

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

In [3]: startups = pd.read_csv('50_Startups.csv')
df = startups.copy()

Out[3]:
   R&D Spend  Administration  Marketing Spend   State   Profit
0    165349.20    136897.80      471784.10  New York    182261.83
1    162997.70    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
5    138176.90    99814.71      362861.36  New York    182901.99
6    134615.46    147198.67      127716.82  California    156187.94
7    130291.13    145530.06      329876.68  Florida    156762.60
8    120542.52    148718.95      311613.29  New York    152211.77
9    123334.88    108679.17      304861.62  California    149799.96
10   101913.08    110994.11      229160.95  Florida    146121.95
11   100671.96    91790.61      249744.55  California    144259.40
12   93863.75    127320.38      249639.44  Florida    141585.52
13   91992.39    135495.07      252864.93  California    134307.35
14   119943.24    156547.42      255612.92  Florida    132602.65
15   114523.61    122616.84      261776.23  New York    129917.04
16   70511.11    121997.95      264346.06  California    159992.93
17   64657.16    145077.58      282574.31  New York    125370.37
18   51749.16    114175.79      294919.57  Florida    124266.90
19   86419.79    153514.11      0.00  New York    122776.86
20   76253.86    113967.30      298664.47  California    118474.03
21   78399.47    153773.43      299737.29  New York    111313.02
22   73994.56    122782.75      303319.26  Florida    110392.25
23   67532.53    105751.03      304768.73  Florida    108739.99
24   77044.01    99281.34      140674.81  New York    108552.40
25   64664.71    139553.16      137962.62  California    107404.34
26   75328.87    144135.98      134060.07  Florida    106733.54
27   72107.60    127864.55      353183.81  New York    10506.31
28   66051.52    182645.56      119148.20  Florida    102004.28
29   65025.48    153032.06      107138.38  New York    101004.64
30   61954.48    115641.28      383199.62  New York    99877.99
31   11136.38    152701.92      88218.23  New York    97463.56
32   63408.86    129219.61      46085.25  California    97427.84
33   54893.95    103057.49      214634.81  Florida    97427.84
34   46426.07    157693.92      210797.67  California    96712.80
35   46014.02    85047.44      205517.64  New York    96479.51
36   28663.76    127056.21      201126.82  Florida    96708.19
37   44069.95    51283.14      197029.42  California    89949.14
38   20229.59    65947.93      185265.10  New York    81229.06
39   89587.51    82982.09      174999.30  California    81005.76
40   89793.37    118546.05      177795.67  California    79239.91
41   27982.92    84710.77      164470.71  Florida    77798.83
42   23640.93    96189.63      148091.11  California    71498.49
43   15505.73    127382.30      35534.17  New York    69758.98
44   22177.74    154906.14      28334.72  California    65200.33
45   1000.23    124153.04      1903.93  New York    64826.08
46   1315.46    115816.21      297114.46  Florida    48490.75
47   0.00    135426.92      0.00  California    42599.73
48   542.05    51743.15      0.00  New York    36673.41
49   0.00    116983.80      45173.06  California    14681.40

In [4]: df.head()

Out[4]:
   R&D Spend  Administration  Marketing Spend   State   Profit
0    165349.20    136897.80      471784.10  New York    182261.83
1    162997.70    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

In [5]: df.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 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+ MB

In [7]: df.shape

Out[7]: (50, 5)

In [8]: df.isna().sum()

Out[8]:
R&D Spend      0
Administration  0
Marketing Spend  0
State          0
Profit         0
dtype: int64

In [9]: corr = df.corr()
corr
Out[9]:
           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      0.200717
Marketing Spend  0.724248     -0.032154      1.000000      0.747766
Profit          0.972900      0.200717      0.747766      1.000000

In [10]: sns.heatmap(corr, xticklabels=corr.columns.values,
                    yticklabels=corr.columns.values);

In [12]: sns.scatterplot(x="R&D Spend", y = "Profit", data=df, color="blue")

Out[12]: <AxesSubplot: xlabel='R&D Spend', ylabel='Profit'>

In [13]: df.hist(figsize = (13,10))
plt.show()

In [14]: df.describe().T
count      mean      std      min      25%      50%      75%      max
R&D Spend    50.0    73721.6156    45002.26482    0.00    39936.3700    73051.080    101802.8000    165349.20
Administration    50.0    121344.6396    28017.802755    51283.14    103730.8750    122699.795    144842.1800    182645.56
Marketing Spend    50.0    111025.0978    122290.310726    0.00    129800.1325    121716.240    299469.0850    471784.10
Profit          50.0    112012.6392    40306.180338    14681.40    90138.9025    107978.190    139765.9775    192261.83

In [16]: df_State = pd.get_dummies(df["State"])

In [18]: dfDummies = pd.get_dummies(df["State"], prefix="State")

In [19]: dfDummies

Out[19]:
   State_California  State_Florida  State_New York
0      0      0      1
1      1      0      0
2      0      1      0
3      0      0      1
4      0      1      0
5      0      0      1
6      1      0      0
7      0      1      0
8      0      0      1
9      1      0      0
10     0      1      0
11     1      0      0
12     0      1      0
13     1      0      0
14     0      1      0
15     0      0      1
16     1      0      0
17     0      0      1
18     0      1      0
19     0      0      1
20     1      0      0
21     0      0      1
22     0      1      0
23     0      1      0
24     0      0      1
25     0      1      0
26     0      1      0
27     0      0      1
28     0      1      0
29     0      0      1
30     0      1      0
31     0      0      1
32     1      0      0
33     0      1      0
34     1      0      0
35     0      0      1
36     0      1      0
37     1      0      0
38     0      0      1
39     1      0      0
40     1      0      0
41     0      1      0
42     0      0      1
43     0      0      1
44     1      0      0
45     0      0      1
46     0      1      0
47     1      0      0
48     0      0      1
49     1      0      0

In [21]: df = pd.concat([df, dfDummies], axis=1)
df = df.drop(["State_California", axis = 1)
df = df.drop(["State_Florida", axis = 1)
df.head()

Out[21]:
   R&D Spend  Administration  Marketing Spend   Profit  State_Florida  State_New York
0    165349.20    136897.80      471784.10    182261.83      0      1
1    162997.70    151377.59      443898.53    191792.06      0      0
2    153441.51    101145.55      407934.54    191050.39      1      0
3    144372.41    118671.85      383199.62    182901.99      0      1
4    142107.34     91391.77      366168.42    166187.94      1      0
5    138176.90    99814.71      362861.36      0      1
6    134615.46    147198.67      127716.82      0      0
7    130291.13    145530.06      329876.68      1      0
8    120542.52    148718.95      311613.29      0      1
9    123334.88    108679.17      304861.62      0      0
10   101913.08    110994.11      229160.95      1      0
11   100671.96    91790.61      249744.55      0      0
12   93863.75    127320.38      249639.44      1      0
13   91992.39    135495.07      252864.93      0      0
14   119943.24    156547.42      255612.92      1      0
15   114523.61    122616.84      261776.23      0      1
16   78013.11    121597.55      264346.06      0      0
17   64057.16    145077.58      282574.31      0      1
18   51749.16    114175.79      294919.57      1      0
19   86419.79    153514.11      0.00      0      1
20   76253.86    113967.30      298664.47      0      0
21   78399.47    153773.43      299737.29      0      1
22   73994.56    122782.75      303319.26      1      0
23   67532.53    105751.03      304768.73      1      0
24   77044.01    99281.34      140674.81      0      1
25   64664.71    139553.16      137962.62      0      0
26   75328.87    144135.98      134060.07      1      0
27   72107.60    127864.55      353183.81      0      1
28   66051.52    182645.56      119148.20      1      0
29   65025.48    153032.06      107138.38      0      1
30   61954.48    115641.28      383199.62      1      0
31   11136.38    152701.92      88218.23      0      0
32   63408.86    129219.61      46085.25      0      0
33   54893.95    103057.49      214634.81      1      0
34   46426.07    157693.92      210797.67      0      0
35   46014.02    85047.44      205517.64      0      1
36   28663.76    127056.21      201126.82      1      0
37   44069.95    51283.14      197029.42      0      0
38   20229.59    65947.93      185265.10      0      1
39   89587.51    82982.09      174999.30      0      0
40   89793.37    118546.05      177795.67      0      0
41   27982.92    84710.77      164470.71      1      0
42   23640.93    96189.63      148091.11      0      0
43   15505.73    127382.30      35534.17      0      1
44   22177.74    154906.14      28334.72      0      0
45   1000.23    124153.04      1903.93      0      1
46   1315.46    115816.21      297114.46      1      0
47   0.00    135426.92      0.00      0      0
48   542.05    51743.15      0.00      0      1
49   0.00    116983.80      45173.06      0      0

In [22]: x = df.drop("Profit", axis = 1)
y = df["Profit"]

In [23]: y

Out[23]:
0    192781.83
1    191792.06
2    191050.39
3    182901.99
4    166187.94
5    182901.12
6    156122.91
7    155725.68
8    152711.77
9    149799.96
10   148121.95
11   144259.40
12   141585.52
13   134807.35
14   132560.65
15   125917.04
16   124992.93
17   125379.37
18   124266.90
19   122776.86
20   118474.03
21   113173.02
22   110392.25
23   108739.99
24   108552.40
25   107484.34
26   105733.54
27   10506.31
28   102004.28
29   101004.64
30   99877.99
31   97463.56
32   97427.84
33   96779.51
34   96712.80
35   96479.51
36   96712.80
37   96479.51
38   96712.80
39   81229.06
40   81005.76
41   79239.91
42   77798.83
43   71498.49
44   69758.98
45   65200.33
46   64826.08
47   48490.75
48   42599.73
49   36673.41
Name: Profit, dtype: float64

In [24]: x

Out[24]:
   R&D Spend  Administration  Marketing Spend   State_Florida  State_New York
0    165349.20    136897.80      471784.10      0      1
1    162997.70    151377.59      443898.53      0      0
2    153441.51    101145.55      407934.54      1      0
3    144372.41    118671.85      383199.62      0      1
4    142107.34     91391.77      366168.42      1      0
5    138176.90    99814.71      362861.36      0      1
6    134615.46    147198.67      127716.82      0      0
7    130291.13    145530.06      329876.68      1      0
8    120542.52    148718.95      311613.29      0      1
9    123334.88    108679.17      304861.62      0      0
10   101913.08    110994.11      229160.95      1      0
11   100671.96    91790.61      249744.55      0      0
12   93863.75    127320.38      249639.44      1      0
13   91992.39    135495.07      252864.93      0      0
14   119943.24    156547.42      255612.92      1      0
15   114523.61    122616.84      261776.23      0      1
16   78013.11    121597.55      264346.06      0      0
17   64057.16    145077.58      282574.31      0      1
18   51749.16    114175.79      294919.57      1      0
19   86419.79    153514.11      0.00      0      1
20   76253.86    113967.30      298664.47      0      0
21   78399.47    153773.43      299737.29      0      1
22   73994.56    122782.75      303319.26      1      0
23   67532.53    105751.03      304768.73      1      0
24   77044.01    99281.34      140674.81      0      1
25   64664.71    139553.16      137962.62      0      0
26   75328.87    144135.98      134060.07      1      0
27   72107.60    127864.55      353183.81      0      1
28   66051.52    182645.56      119148.20      1      0
29   65025.48    153032.06      107138.38      0      1
30   61954.48    115641.28      383199.62      1      0
31   11136.38    152701.92      88218.23      0      0
32   63408.86    129219.61      46085.25      0      0
33   54893.95    103057.49      214634.81      1      0
34   46426.07    157693.92      210797.67      0      0
35   46014.02    85047.44      205517.64      0      1
36   28663.76    127056.21      201126.82      1      0
37   44069.95    51283.14      197029.42      0      0
38   20229.59    65947.93      185265.10      0      1
39   89587.51    82982.09      174999.30      0      0
40   89793.37    118546.05      177795.67      0      0
41   27982.92    84710.77      164470.71      1      0
42   23640.93    96189.63      148091.11      0      0
43   15505.73    127382.30      35534.17      0      1
44   22177.74    154906.14      28334.72      0      0
45   1000.23    124153.04      1903.93      0      1
46   1315.46    115816.21      297114.46      1      0
47   0.00    135426.92      0.00      0      0
48   542.05    51743.15      0.00      0      1
49   0.00    116983.80      45173.06      0      0

In [26]: from sklearn.model_selection import train_test_split

In [27]: x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.28, random_state = 42)

In [28]: x_train

Out[28]:
   R&D Spend  Administration  Marketing Spend   State_Florida  State_New York
12    93863.75    127320.38      249639.44      1      0
4    142107.34     91391.77      366168.42      1      0
7    44069.95    51283.14      197029.42      0      0
8    120542.52    148718.95      311613.29      0      0
3    144372.41    118671.85      383199.62      0      1
6    134615.46    147198.67      127716.82      0      0
18   51749.16    114175.79      294919.57      1      0
46   1315.46    115816.21      297114.46      1      0
47   0.00    135426.92      0.00      0      0
15   114523.61    122616.84      261776.23      0      1
9    123334.88    108679.17      304861.62      0      0
16   78013.11    121597.55      264346.06      0      0
24   77044.01    99281.34      140674.81      0      1
34   46426.07    157693.92      210797.67      0      0
10   101913.08    110994.11      229160.95      1      0
11   100671.96    91790.61      249744.55      0      0
13   93863.75    127320.38      249639.44      1      0
12   91992.39    135495.07      252864.93      0      0
14   119943.24    156547.42      255612.92      1      0
17   64057.16    145077.58      282574.31      0      1
18   51749.16    114175.79      294919.57      1      0
19   86419.79    153514.11      0.00      0      1
20   76253.86    113967.30      298664.47      0      0
21   78399.47    153773.43      299737.29      0      1
22   73994.56    122782.75      303319.26      1      0
23   67532.53    105751.03      304768.73      1      0
24   77044.01    99281.34      140674.81      0      1
25   64664.71    139553.16      137962.62      0      0
26   75328.87    144135.98      134060.07      1      0
27   72107.60    127864.55      353183.81      0      1
28   66051.52    182645.56      119148.20      1      0
29   65025.48    153032.06      107138.38      0      1
30   61954.48    115641.28      383199.62      1      0
31   11136.38    152701.92      88218.23      0      0
32   63408.86    129219.61      46085.25      0      0
33   54893.95    103057.49      214634.81      1      0
34   46426.07    157693.92      210797.67      0      0
35   46014.02    85047.44      205517.64      0      1
36   28663.76    127056.21      201126.82      1      0
37   44069.95    51283.14      197029.42      0      0
38   20229.59    65947.93      185265.10      0      1
39   89587.51    82982.09      174999.30      0      0
40   89793.37    118546.05      177795.67      0      0
41   27982.92    84710.77      164470.71      1      0
42   23640.93    96189.63      148091.11      0      0
43   15505.73    127382.30      35534.17      0      1
44   22177.74    154906.14      28334.72      0      0
45   1000.23    124153.04      1903.93      0      1
46   1315.46    115816.21      297114.46      1      0
47   0.00    135426.92      0.00      0      0

In [29]: x_test

Out[29]:
   R&D Spend  Administration  Marketing Spend   State_Florida  State_New York
13   134307.35    126362.87083    78447.0917      0      0
39   81005.76    84608.453036    3602.693636      0      0
30   99937.59    99677.494251    260.095749      0      0
45   64826.08    46357.400686    38568.619314      0      0
17   125370.37    128750.482885    3380.112885      0      0
48   35673.41    50932.417419    135239.007419      0      0
25   107404.34    100643.242816    6761.097184      0      0
32   94227.64    97596.275748    171.435746      0      0
19   122776.86    113997.425344    9679.437566      0      0
Name: Profit, dtype: float64

In [42]: from sklearn.linear_model import LinearRegression
lm = LinearRegression()

In [43]: model = lm.fit(x_train, y_train)

In [45]: y_pred = model.predict(x_test)
y_pred
Out[45]:
array([1192362.87998755, 84668.45383564, 96677.49425147, 46357.40068682,
       128750.48288504, 50912.417418, 109741.36827202, 106643.24281647,
       97596.27574854, 113997.42534321])

In [47]: df['pd.DataFrame({"Main values": y_test, "Predict values": y_pred, "Difference": abs(y_pred-y_test)})']
df
Out[47]:
   Main values  Predict values  Difference
13    134307.35    126362.87083    7844.70917
39     81005.76    84608.453036    3602.693636
30     99937.59    99677.494251     260.095749
45     64826.08    46357.400686    18468.619314
17    125370.37    128750.482885     3380.112885
48     35673.41    50932.417419    15259.007419
25    107404.34    100643.242816     6761.097184
32     94227.64
```