

```

1 # TASK 1
2 # 1. Import thư viện cần thiết
3 import os
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7
8 # 2. Định nghĩa tên file
9 filename = "50_Startups.csv"
10
11 # 3. Kiểm tra nếu file chưa tồn tại thì yêu cầu tải lên
12 if not os.path.exists(filename):
13     from google.colab import files
14     uploaded = files.upload()
15     filename = list(uploaded.keys())[0] # Lấy tên file vừa tải lên
16
17 # 4. Đọc file CSV vào DataFrame
18 df = pd.read_csv(filename)
19
20 # 5. Loại bỏ các cột chứa dữ liệu không phải số
21 df_numeric = df.select_dtypes(include=[np.number]) # Chỉ giữ lại các cột số
22
23 # 6. Hiển thị 5 dòng đầu tiên để kiểm tra dữ liệu
24 print("\n💡 Dữ liệu số ban đầu:")
25 print(df_numeric.head())
26
27 # 7. Tính toán các thống kê cơ bản
28 summary = df_numeric.describe()
29 print("\n💡 Thống kê dữ liệu:\n", summary)
30
31 # 8. Vẽ scatter plot giữa từng feature và biến phụ thuộc y
32 y = df_numeric.iloc[:, -1] # Cột cuối cùng là biến phụ thuộc (y)
33 features = df_numeric.columns[:-1] # Các cột còn lại là feature
34
35 for feature in features:
36     plt.scatter(df_numeric[feature], y)
37     plt.xlabel(feature)
38     plt.ylabel("Output (y)")
39     plt.title(f'Scatter plot: {feature} vs Output')
40     plt.show()
41
42 # 9. Vẽ histogram của từng feature
43 for feature in df_numeric.columns:
44     mean = df_numeric[feature].mean()
45     var = df_numeric[feature].var()
46     plt.hist(df_numeric[feature], bins=20, alpha=0.7, color='blue')
47     plt.title(f'Histogram of {feature} (mean={mean:.1f}, var={var:.1f})')
48     plt.xlabel(feature)
49     plt.ylabel("Frequency")
50     plt.show()
51
52 # 10. Tính hệ số tương quan giữa các biến
53 correlations = df_numeric.corr()
54 print("\n💡 Hệ số tương quan giữa y và các feature:")
55 print(correlations.iloc[-1, :-1]) # In hệ số tương quan của y với các feature
56
57 # 11. Chuẩn hóa dữ liệu
58 # a. Minmax scaling (đưa dữ liệu về khoảng [0,1])
59 X_minmax = (df_numeric.iloc[:, :-1] - df_numeric.iloc[:, :-1].min()) / (df_numeric.iloc[:, :-1].max() - df_numeric.iloc[:, :-1].min())
60
61 # b. Standardization (chuẩn hóa sao cho mean=0, variance=1)
62 X_standardized = (df_numeric.iloc[:, :-1] - df_numeric.iloc[:, :-1].mean()) / df_numeric.iloc[:, :-1].std()
63
64 # 12. Chuyển đổi dữ liệu thành X1 (min=1, max=10)
65 X1 = df_numeric.iloc[:, :-1].apply(lambda x: 1 + 9 * (x - x.min()) / (x.max() - x.min()))
66
67 # 13. Chuyển đổi dữ liệu thành X2 bằng chuẩn hóa (standardize)
68 X2 = X_standardized.copy()
69
70 # 14. In kết quả để kiểm tra
71 print("\n💡 Dữ liệu sau MinMax Scaling:\n", X_minmax.head())
72 print("\n💡 Dữ liệu sau Standardization:\n", X_standardized.head())
73 print("\n💡 Dữ liệu sau chuyển đổi thành X1:\n", X1.head())
74 print("\n💡 Dữ liệu sau chuẩn hóa thành X2:\n", X2.head())
75

```



Choose Files 50\_Startups.csv

• **50\_Startups.csv**(text/csv) - 2436 bytes, last modified: 2/14/2025 - 100% done  
 Saving 50\_Startups.csv to 50\_Startups.csv

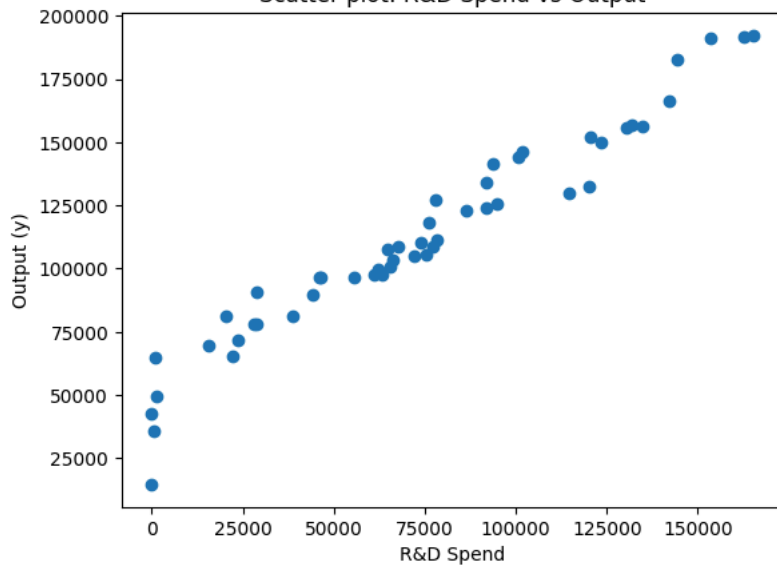
🔴 Dữ liệu số ban đầu:

	R&D Spend	Administration	Marketing Spend	Profit
0	165349.20	136897.80	471784.10	192261.83
1	162597.70	151377.59	443898.53	191792.06
2	153441.51	101145.55	407934.54	191050.39
3	144372.41	118671.85	383199.62	182901.99
4	142107.34	91391.77	366168.42	166187.94

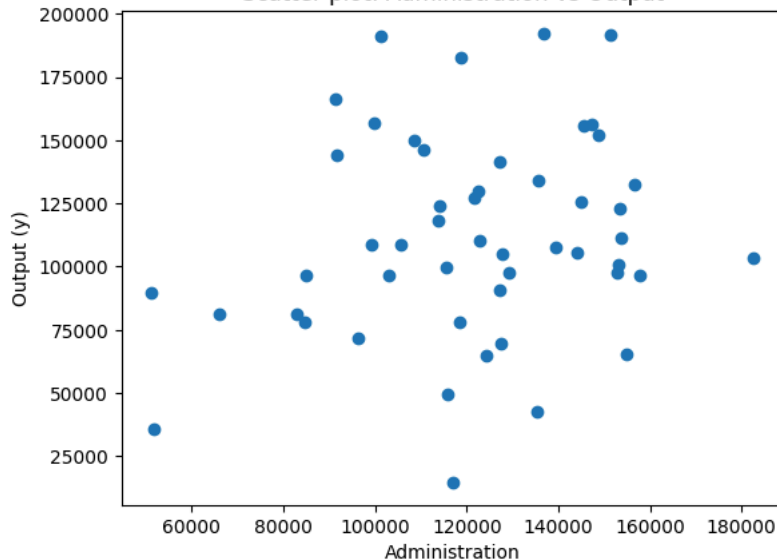
🔴 Thống kê dữ liệu:

	R&D Spend	Administration	Marketing Spend	Profit
count	50.000000	50.000000	50.000000	50.000000
mean	73721.615600	121344.639600	211025.097800	112012.639200
std	45902.256482	28017.802755	122290.310726	40306.180338
min	0.000000	51283.140000	0.000000	14681.400000
25%	39936.370000	103730.875000	129300.132500	90138.902500
50%	73051.080000	122699.795000	212716.240000	107978.190000
75%	101602.800000	144842.180000	299469.085000	139765.977500
max	165349.200000	182645.560000	471784.100000	192261.830000

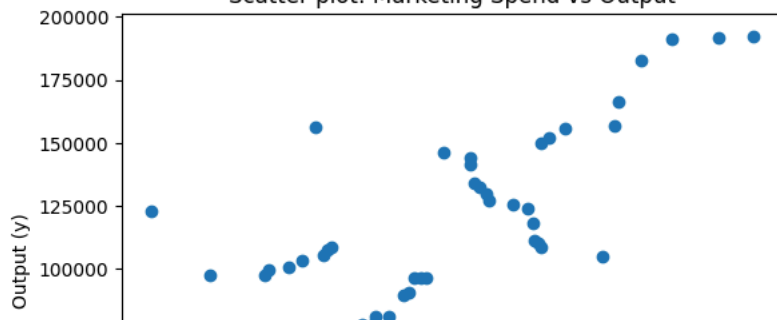
Scatter plot: R&amp;D Spend vs Output

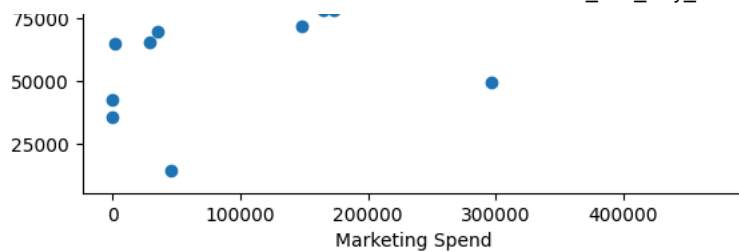


Scatter plot: Administration vs Output

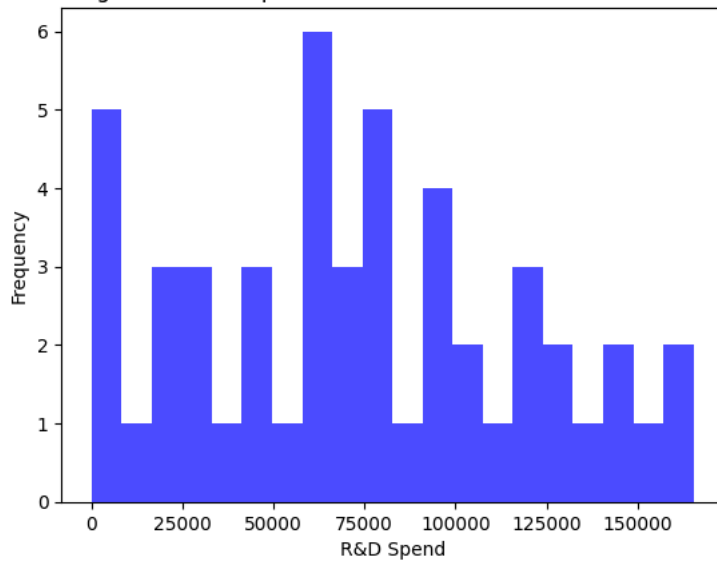


Scatter plot: Marketing Spend vs Output

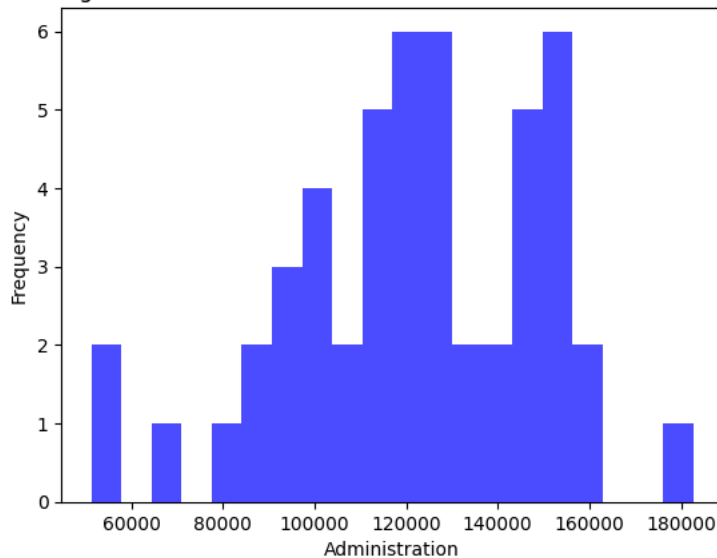




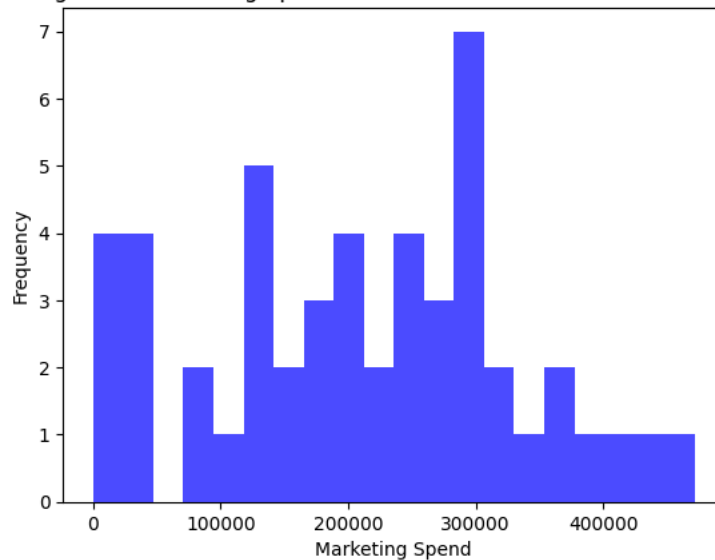
Histogram of R&amp;D Spend (mean=73721.6, var=2107017150.2)



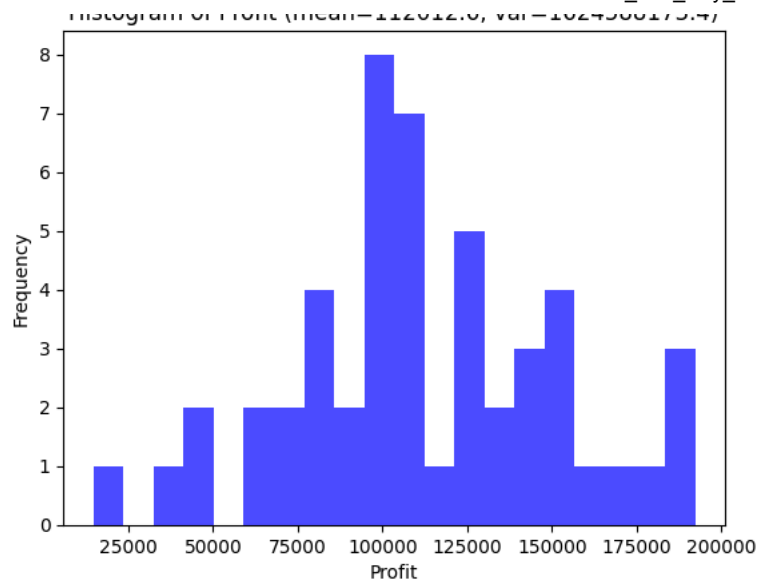
Histogram of Administration (mean=121344.6, var=784997271.2)



Histogram of Marketing Spend (mean=211025.1, var=14954920097.4)



Histogram of Profit (mean=112012.6, var=1624588173.4)



✦ Hệ số tương quan giữa y và các feature:

R&D Spend 0.972900

Administration 0.200717

Marketing Spend 0.747766

Name: Profit, dtype: float64

✦ Dữ liệu sau MinMax Scaling:

	R&D Spend	Administration	Marketing Spend
0	1.000000	0.651744	1.000000
1	0.983359	0.761972	0.940893
2	0.927985	0.379579	0.864664
3	0.873136	0.512998	0.812235
4	0.859438	0.305328	0.776136

✦ Dữ liệu sau Standardization:

	R&D Spend	Administration	Marketing Spend
0	1.996146	0.555117	2.132295
1	1.936203	1.071924	1.904267
2	1.736731	-0.720938	1.610180
3	1.539157	-0.095396	1.407916
4	1.489812	-1.069066	1.268648

✦ Dữ liệu sau chuyển đổi thành X1:

	R&D Spend	Administration	Marketing Spend
0	10.000000	6.865695	10.000000
1	9.850235	7.857746	9.468040
2	9.351861	4.416211	8.781972
3	8.858228	5.616986	8.310116
4	8.734940	3.747952	7.985220

✦ Dữ liệu sau chuẩn hóa thành X2:

	R&D Spend	Administration	Marketing Spend
0	1.996146	0.555117	2.132295
1	1.936203	1.071924	1.904267
2	1.736731	-0.720938	1.610180
3	1.539157	-0.095396	1.407916
4	1.489812	-1.069066	1.268648