

TI35_DSBDA_5th_HeartDisease

January 30, 2025

```
[1]: import pandas as pd
```

```
[3]: df=pd.read_csv("/home/bcl07/heart_disease.csv")
```

```
[4]: df
```

```
[4]:
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	\
0	52	1	0	125	212	0	1	168	0	1.0	
1	53	1	0	140	203	1	0	155	1	3.1	
2	70	1	0	145	174	0	1	125	1	2.6	
3	61	1	0	148	203	0	1	161	0	0.0	
4	62	0	0	138	294	1	1	106	0	1.9	
...
1020	59	1	1	140	221	0	1	164	1	0.0	
1021	60	1	0	125	258	0	0	141	1	2.8	
1022	47	1	0	110	275	0	0	118	1	1.0	
1023	50	0	0	110	254	0	0	159	0	0.0	
1024	54	1	0	120	188	0	1	113	0	1.4	

	slope	ca	thal	target
0	2	2	3	0
1	0	0	3	0
2	0	0	3	0
3	2	1	3	0
4	1	3	2	0
...
1020	2	0	2	1
1021	1	1	3	0
1022	1	1	2	0
1023	2	0	2	1
1024	1	1	3	0

[1025 rows x 14 columns]

```
[5]: df.columns
```

```
[5]: Index(['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg', 'thalach',
          'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target'],
          dtype='object')
```

```
[6]: df.isnull().sum()
```

```
[6]: age      0
     sex      0
     cp       0
     trestbps  0
     chol     0
     fbs      0
     restecg  0
     thalach  0
     exang    0
     oldpeak  0
     slope    0
     ca       0
     thal     0
     target   0
     dtype: int64
```

```
[8]: df=df.drop_duplicates()
```

```
[9]: df.describe()
```

```
[9]:
```

	age	sex	cp	trestbps	chol	fbs	\
count	302.00000	302.000000	302.000000	302.000000	302.000000	302.000000	
mean	54.42053	0.682119	0.963576	131.602649	246.500000	0.149007	
std	9.04797	0.466426	1.032044	17.563394	51.753489	0.356686	
min	29.00000	0.000000	0.000000	94.000000	126.000000	0.000000	
25%	48.00000	0.000000	0.000000	120.000000	211.000000	0.000000	
50%	55.50000	1.000000	1.000000	130.000000	240.500000	0.000000	
75%	61.00000	1.000000	2.000000	140.000000	274.750000	0.000000	
max	77.00000	1.000000	3.000000	200.000000	564.000000	1.000000	

	restecg	thalach	exang	oldpeak	slope	ca	\
count	302.000000	302.000000	302.000000	302.000000	302.000000	302.000000	
mean	0.526490	149.569536	0.327815	1.043046	1.397351	0.718543	
std	0.526027	22.903527	0.470196	1.161452	0.616274	1.006748	
min	0.000000	71.000000	0.000000	0.000000	0.000000	0.000000	
25%	0.000000	133.250000	0.000000	0.000000	1.000000	0.000000	
50%	1.000000	152.500000	0.000000	0.800000	1.000000	0.000000	
75%	1.000000	166.000000	1.000000	1.600000	2.000000	1.000000	
max	2.000000	202.000000	1.000000	6.200000	2.000000	4.000000	

	thal	target
--	------	--------

count	302.000000	302.000000
mean	2.314570	0.543046
std	0.613026	0.498970
min	0.000000	0.000000
25%	2.000000	0.000000
50%	2.000000	1.000000
75%	3.000000	1.000000
max	3.000000	1.000000

```
[10]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 302 entries, 0 to 878
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   age         302 non-null    int64
1   sex         302 non-null    int64
2   cp          302 non-null    int64
3   trestbps    302 non-null    int64
4   chol        302 non-null    int64
5   fbs         302 non-null    int64
6   restecg     302 non-null    int64
7   thalach     302 non-null    int64
8   exang       302 non-null    int64
9   oldpeak     302 non-null    float64
10  slope       302 non-null    int64
11  ca          302 non-null    int64
12  thal        302 non-null    int64
13  target      302 non-null    int64
dtypes: float64(1), int64(13)
memory usage: 35.4 KB
```

```
[11]: df.isna().sum()
```

```
[11]: age         0
sex         0
cp          0
trestbps    0
chol        0
fbs         0
restecg     0
thalach     0
exang       0
oldpeak     0
slope       0
ca          0
```

```

thal      0
target    0
dtype: int64

```

```
[12]: df.head()
```

```

[12]:   age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope  \
0    52   1   0     125    212   0         1     168     0       1.0     2
1    53   1   0     140    203   1         0     155     1       3.1     0
2    70   1   0     145    174   0         1     125     1       2.6     0
3    61   1   0     148    203   0         1     161     0       0.0     2
4    62   0   0     138    294   1         1     106     0       1.9     1

      ca  thal  target
0     2     3        0
1     0     3        0
2     0     3        0
3     1     3        0
4     3     2        0

```

```
[13]: df.fbs.unique()
```

```
[13]: array([0, 1])
```

```
[16]: subSet1 = df[['age','cp','chol','thal']]
```

```
[18]: subSet2 = df[['exang','slope','target']]
```

```
[19]: merged_df = subSet1.merge(right=subSet2,how='cross')
merged_df.head()
```

```

[19]:   age  cp  chol  thal  exang  slope  target
0    52   0  212     3     0     2         0
1    52   0  212     3     1     0         0
2    52   0  212     3     1     0         0
3    52   0  212     3     0     2         0
4    52   0  212     3     0     1         0

```

```
[20]: df.columns
```

```

[20]: Index(['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg', 'thalach',
            'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target'],
            dtype='object')

```

```

[21]: def remove_outliers(column):
      Q1 = column.quantile(0.25)
      Q3 = column.quantile(0.75)

```

```

IQR = Q3 - Q1
threshold = 1.5 * IQR
outlier_mask = (column < Q1 - threshold) | (column > Q3 + threshold)
return column[~outlier_mask]

```

```

[23]: col_name = ['cp', 'thal', 'exang', 'oldpeak', 'slope', 'ca']
for col in col_name:
    df[col] = remove_outliers(df[col])

```

/tmp/ipykernel_10564/1228815343.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
df[col] = remove_outliers(df[col])
```

/tmp/ipykernel_10564/1228815343.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
df[col] = remove_outliers(df[col])
```

/tmp/ipykernel_10564/1228815343.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
df[col] = remove_outliers(df[col])
```

/tmp/ipykernel_10564/1228815343.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
df[col] = remove_outliers(df[col])
```

/tmp/ipykernel_10564/1228815343.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
df[col] = remove_outliers(df[col])
```

/tmp/ipykernel_10564/1228815343.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

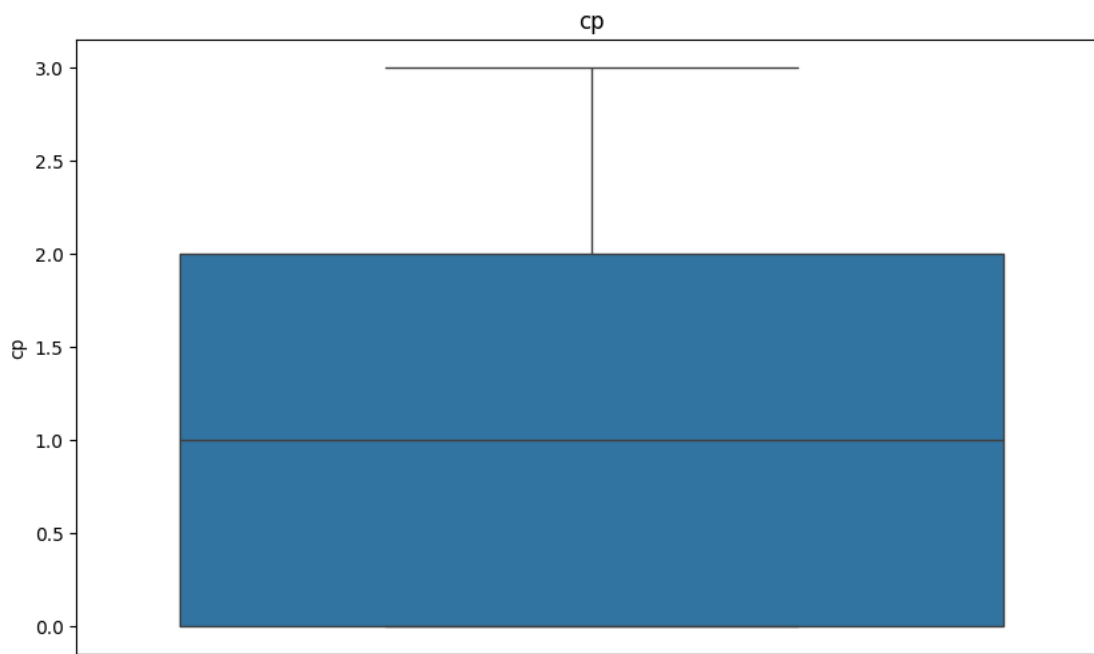
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

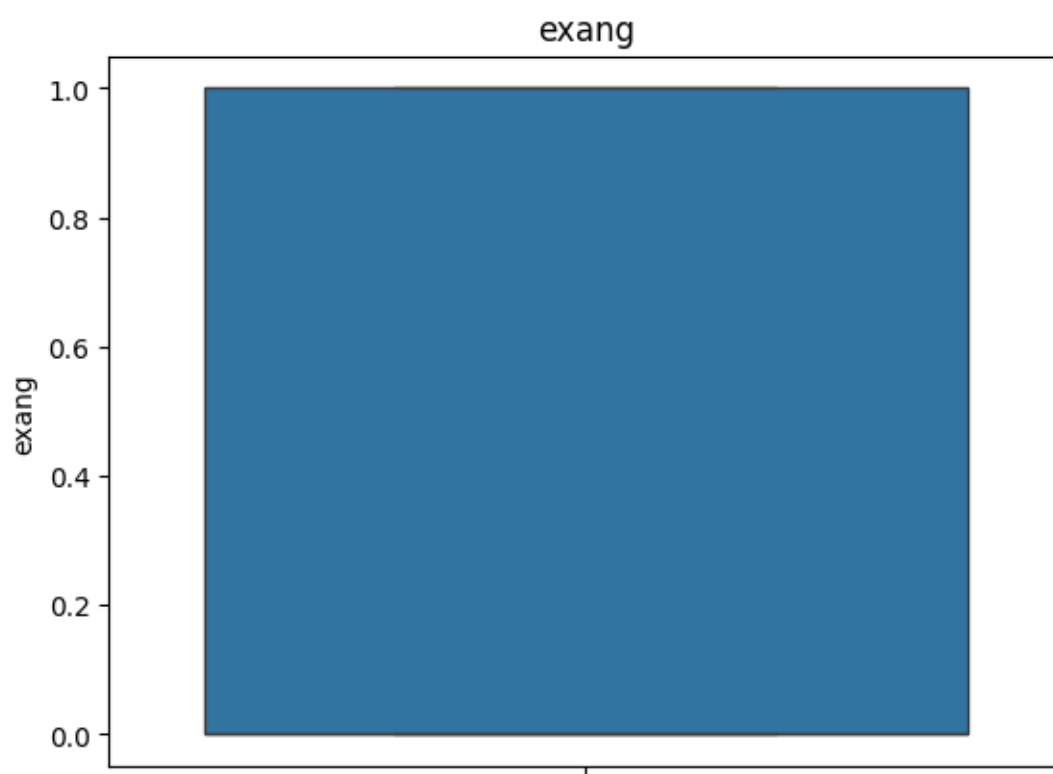
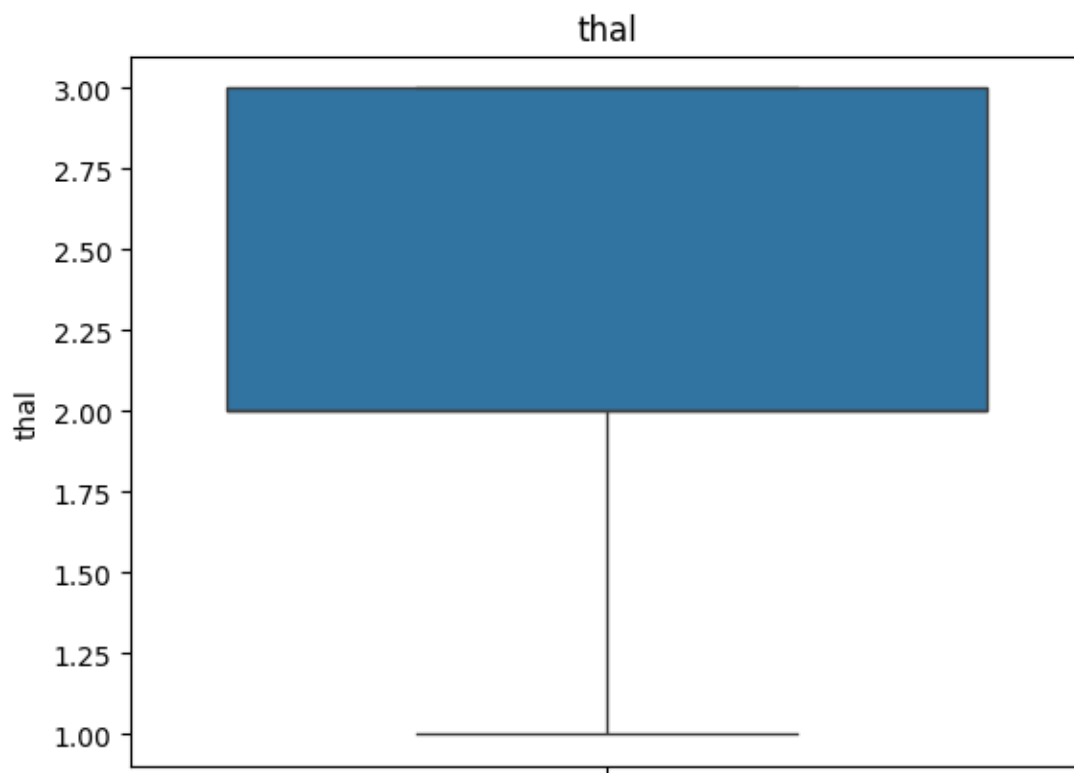
```
df[col] = remove_outliers(df[col])
```

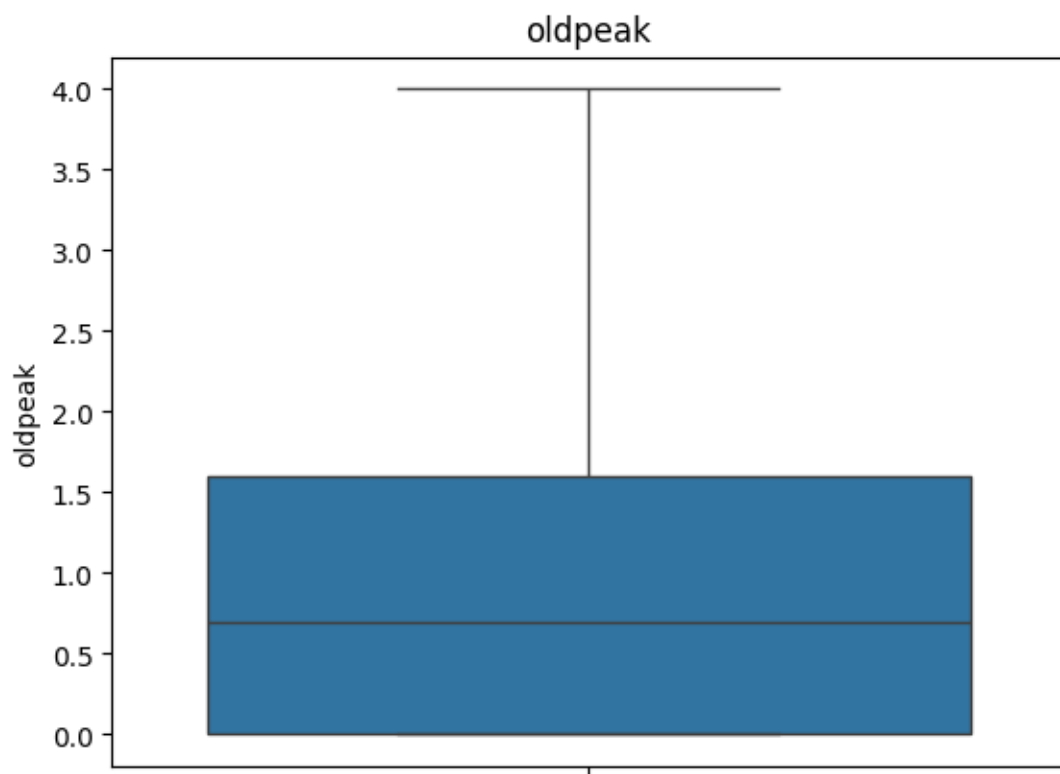
```
[28]: from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import accuracy_score, confusion_matrix
from sklearn.linear_model import LogisticRegression
import seaborn as sns
import matplotlib.pyplot as plt
```

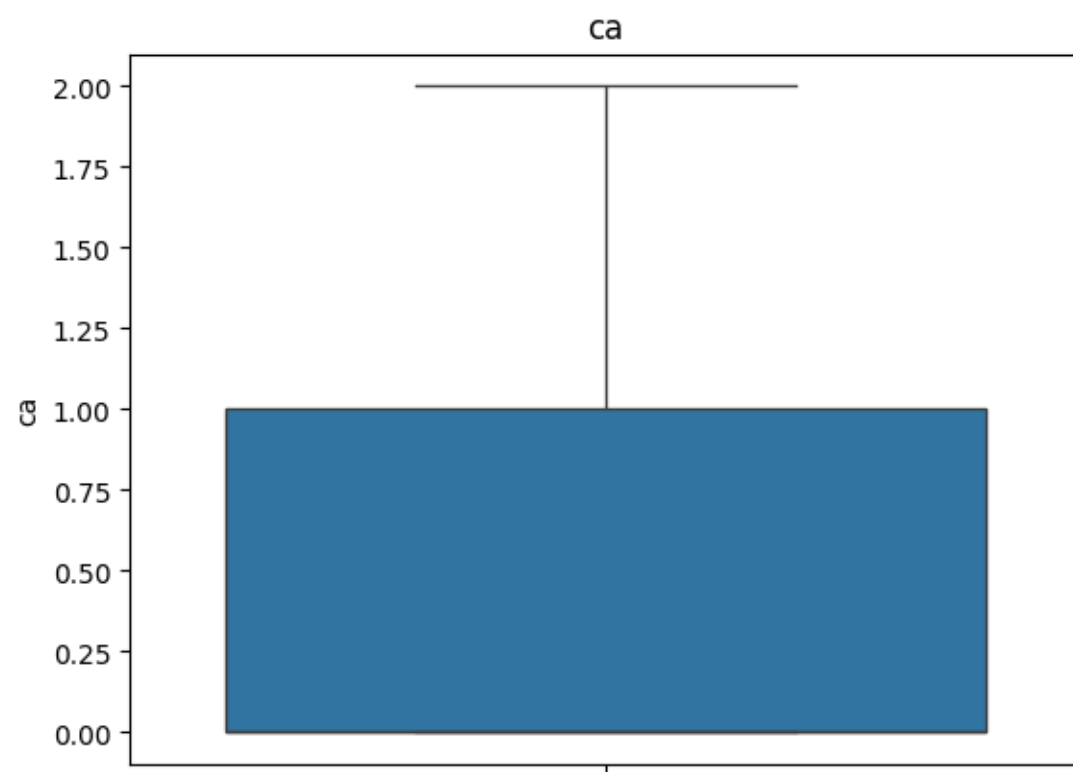
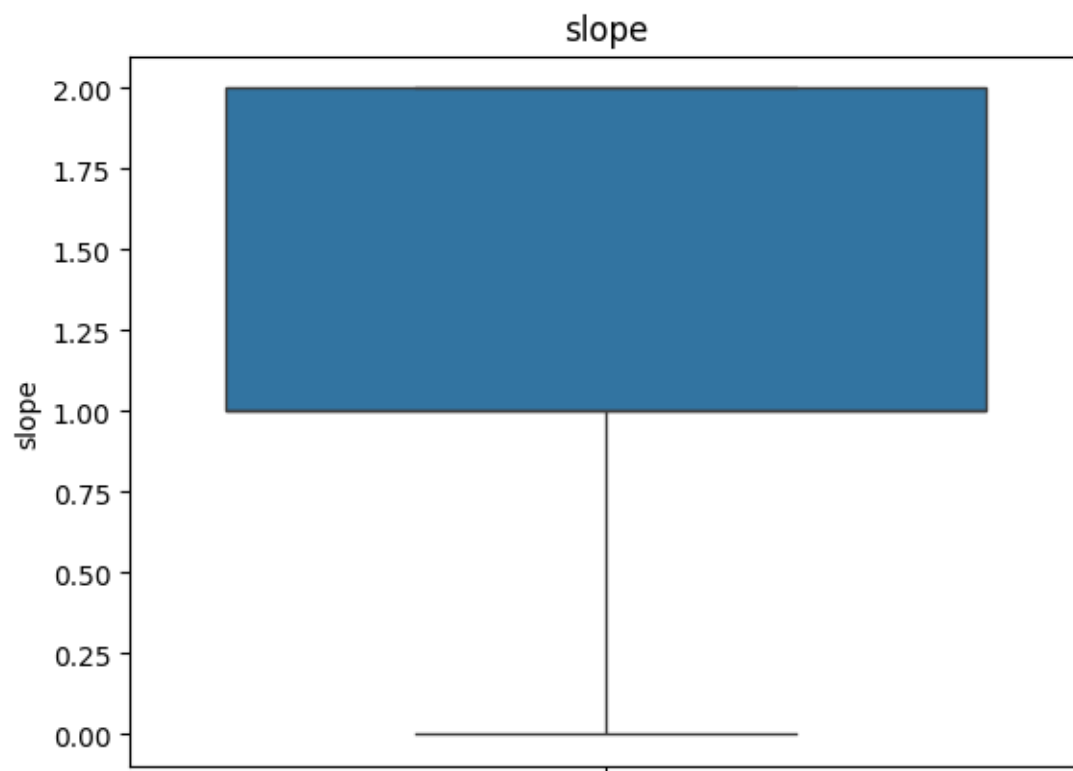
```
[29]: plt.figure(figsize=(10, 6)) # Adjust the figure size if needed

for col in col_name:
    sns.boxplot(data=df[col])
    plt.title(col)
    plt.show()
```









```
[30]: df = df.dropna()
```

```
[31]: df.isna().sum()
```

```
[31]: age          0
      sex          0
      cp           0
      trestbps     0
      chol         0
      fbs          0
      restecg      0
      thalach      0
      exang        0
      oldpeak      0
      slope        0
      ca           0
      thal         0
      target       0
      dtype: int64
```

```
[32]: df = df.drop('fbs',axis=1)
```

```
[34]: correlations = df.corr()['target'].drop('target')

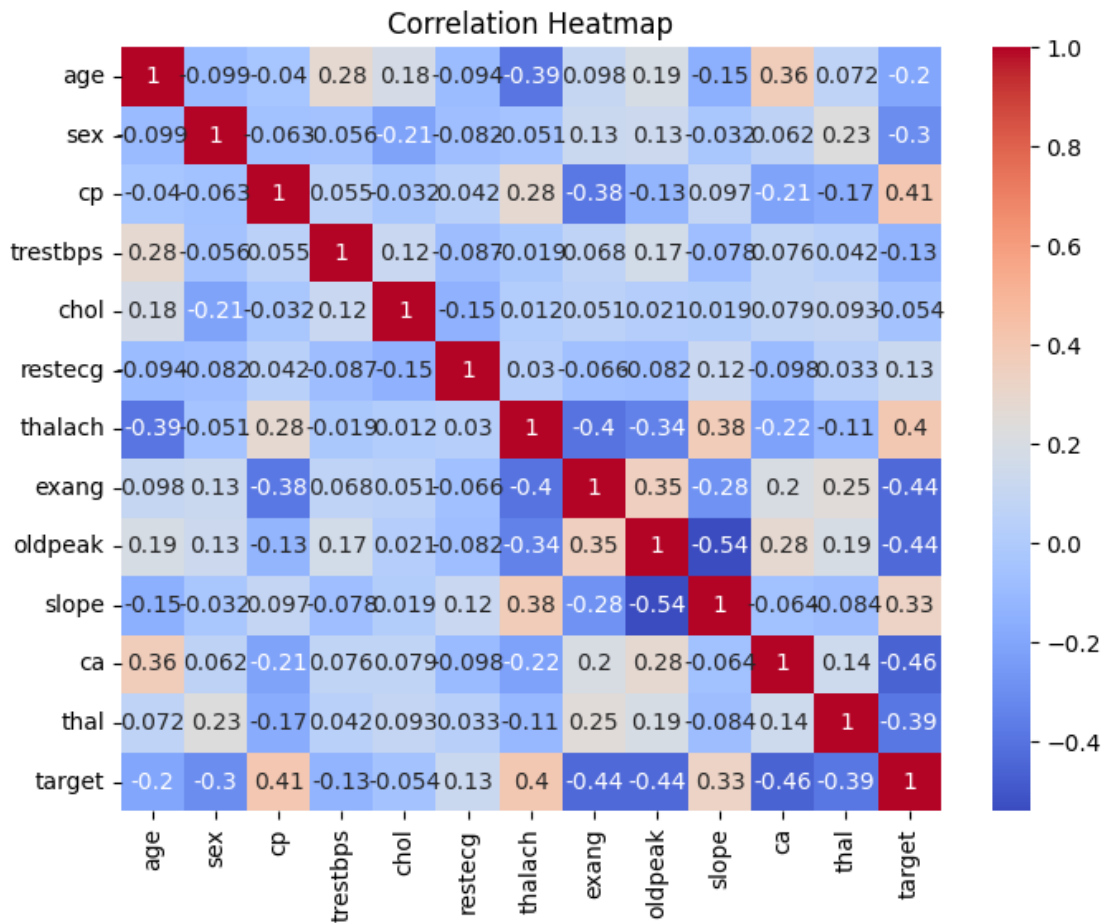
      # Print correlations
      print("Correlation with the Target:")
      print(correlations)
      print()

      # Plot correlation heatmap
      plt.figure(figsize=(8, 6))
      sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
      plt.title('Correlation Heatmap')
      plt.show()
```

Correlation with the Target:

```
age          -0.199970
sex          -0.300311
cp            0.408985
trestbps     -0.132882
chol         -0.053834
restecg       0.125710
thalach       0.398870
exang        -0.435511
oldpeak      -0.436247
```

```
slope      0.327420
ca         -0.459629
thal      -0.389514
Name: target, dtype: float64
```



```
[37]: x = df[['cp', 'thal', 'exang', 'oldpeak', 'slope', 'ca']]
      y = df.target
      x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=0)

      x_train.shape, x_test.shape, y_train.shape, y_test.shape
```

```
[37]: ((219, 6), (55, 6), (219,), (55,))
```

```
[38]: from sklearn.preprocessing import StandardScaler
```

```
[39]: scaler = StandardScaler()
```

```
[40]: x_train_scaled = scaler.fit_transform(x_train)
      x_test_scaled = scaler.transform(x_test)
```

```
[42]: import numpy as np
```

```
[43]: y_train= np.array(y_train).reshape(-1, 1)
      y_test= np.array(y_test).reshape(-1, 1)
```

```
[44]: y_train.shape
```

```
[44]: (219, 1)
```

```
[45]: model = LogisticRegression()
      model.fit(x_train_scaled, y_train)

      # Make predictions on the test set
      y_pred = model.predict(x_test_scaled)

      # Evaluate the model's accuracy
      accuracy = accuracy_score(y_test, y_pred)
      print("Accuracy:", accuracy)
```

Accuracy: 0.8181818181818182

/home/bcl07/.local/lib/python3.8/site-packages/sklearn/utils/validation.py:1183:
DataConversionWarning: A column-vector y was passed when a 1d array was
expected. Please change the shape of y to (n_samples,), for example using
ravel().

```
y = column_or_1d(y, warn=True)
```

```
[46]: #Classification model using Decision Tree
      from sklearn.tree import DecisionTreeClassifier
      tc=DecisionTreeClassifier(criterion='entropy')
      tc.fit(x_train_scaled,y_train)
      y_pred=tc.predict(x_test_scaled)

      print("Training Accuracy Score :",accuracy_score(y_pred,y_test))
      print("Training Confusion Matrix :",confusion_matrix(y_pred,y_test))
```

Training Accuracy Score : 0.7454545454545455

Training Confusion Matrix : [[21 5]
[9 20]]

```
[ ]:
```