

```
import numpy as np # for performing mathematical calculations behind ML algorithms
import matplotlib.pyplot as plt # for visualization
import pandas as pd # for handling and cleaning the dataset
import seaborn as sns # for visualization
import sklearn # for model evaluation and development
```

```
from google.colab import drive
drive.mount('/content/gdrive')
```

Mounted at /content/gdrive

```
dataset = pd.read_csv('/content/gdrive/MyDrive/50_Startups.csv')
```

```
-----
NameError                                Traceback (most recent call last)
<ipython-input-1-4cd1d9a25275> in <module>()
----> 1 dataset = pd.read_csv('/content/gdrive/MyDrive/50_Startups.csv')

NameError: name 'pd' is not defined
```

SEARCH STACK OVERFLOW

```
dataset.head()
```

	R&D Spend	Administration	Marketing Spend	State	Profit
0	165349.20	136897.80	471784.10	New York	192261.83
1	162597.70	151377.59	443898.53	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

```
dataset.tail()
```

	R&D Spend	Administration	Marketing Spend	State	Profit
45	1000.23	124153.04	1903.93	New York	64926.08
46	1315.46	115816.21	297114.46	Florida	49490.75
47	0.00	135426.92	0.00	California	42559.73
48	542.05	51743.15	0.00	New York	35673.41
49	0.00	116983.80	45173.06	California	14681.40

```
print('There are ',dataset.shape[0], 'rows and ',dataset.shape[1], 'columns in the dataset.')
```

There are 50 rows and 5 columns in the dataset.

```
dataset.isnull().sum()
```

```
R&D Spend      0
Administration  0
Marketing Spend  0
State           0
Profit          0
dtype: int64
```

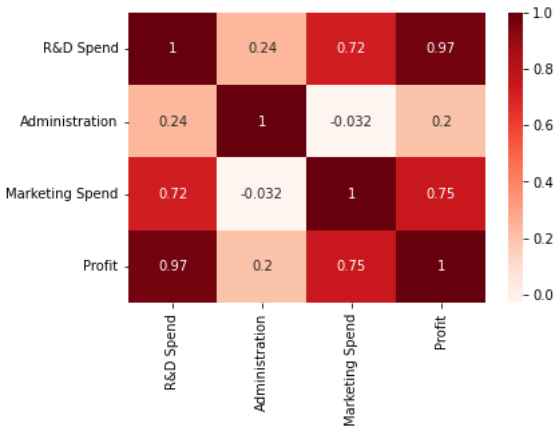
```
dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50 entries, 0 to 49
Data columns (total 5 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   State                50 non-null    object
4   Profit               50 non-null    float64
dtypes: float64(4), object(1)
memory usage: 2.1+ KB
```

```
c = dataset.corr()
c
```

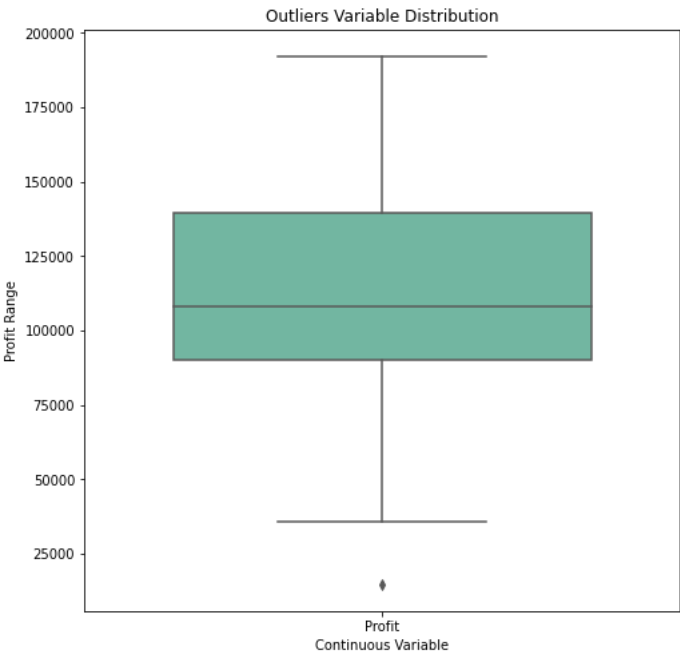
	R&D Spend	Administration	Marketing Spend	Profit
R&D Spend	1.000000	0.241955	0.724248	0.972900
Administration	0.241955	1.000000	0.032154	0.000747
Marketing Spend	0.724248	0.032154	1.000000	0.750000
Profit	0.972900	0.000747	0.750000	1.000000

```
#EDA
sns.heatmap(c,annot=True,cmap='Reds')
plt.show()
```

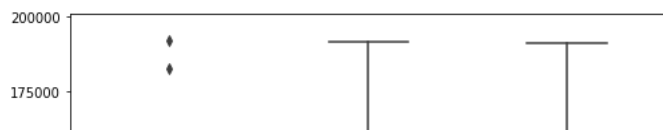


```
outliers = ['Profit']
plt.rcParams['figure.figsize'] = [8,8]
sns.boxplot(data=dataset[outliers], orient="v", palette="Set2" , width=0.7)
plt.title("Outliers Variable Distribution")
plt.ylabel("Profit Range")
plt.xlabel("Continuous Variable")

plt.show()
```

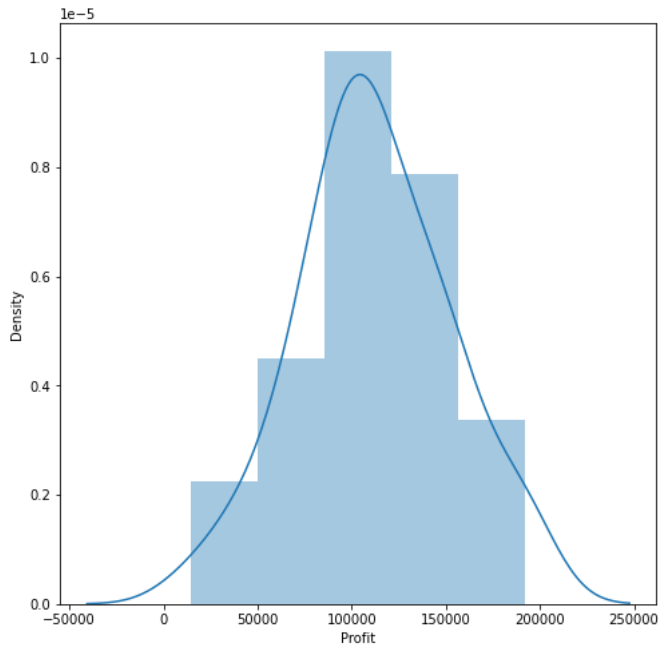


```
sns.boxplot(x = 'State', y = 'Profit', data = dataset)
plt.show()
```



```
#Histogram
sns.distplot(dataset['Profit'],bins=5,kde=True)
plt.show()
```

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning
warnings.warn(msg, FutureWarning)



```
sns.pairplot(dataset)
plt.show()
```

