**1. Preliminary analysis:**

**a) Perform preliminary data inspection and report the findings on the structure of the data, missing values, duplicates, etc.**

`df.head()`

| | age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | target |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 63 | 1 | 3 | 145 | 233 | 1 | 0 | 150 | 0 | 2.3 | 0 | 0 | 1 | 1 |
| 1 | 37 | 1 | 2 | 130 | 250 | 0 | 1 | 187 | 0 | 3.5 | 0 | 0 | 2 | 1 |
| 2 | 41 | 0 | 1 | 130 | 204 | 0 | 0 | 172 | 0 | 1.4 | 2 | 0 | 2 | 1 |
| 3 | 56 | 1 | 1 | 120 | 236 | 0 | 1 | 178 | 0 | 0.8 | 2 | 0 | 2 | 1 |
| 4 | 57 | 0 | 0 | 120 | 354 | 0 | 1 | 163 | 1 | 0.6 | 2 | 0 | 2 | 1 |

`df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   age       303 non-null    int64
 1   sex       303 non-null    int64
 2   cp        303 non-null    int64
 3   trestbps  303 non-null    int64
 4   chol      303 non-null    int64
 5   fbs       303 non-null    int64
 6   restecg   303 non-null    int64
 7   thalach   303 non-null    int64
 8   exang     303 non-null    int64
 9   oldpeak   303 non-null    float64
 10  slope     303 non-null    int64
 11  ca        303 non-null    int64
 12  thal      303 non-null    int64
 13  target    303 non-null    int64
dtypes: float64(1), int64(13)
memory usage: 33.3 KB
```

`df.shape`

```
(303, 14)
```

```
# Missing values
df.isnull().sum()
```

```
age           0
sex           0
cp            0
trestbps      0
chol         71
fbs           0
restecg       0
thalach       0
exang         0
oldpeak       0
slope         0
ca            0
thal          0
dtype: int64
```

: `df.duplicated().sum()`

: 1

b) **Based on these findings, remove duplicates (if any) and treat missing values using an appropriate strategy**

```
df2=df[df.duplicated()]
df2
```

|     | age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | target |
|-----|-----|-----|----|----------|------|-----|---------|---------|-------|---------|-------|----|------|--------|
| 164 | 38  | 1   | 2  | 138      | 175  | 0   | 1       | 173     | 0     | 0.0     | 2     | 4  | 2    | 1      |

```
df=df.drop_duplicates()
df.duplicated().sum()
```
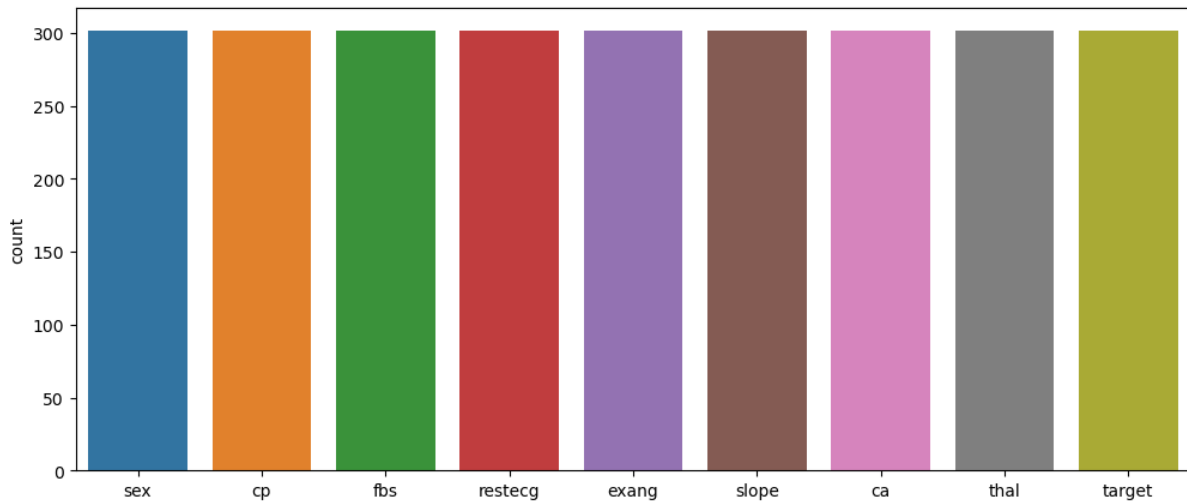
0

```
# removed duplicate value
```

2. **Prepare a report about the data explaining the distribution of the disease and the related factors using the steps listed below:**

   a. **Get a preliminary statistical summary of the data and explore the measures of central tendencies and spread of the data**

```
df.describe().T
```

|         | count | mean       | std       | min   | 25%   | 50%   | 75%   | max   |
|---------|-------|------------|-----------|-------|-------|-------|-------|-------|
| age     | 303.0 | 54.366337  | 9.082101  | 29.0  | 47.5  | 55.0  | 61.0  | 77.0  |
| sex     | 303.0 | 0.683168   | 0.466011  | 0.0   | 0.0   | 1.0   | 1.0   | 1.0   |
| cp      | 303.0 | 0.966997   | 1.032052  | 0.0   | 0.0   | 1.0   | 2.0   | 3.0   |
| trestbps| 303.0 | 131.623762 | 17.538143 | 94.0  | 120.0 | 130.0 | 140.0 | 200.0 |
| chol    | 303.0 | 246.264026 | 51.830751 | 126.0 | 211.0 | 240.0 | 274.5 | 564.0 |
| fbs     | 303.0 | 0.148515   | 0.356198  | 0.0   | 0.0   | 0.0   | 0.0   | 1.0   |
| restecg | 303.0 | 0.528053   | 0.525860  | 0.0   | 0.0   | 1.0   | 1.0   | 2.0   |
| thalach | 303.0 | 149.646865 | 22.905161 | 71.0  | 133.5 | 153.0 | 166.0 | 202.0 |
| exang   | 303.0 | 0.326733   | 0.469794  | 0.0   | 0.0   | 0.0   | 1.0   | 1.0   |
| oldpeak | 303.0 | 1.039604   | 1.161075  | 0.0   | 0.0   | 0.8   | 1.6   | 6.2   |
| slope   | 303.0 | 1.399340   | 0.616226  | 0.0   | 1.0   | 1.0   | 2.0   | 2.0   |
| ca      | 303.0 | 0.729373   | 1.022606  | 0.0   | 0.0   | 0.0   | 1.0   | 4.0   |
| thal    | 303.0 | 2.313531   | 0.612277  | 0.0   | 2.0   | 2.0   | 3.0   | 3.0   |
| target  | 303.0 | 0.544554   | 0.498835  | 0.0   | 0.0   | 1.0   | 1.0   | 1.0   |

**b. Identify the data variables which are categorical and describe and explore these variables using the appropriate tools, such as count plot**



**c. Study the occurrence of CVD across the Age category**

```
age_df=df[['age', 'target']]
age_df.groupby(['target']).mean()
```

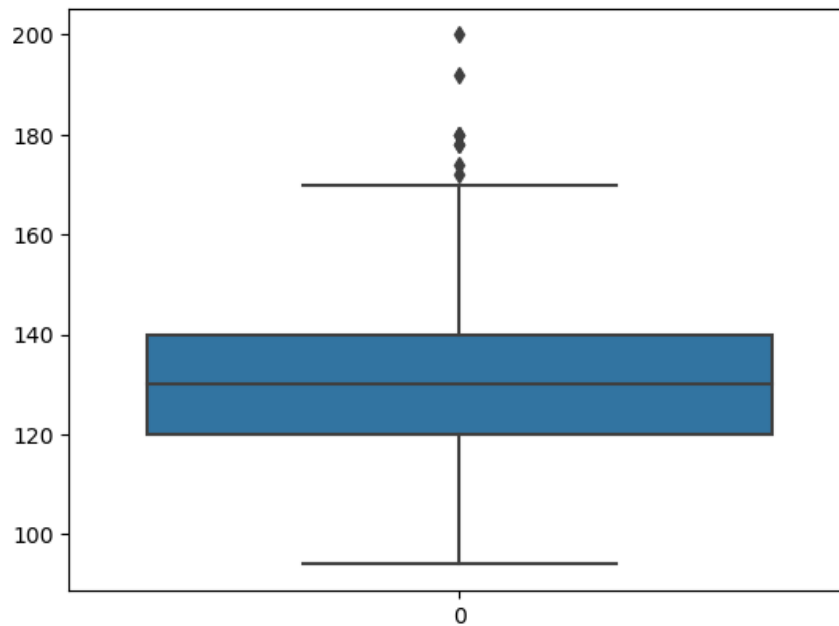|        | age       |
|--------|-----------|
| target |           |
| 0      | 56.601449 |
| 1      | 52.496970 |

### d. Study the composition of all patients with respect to the Sex category

```
df.groupby(['sex']).mean()
```

|     | age       | cp       | trestbps   | chol       | fbs     | restecg  | thalach    | exang    | oldpeak  | slope    | ca       | thal     | target   |
|-----|-----------|----------|------------|------------|---------|----------|------------|----------|----------|----------|----------|----------|----------|
| sex |           |          |            |            |         |          |            |          |          |          |          |          |          |
| 0   | 55.677083 | 1.041667 | 133.083333 | 261.302083 | 0.12500 | 0.572917 | 151.125000 | 0.229167 | 0.876042 | 1.427083 | 0.552083 | 2.125000 | 0.750000 |
| 1   | 53.758454 | 0.932367 | 130.946860 | 239.289855 | 0.15942 | 0.507246 | 148.961353 | 0.371981 | 1.115459 | 1.386473 | 0.811594 | 2.400966 | 0.449275 |

```
#sex has more impact on target 0
```

### e. Study if one can detect heart attacks based on anomalies in the resting blood pressure (trestbps) of a patient.

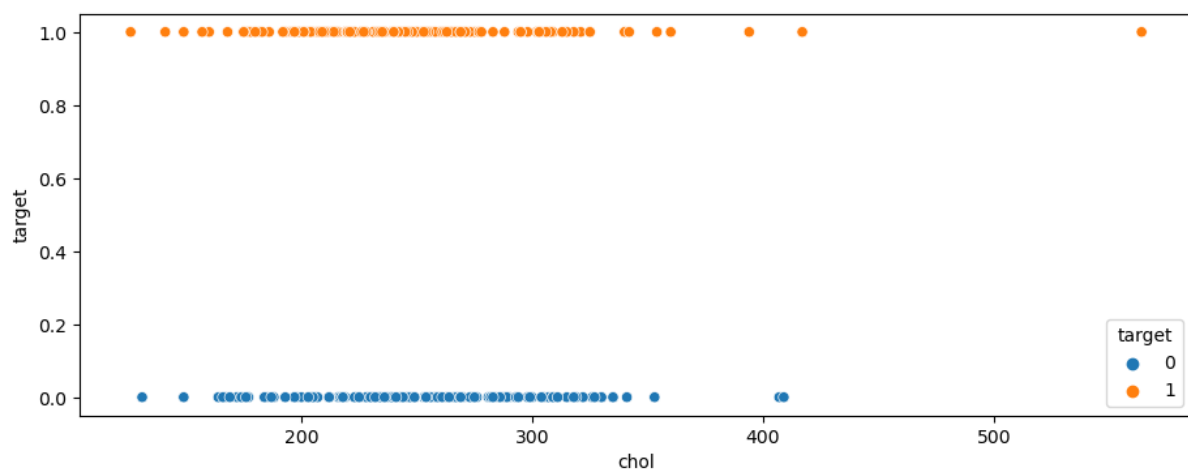```
df[df['trestbps']>170]['target'].value_counts()
```

```
0    6
1    3
Name: target, dtype: int64
```

```
df[df['trestbps']>180]['target'].value_counts()
```

```
0    2
Name: target, dtype: int64
```

```
#it has high impact on target 0
```

### f. Describe the relationship between cholesterol levels and a target variable

```
chol_df=df[['chol', 'target']]
chol_df.groupby(['target']).mean()
```

|  | chol |
|---|---|
| **target** | |
| **0** | 251.086957 |
| **1** | 242.230303 |

```
#chol has more impact on target 0
```

g. **State what relationship exists between peak exercising and the occurrence of a heart attack**

```
slope_df=df[['slope', 'target']]
slope_df.groupby(['target']).mean()
```
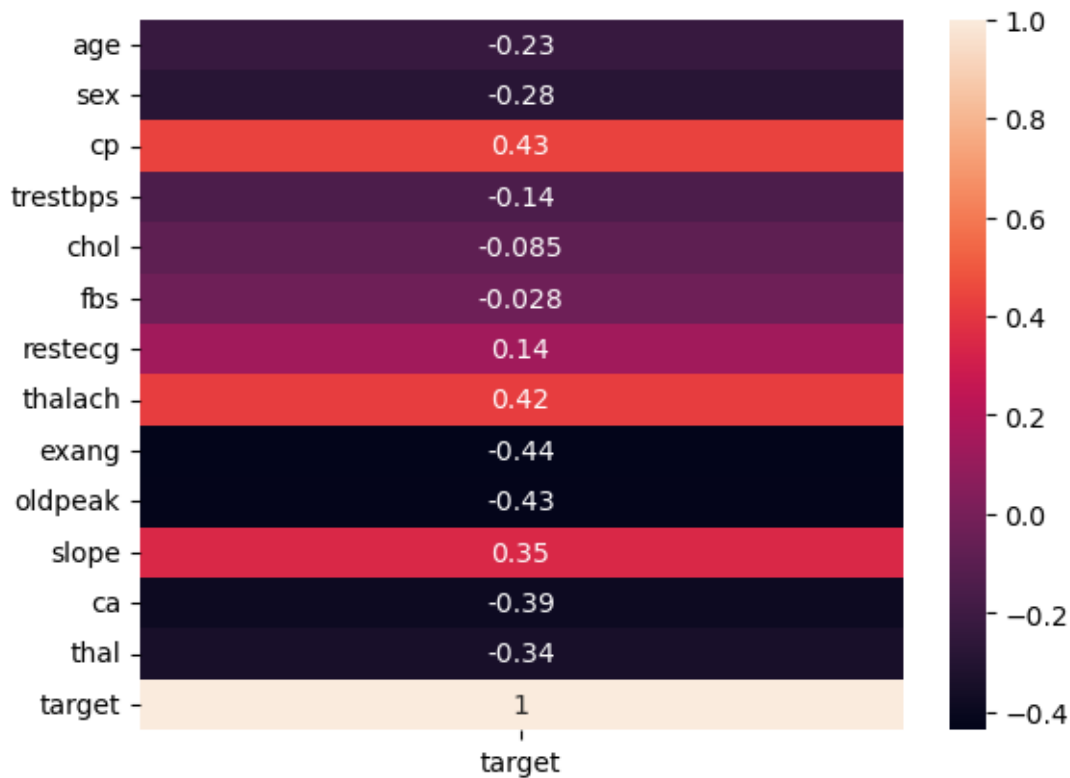
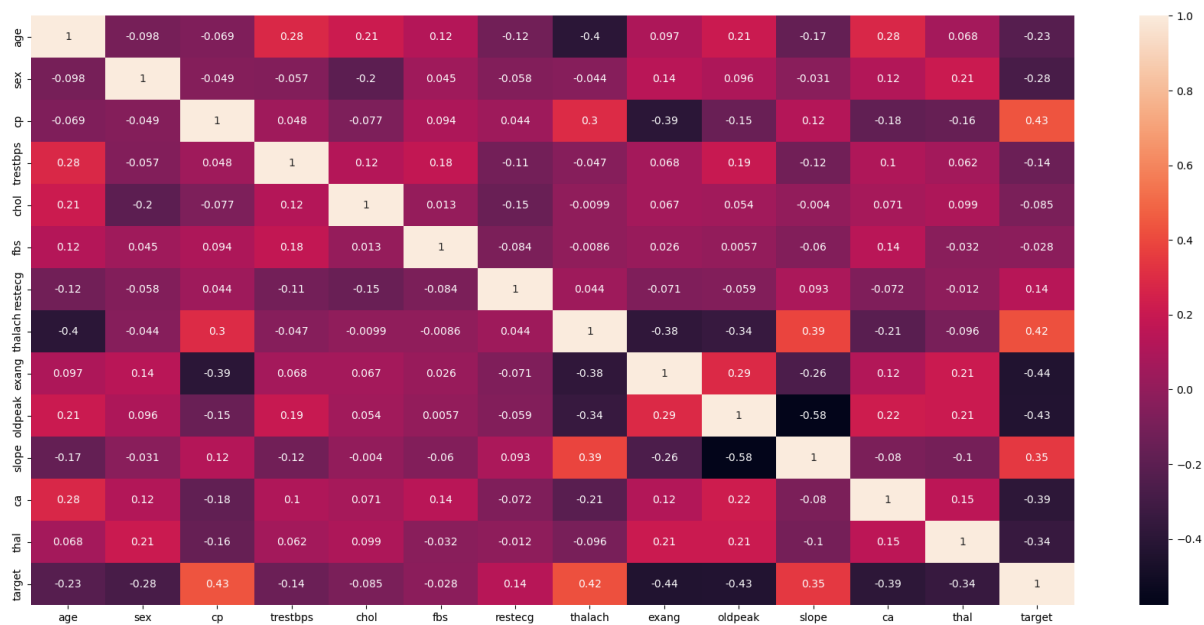|  | slope |
|---|---|
| **target** | |
| **0** | 1.166667 |
| **1** | 1.593939 |

```
plt.figure(figsize=(11,4))
sns.pairplot(data=df,x_vars='slope',y_vars='target',hue='target')
plt.show()
```

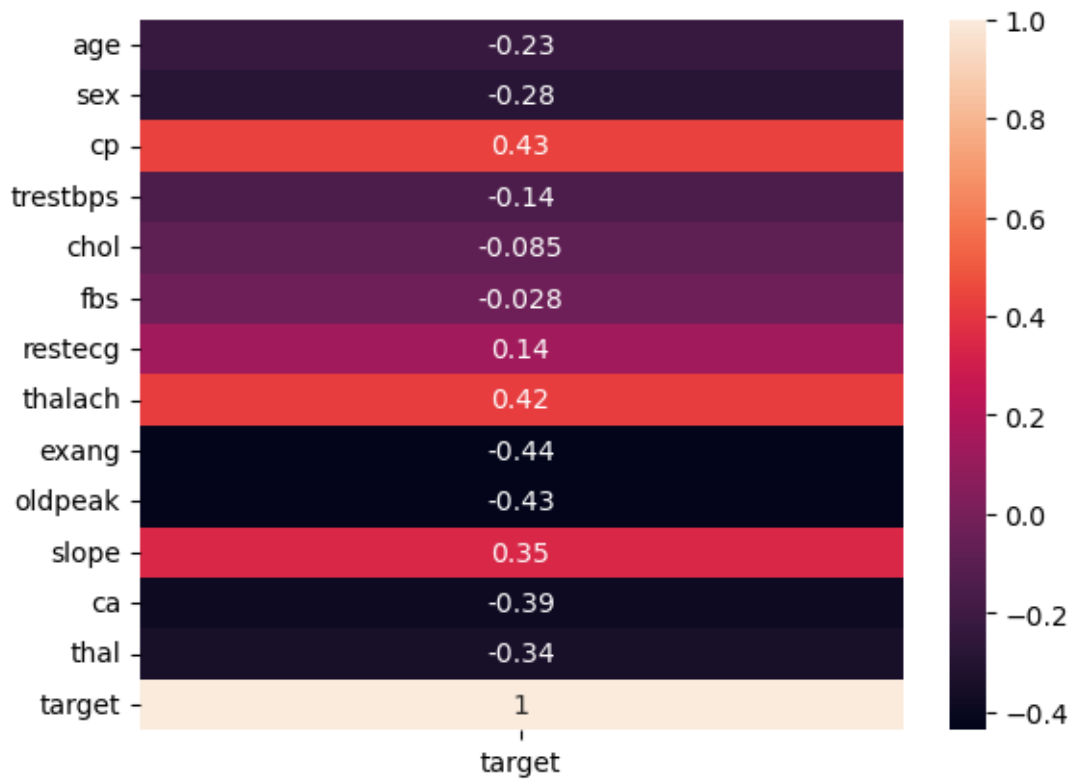<Figure size 1100x400 with 0 Axes>



```
# slope that us peak exercising has very less impact on both target 0 and 1
```

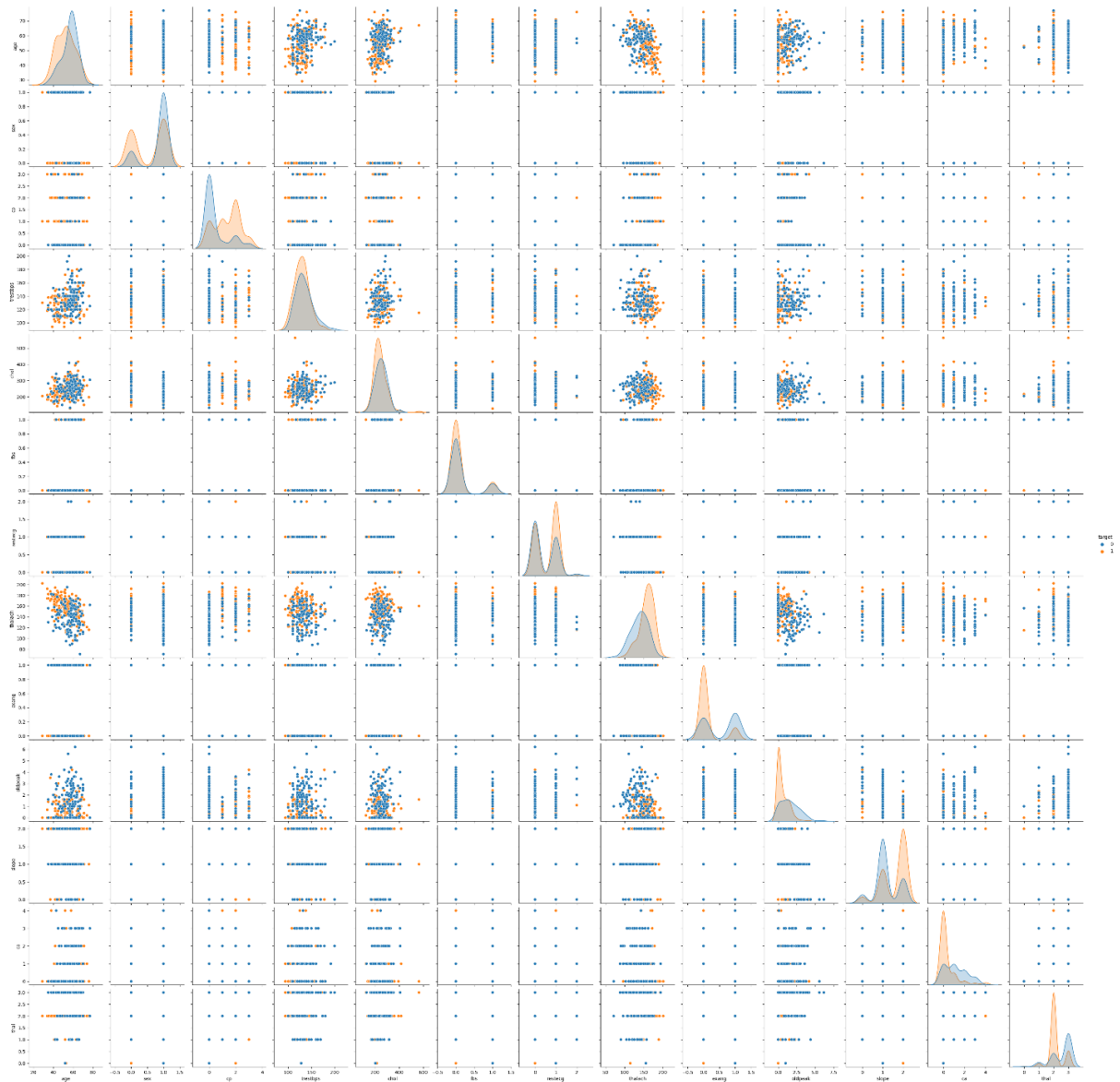h. **Check if thalassemia is a major cause of CVD**

**i.** **List how the other factors determine the occurrence of CVD**

**j. Use a pair plot to understand the relationship between all the given variables**

3. Build a baseline model to predict the risk of a heart attack using a logistic regression and random forest and explore the results while using correlation analysis and logistic regression (leveraging standard error and p-values from statsmodels) for feature selection

```python
# model building
```

```python
from sklearn.linear_model import LogisticRegression
```

```python
logreg=LogisticRegression()
```

```python
logreg.fit(X_train,y_train)
```

```
▾ LogisticRegression
LogisticRegression()
```

```python
pred=logreg.predict(X_test)
```

```python
logreg.score(X_train,y_train)
```

```
0.8613861386138614
```

```python
logreg.score(X_test,y_test)
```

```
0.82
```

```python
#testing
from sklearn.metrics import confusion_matrix, classification_report
confusion_matrix(y_test, pred)
```

```
array([[35,  8],
       [10, 47]], dtype=int64)
```

```python
print(classification_report(y_test, pred))

print(classification_report())
```

```
              precision    recall  f1-score   support

           0       0.77      0.79      0.78        42
           1       0.84      0.83      0.83        58

    accuracy                           0.81       100
   macro avg       0.80      0.81      0.81       100
weighted avg       0.81      0.81      0.81       100
```

```python
#random forest model
```

```python
from sklearn.ensemble import RandomForestClassifier
```

```python
clf_rf = RandomForestClassifier()
```

```python
clf_rf.fit(X_train, y_train)
```

▾ RandomForestClassifier
RandomForestClassifier()

```python
clf_rf.score(X_test, y_test)
```

0.81

```python
clf_rf.score(X_train, y_train)
```

0.9900990099009901

```python
predictions = clf_rf.predict(X_test)
```

```python
confusion_matrix(y_test, predictions)
```

```
array([[32, 10],
       [10, 48]], dtype=int64)
```

```python
print(classification_report(y_test, predictions))
```

```
              precision    recall  f1-score   support

           0       0.76      0.76      0.76        42
           1       0.83      0.83      0.83        58

    accuracy                           0.80       100
   macro avg       0.79      0.79      0.79       100
weighted avg       0.80      0.80      0.80       100
```