

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: df=pd.read_csv(r'C:\Users\sharma17\Downloads\ml\regression dataset\Startups.csv')
```

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

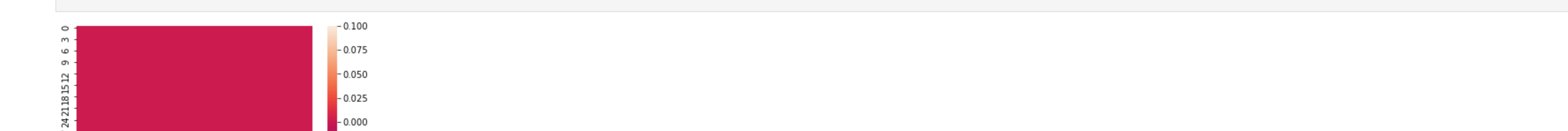
	R&D Spend	Administration	Marketing Spend	State	Profit
0	165349.20	136897.80	471784.10	New York	182261.83
1	162597.70	151377.59	448988.53	California	191792.06
2	153441.51	101145.55	407934.54	Florida	191050.39
3	144372.41	118671.85	383199.62	New York	182901.99
4	142107.34	91391.77	366168.42	Florida	166187.94

```
In [5]: df.drop('State',inplace=True,axis=1)
```

```
In [6]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50 entries, 0 to 49
Data columns (total 4 columns):
 #   Column              Non-Null Count  Dtype  
---  --
 0   R&D Spend           50 non-null     float64
 1   Administration      50 non-null     float64
 2   Marketing Spend     50 non-null     float64
 3   Profit              50 non-null     float64
dtypes: float64(4)
memory usage: 1.7 KB
```

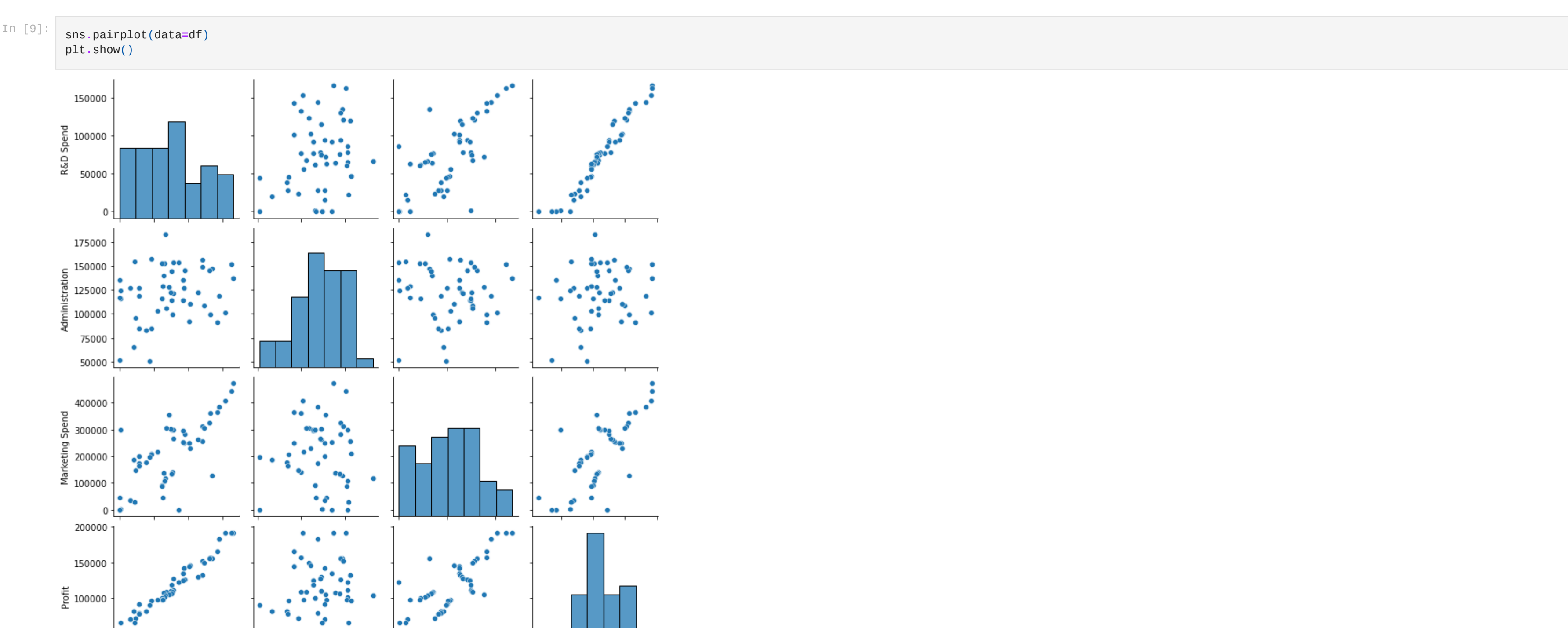
```
In [7]: sns.heatmap(df.isnull())
plt.show()
```



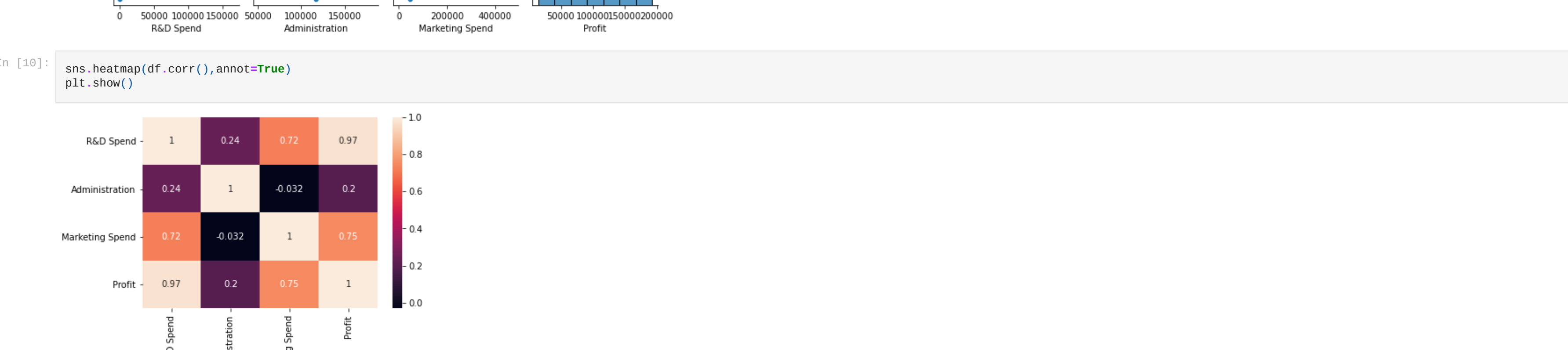
```
In [8]: df.describe()
```

	R&D Spend	Administration	Marketing Spend	Profit
count	50.000000	50.000000	50.000000	50.000000
mean	73721.615600	121344.629600	231025.097800	112012.620200
std	45902.256482	28017.802755	122290.310726	40306.180338
min	0.000000	51283.140000	0.000000	14681.400000
25%	39935.370000	103730.875000	129300.132500	90139.902500
50%	73951.080000	122688.795000	212716.240000	107978.190000
75%	101692.800000	144842.180000	294469.085000	139765.977500
max	165349.200000	182645.560000	471784.100000	182261.830000

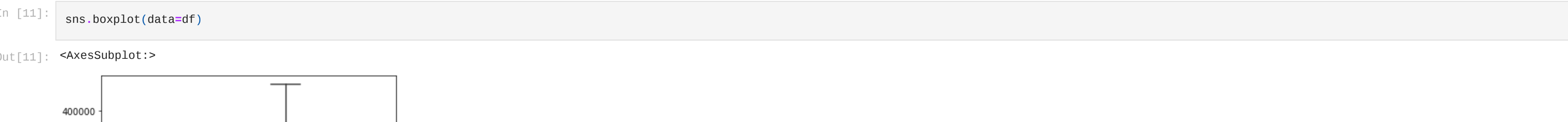
```
In [9]: sns.pairplot(data=df)
plt.show()
```



```
In [10]: sns.heatmap(df.corr(),annot=True)
plt.show()
```



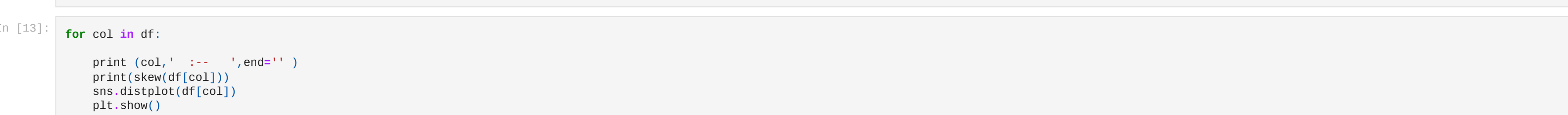
```
In [11]: sns.boxplot(data=df)
```



```
In [12]: from scipy.stats import skew
```

```
In [13]: for col in df:
```

```
    print(col,' :-- ',end='')
    print(skew(df[col]))
    sns.distplot(df[col])
    plt.show()
```



```
In [64]: col=df.drop('Profit',axis=1).columns
```

```
x = df.iloc[:,1:3].values
```

```
y = df.iloc[:,1].values
```

```
In [65]: from sklearn.linear_model import LinearRegression, Lasso,Ridge
```

```
from sklearn.tree import DecisionTreeRegressor
```

```
from sklearn.svm import svm
```

```
from sklearn.neighbors import KNeighborsRegressor
```

```
from sklearn.ensemble import RandomForestRegressor, AdaBoostRegressor, GradientBoostingRegressor
```

```
from xgboost import XGBRegressor
```

```
from sklearn.model_selection import StratifiedFold
```

```
from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
```

```
from sklearn.model_selection import train_test_split, cross_val_score
```

```
In [66]: models=[('Linear Regression',LinearRegression()),
              ('Decision Tree Regressor',DecisionTreeRegressor()),
              ('Lasso Regression ',Lasso()),
              ('Ridge Regression ',Ridge()),
              ('Support vector Regression ',svm(kernel='linear')),
              ('K Nearest Neighbors Regressor',KNeighborsRegressor()),
              ('Random Forest Regressor',RandomForestRegressor()),
              ('AdaBoost Regressor',AdaBoostRegressor()),
              ('Gradient Boosting Regressor', GradientBoostingRegressor()),
              ('Xtreme Gradient Boosting Regressor',XGBRegressor())
            ]
```

```
In [67]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=.28,random_state=9)
```

```
In [68]: accuracy=[]
```

```
for name,model in models:
    print(name,' :-- ')
    print()
    model.fit(x_train,y_train)
    ypred=model.predict(x_test)
    print('      mean absolute error      :-- ',mean_absolute_error(ytest,ypred))
    print()
    print('      mean squared error      :-- ',mean_squared_error(ytest,ypred))
    print()
    print('      squar root mean squared error :-- ',np.sqrt(mean_squared_error(ytest,ypred)))

    print()
    print('      accuracy      :-- ',r2_score(ytest,ypred)*100)
    print()
    accuracy.append(r2_score(ytest,ypred))

Linear Regression :--
mean absolute error      :-- 7328.441614848123
mean squared error      :-- 77566468.16885388
squar root mean squared error :-- 8883.77579046933
accuracy      :-- 93.93955917828573

Decision Tree Regressor :--
mean absolute error      :-- 5232.612080809064
mean squared error      :-- 49697386.56916062
squar root mean squared error :-- 7843.245594834814
accuracy      :-- 96.1218708064904

Lasso Regression :--
mean absolute error      :-- 7328.441565311366
mean squared error      :-- 77566467.1863158
squar root mean squared error :-- 8883.775734667246
accuracy      :-- 93.93955925563303

Ridge Regression :--
mean absolute error      :-- 7328.4416150478255
mean squared error      :-- 77566468.17418559
squar root mean squared error :-- 8883.77579072139
accuracy      :-- 93.93955917778882

Support vector Regression :--
mean absolute error      :-- 11583.344547538087
mean squared error      :-- 168811542.7753044
squar root mean squared error :-- 12881.1481189959597
accuracy      :-- 87.42578966235275

K Nearest Neighbors Regressor :--
mean absolute error      :-- 19447.829480809087
mean squared error      :-- 143358751.24623257
squar root mean squared error :-- 11973.251496143877
accuracy      :-- 88.79639967653898

Random Forest Regressor :--
mean absolute error      :-- 5382.992028686825
mean squared error      :-- 46897136.1528224
squar root mean squared error :-- 6848.148373978356
accuracy      :-- 96.33298581298052

AdaBoost Regressor :--
mean absolute error      :-- 8858.55568080908015
mean squared error      :-- 86273881.35834625
squar root mean squared error :-- 9288.378833781081
accuracy      :-- 93.25480355828945

Gradient Boosting Regressor :--
mean absolute error      :-- 7843.152119635384
mean squared error      :-- 162277207.79570816
squar root mean squared error :-- 12713.219457597899
accuracy      :-- 92.6026679674859

Xtreme Gradient Boosting Regressor :--
mean absolute error      :-- 7916.362283125002
mean squared error      :-- 94474686.90564486
squar root mean squared error :-- 9719.808959327267
accuracy      :-- 92.61276816404741
```

```
In [69]: print('this is the mean accuracy all models we have used :-- ',np.array(accuracy).mean()*100)
```

```
this is the mean accuracy all models we have used :-- 92.83582734232925
```

```
In [70]: model_names=[]
```

```
for name,model in models:
    model_names.append(name)
```

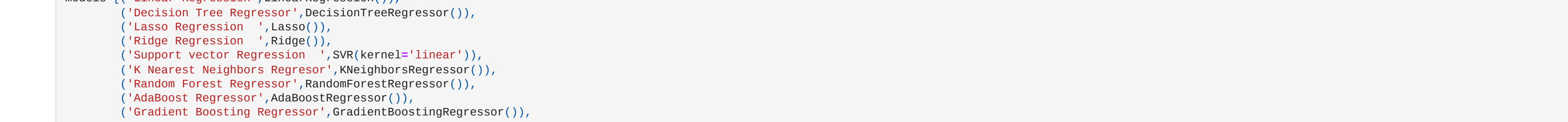
```
plt.figure(figsize=(8,8))
```

```
plt.plot(model_names,accuracy,marker='o')
```

```
plt.grid()
```

```
plt.xticks(rotation=90)
```

```
plt.show()
```



```
In [71]: Random=RandomForestRegressor()
```

```
Random.fit(x_train,y_train)
```

```
ypred=Random.predict(x_test)
```

```
print('training error      :-- ',Random.score(x_train,y_train)*100)
```

```
print('testing error      :-- ',Random.score(x_test,ytest)*100)
```

```
print('square root of mean squared error :-- ',np.sqrt(mean_squared_error(ytest,ypred)))
```

```
print('mean absolute error :-- ',mean_absolute_error(ytest,ypred))
```

```
training error      :-- 98.81488089286113
```

```
testing error      :-- 95.96759184779739
```

```
square root of mean squared error :-- 7382.234464069621
```

```
mean absolute error :-- 5614.063823068934
```

```
In [72]: accuracies = cross_val_score(estimator = RandomForestRegressor(), X = x_train, y = y_train, cv = 5,verbose = 1)
```

```
accuracies
```

```
[Parallel(n_jobs=1): Using backend SequentialBackend with 1 concurrent workers.
```

```
Parallel(n_jobs=1): Done 5 out of 5 | elapsed: 0.4s finished
```

```
array([0.7840862, 0.8894862, 0.8973132, 0.8885861, 0.93883515])
```

```
In [73]: print('this is the minimum accuracy we can get from this model :-- ',accuracies.mean()*100)
```

```
this is the minimum accuracy we can get from this model :-- 87.9229839218835
```

```
In [74]: Random.feature_importances_
```

```
array([0.92788124, 0.06682742])
```

```
In [75]: pd.DataFrame(Random.feature_importances_.col,columns=['co-efficient'])
```

```
Out[75]:
```

	co-efficient
R&D Spend	0.927881
Administration	0.066829
Marketing Spend	0.066027

Above interpretation tell as about

```
In [ ]: setting all values fixed if we increase 1 unit of R&D Spend, then it will increase 0.927881 unit in Profit
```

```
setting all values fixed if we increase 1 unit of Administration, then it will increase 0.066829 unit in Profit
```

```
setting all values fixed if we increase 1 unit of Marketing Spend, then it will increase 0.066027 unit in Profit
```

```
In [ ]:
```