

1. Perform a complete data inspection including – (20 marks)
 - a. Missing Data Treatment
 - b. Descriptive Statistics of each variable (Eg. Boxplot, Histogram etc.)
 - c. Visualization of all continuous variables

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

"""
1. Perform a complete data inspection including
a. Missing Data Treatment
b. Descriptive Statistics of each variable (Eg.
Boxplot, Histogram etc.)
c. Visualization of all continuous variables
"""

df = pd.read_csv('Dataset_Day5.csv')
print(df.info())
skewness = df.skew()
print(skewness)
missing_value_percent = df.isna().sum() / len(df) *
100
print(missing_value_percent)
descriptive_statistics = df.describe()
print(descriptive_statistics)
# proportion of non-retail business acres per town
(boxplot is also done for continuous variables)
sns.boxplot(data=df["INDUS"], orient='v')
plt.plot()
# sns.histplot(df, x='NOX')
# plt.plot()
# for continuous variables
sns.pairplot(df[['CRIM', 'ZN', 'INDUS', 'NOX',
'RM', 'AGE', 'DIS', 'PTRATIO', 'B', 'LSTAT',
'MEDV']])
plt.show()
# plot between CRIM(per capita crime rate by town)
and B(1000(bk-0.63)^2 bk is proportion of black
people in town)
```

```
plt.scatter(df['B'], df['CRIM'])
plt.xlabel('B')
plt.ylabel('CRIM')
plt.title('Scatter Plot of B vs CRIM')
plt.show()
sns.heatmap(df.corr())
plt.show()
```

```

1 C:\Users\tejas\PycharmProjects\pythonProject\venv\Scripts\python.exe C:\Users\
tejas\PycharmProjects\pythonProject\START\Q1.py
2 <class 'pandas.core.frame.DataFrame'>
3 RangeIndex: 506 entries, 0 to 505
4 Data columns (total 14 columns):
5 #   Column      Non-Null Count  Dtype
6 ---  -
7 0   CRIM        506 non-null    float64
8 1   ZN          506 non-null    float64
9 2   INDUS       506 non-null    float64
10 3   CHAS        506 non-null    float64
11 4   NOX         506 non-null    float64
12 5   RM          506 non-null    float64
13 6   AGE         506 non-null    float64
14 7   DIS         506 non-null    float64
15 8   RAD         506 non-null    float64
16 9   TAX         506 non-null    float64
17 10  PTRATIO     506 non-null    float64
18 11  B           506 non-null    float64
19 12  LSTAT       506 non-null    float64
20 13  MEDV        506 non-null    float64
21 dtypes: float64(14)
22 memory usage: 55.5 KB
23 None
24 CRIM          5.223149
25 ZN            2.225666
26 INDUS        0.295022
27 CHAS          3.405904
28 NOX           0.729308
29 RM            0.403612
30 AGE          -0.598963
31 DIS           1.011781
32 RAD           1.004815
33 TAX           0.669956
34 PTRATIO      -0.802325
35 B            -2.890374
36 LSTAT         0.906460
37 MEDV          1.108098
38 dtype: float64
39 CRIM          0.0
40 ZN            0.0
41 INDUS         0.0
42 CHAS          0.0
43 NOX           0.0
44 RM            0.0
45 AGE           0.0
46 DIS           0.0
47 RAD           0.0
48 TAX           0.0
49 PTRATIO       0.0
50 B             0.0
51 LSTAT         0.0
52 MEDV          0.0
53 dtype: float64
54      CRIM      ZN      INDUS  ...      B      LSTAT
55 count  506.000000  506.000000  506.000000  ...  506.000000  506.000000  506.
000000
56 mean    3.613524   11.363636   11.136779  ...  356.674032   12.653063   22.

```

```

56 532806
57 std      8.601545    23.322453    6.860353    ...    91.294864    7.141062    9.
   197104
58 min      0.006320    0.000000    0.460000    ...    0.320000    1.730000    5.
   000000
59 25%      0.082045    0.000000    5.190000    ...    375.377500    6.950000    17.
   025000
60 50%      0.256510    0.000000    9.690000    ...    391.440000    11.360000    21.
   200000
61 75%      3.677083    12.500000    18.100000    ...    396.225000    16.955000    25.
   000000
62 max      88.976200    100.000000    27.740000    ...    396.900000    37.970000    50.
   000000
63
64 [8 rows x 14 columns]
65

```

Figure 1

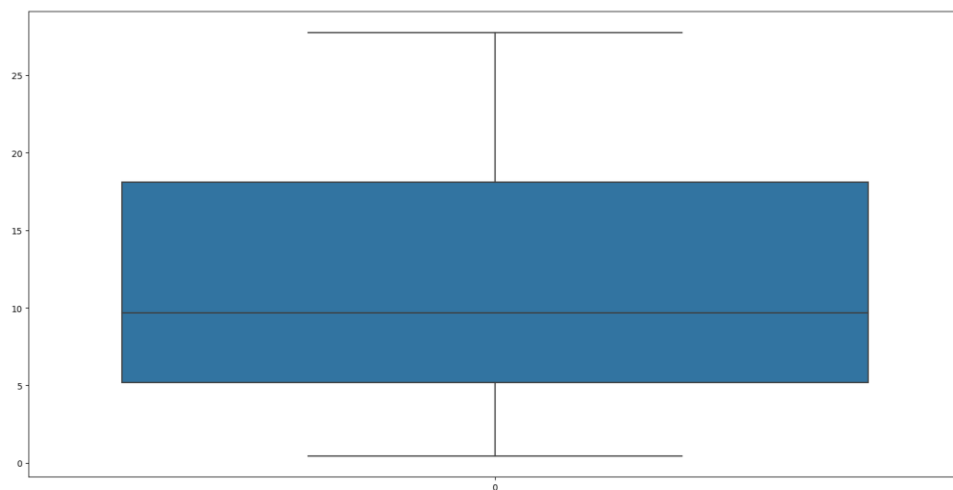


Figure 2

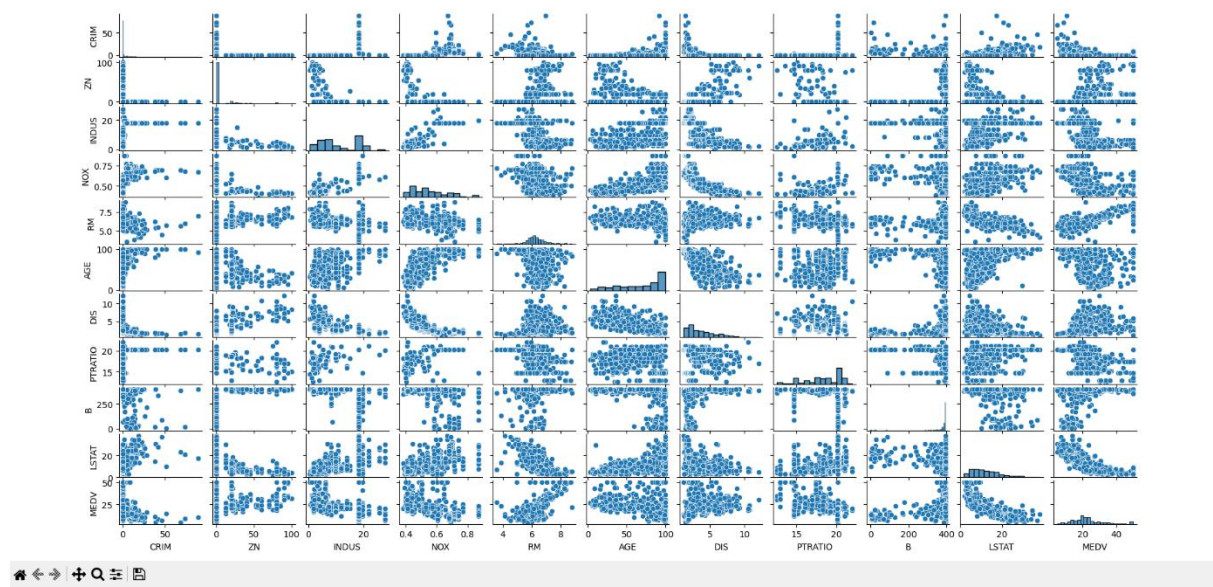


Figure 1

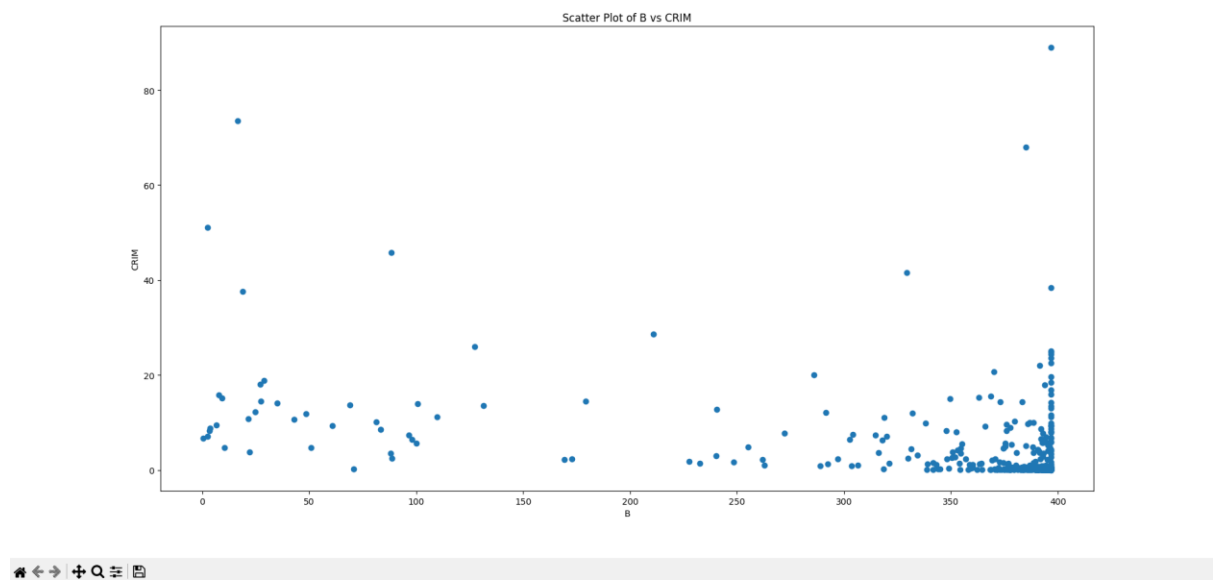
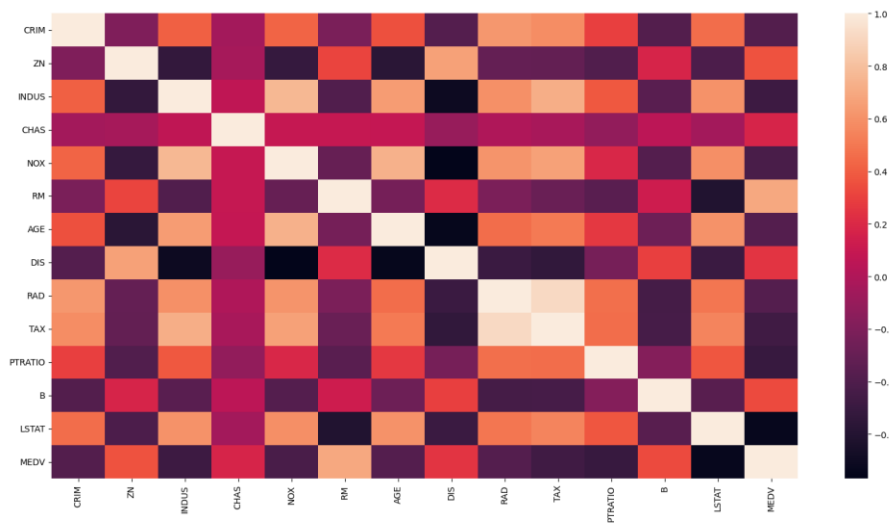


Figure 1



2. Create a simple linear regression model that quantitatively relates 'MEDV' with 'RM'. (10 marks)

- Share the model performance metrics and print the full regression model with coefficients.
- Use the model to predict the price of the house for 'RM' = 7

```
2. import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import
LinearRegression
from sklearn.model_selection import
train_test_split
from sklearn.metrics import mean_squared_error,
r2_score, mean_absolute_error

"""
2. Create a simple linear regression model
that quantitatively relates 'MEDV' with 'RM'.
(10 marks)
a. Share the model performance metrics and
print the full regression model with
coefficients.
b. Use the model to predict the price of the
house for 'RM' = 7
"""
```

```
df = pd.read_csv('Dataset_Day5.csv')
print(df.info())
X = df[['RM']]
y = df[['MEDV']]
# split
X_train, X_test, y_train, y_test =
train_test_split(X, y, test_size=0.2,
random_state=42)
# calling and fitting the model on training
data
lm = LinearRegression()
lm = lm.fit(X_train, y_train)
# predict using test set
y_pred = lm.predict(X_test)
print(lm.coef_) # scale parameter
print(lm.intercept_) # intercept parameter
print(r2_score(y_test, y_pred))
print(mean_absolute_error(y_test, y_pred))
print(mean_squared_error(y_test, y_pred))
rm_value = 7
price_prediction = lm.predict([[rm_value]])
print("Predicted price for 'RM' = 7:",
price_prediction[0])
```

```

1 C:\Users\tejas\PycharmProjects\pythonProject\venv\Scripts\python.exe C:\Users\
  tejas\PycharmProjects\pythonProject\START\Q2.py
2 <class 'pandas.core.frame.DataFrame'>
3 RangeIndex: 506 entries, 0 to 505
4 Data columns (total 14 columns):
5 #   Column      Non-Null Count  Dtype
6 ---  -
7 0    CRIM        506 non-null    float64
8 1    ZN          506 non-null    float64
9 2    INDUS       506 non-null    float64
10 3    CHAS        506 non-null    float64
11 4    NOX         506 non-null    float64
12 5    RM          506 non-null    float64
13 6    AGE         506 non-null    float64
14 7    DIS         506 non-null    float64
15 8    RAD         506 non-null    float64
16 9    TAX         506 non-null    float64
17 10   PTRATIO     506 non-null    float64
18 11   B           506 non-null    float64
19 12   LSTAT       506 non-null    float64
20 13   MEDV        506 non-null    float64
21 dtypes: float64(14)
22 memory usage: 55.5 KB
23 None
24 [[9.34830141]]
25 [-36.2463189]
26 0.3707569232254778
27 4.478335832064147
28 46.144775347317264
29 Predicted price for 'RM' = 7: [29.19179095]
30 C:\Users\tejas\PycharmProjects\pythonProject\venv\Lib\site-packages\sklearn\
  base.py:464: UserWarning: X does not have valid feature names, but
  LinearRegression was fitted with feature names
31   warnings.warn(
32
33 Process finished with exit code 0

```

3. Create a simple linear regression model that quantitatively relates 'MEDV' with 'DIS'. (10 marks)

- a. Share the model performance metrics and print the full regression model with coefficients.
- b. Use the model to predict the price of the house for 'DIS' = 15

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import
train_test_split
from sklearn.metrics import mean_squared_error,
r2_score, mean_absolute_error

```



```
"""
```

3. Create a simple linear regression model that quantitatively relates 'MEDV' with 'DIS'. (10 marks)

- a. Share the model performance metrics and print the full regression model with coefficients.
- b. Use the model to predict the price of the house for 'DIS' = 15

```
"""
```

```
df = pd.read_csv('Dataset_Day5.csv')
print(df.info())
X = df[['DIS']]
y = df['MEDV']
X_train, X_test, y_train, y_test =
train_test_split(X, y, test_size=0.2,
random_state=42)
lm = LinearRegression()
lm = lm.fit(X_train, y_train)
y_pred = lm.predict(X_test)
print(lm.coef_) # scale parameter
print(lm.intercept_) # intercept parameter
print(r2_score(y_test, y_pred))
print(mean_absolute_error(y_test, y_pred))
print(mean_squared_error(y_test, y_pred))
DIS_value = 7
price_prediction = lm.predict([[DIS_value]])
print("Predicted price for 'DIS' = 15:",
price_prediction[0])
```

```
1 C:\Users\tejas\PycharmProjects\pythonProject\venv\Scripts\python.exe C:\Users\
  tejas\PycharmProjects\pythonProject\START\Q3.py
2 <class 'pandas.core.frame.DataFrame'>
3 RangeIndex: 506 entries, 0 to 505
4 Data columns (total 14 columns):
5 #   Column      Non-Null Count  Dtype
6 ---  -
7 0    CRIM        506 non-null    float64
8 1    ZN          506 non-null    float64
9 2    INDUS       506 non-null    float64
10 3    CHAS        506 non-null    float64
11 4    NOX         506 non-null    float64
12 5    RM          506 non-null    float64
13 6    AGE         506 non-null    float64
14 7    DIS         506 non-null    float64
15 8    RAD         506 non-null    float64
16 9    TAX         506 non-null    float64
17 10   PTRATIO    506 non-null    float64
18 11   B          506 non-null    float64
19 12   LSTAT      506 non-null    float64
20 13   MEDV       506 non-null    float64
21 dtypes: float64(14)
22 memory usage: 55.5 KB
23 None
24 [1.0295094]
25 18.875962058273238
26 0.07332042069244615
27 5.967846118518974
28 67.95691932803946
29 Predicted price for 'DIS' = 15: 26.08252784851876
30 C:\Users\tejas\PycharmProjects\pythonProject\venv\Lib\site-packages\sklearn\
  base.py:464: UserWarning: X does not have valid feature names, but
  LinearRegression was fitted with feature names
31   warnings.warn(
32
33 Process finished with exit code 0
```