows', 'Caption', 'Hashtags']]
          sns.heatmap(dat.corr())
```

Out[11]: <AxesSubplot:>



To train the model-Model Building


```
In [12]: x=a[['Impressions', 'From Home', 'From Hashtags', 'From Explore',
            'From Other', 'Saves', 'Comments', 'Shares', 'Likes', 'Profile Visits',
            ]]
        y=a['Follows']
```

```
In [13]: # to split my dataset into training and test data
        from sklearn.model_selection import train_test_split
        x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

```
In [14]: from sklearn.linear_model import LinearRegression
        lr= LinearRegression()
        lr.fit(x_train,y_train)
```

Out[14]: LinearRegression()

```
In [15]: print(lr.intercept_)
```

-0.4051550278512117

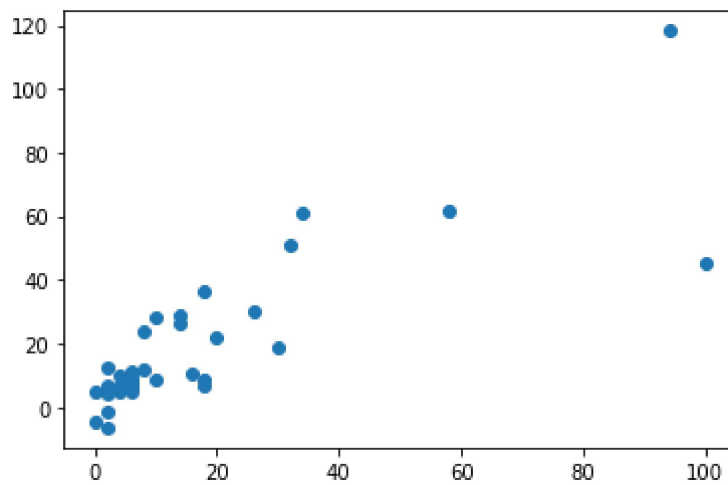
```
In [16]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
        coeff
```

Out[16]:

	Co-efficient
Impressions	-0.002162
From Home	-0.001509
From Hashtags	0.001912
From Explore	0.009281
From Other	-0.019073
Saves	0.020288
Comments	-1.029568
Shares	-0.135297
Likes	0.098460
Profile Visits	0.309381

```
In [17]: prediction=lr.predict(x_test)
        plt.scatter(y_test,prediction)
```

Out[17]: <matplotlib.collections.PathCollection at 0x2f7ac5174f0>



```
In [18]: print(lr.score(x_test,y_test))
```

```
0.6409062978147559
```

```
In [19]: lr.score(x_train,y_train)
```

```
Out[19]: 0.9390755596965219
```

Ridge regression

```
In [20]: from sklearn.linear_model import Ridge,Lasso
```

```
In [21]: rr=Ridge(alpha=10)
rr.fit(x_train,y_train)
rr.score(x_test,y_test)
```

```
Out[21]: 0.6411059062758468
```

```
In [22]: rr.score(x_train,y_train)
```

```
Out[22]: 0.9390747777345017
```

Lasso regression

```
In [25]: la=Lasso(alpha=10)
la.fit(x_train,y_train)
print(la.score(x_train,y_train))
```

```
0.9342791048069933
```

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\_coordinate_descent.py:5
30: ConvergenceWarning: Objective did not converge. You might want to increase the numbe
r of iterations. Duality gap: 909.5892184644026, tolerance: 17.78378795180723
model = cd_fast.enet_coordinate_descent(
```

In [24]: `la.score(x_test,y_test)`

Out[24]: 0.6269762366994706

In []: