

In [1]:

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn import tree
from sklearn.metrics import accuracy_score, confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt
```

In [3]:

```
data = 'heart.csv'
df = pd.read_csv(data)
df.head()
```

Out[3]:

	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

In [4]:

```
df[df['chol']>300].shape
```

Out[4]:

(43, 14)

In [5]:

```
f=df[df['thal']==2]
```

In [6]:

```
f[f['target']==1].shape
```

Out[6]:

(130, 14)

In [7]:

```
df.shape
```

Out[7]:

(303, 14)

In [9]:

```
df.isna().sum()
```

Out[9]:

```
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
```

```
thalach      0
exang        0
oldpeak      0
slope        0
ca           0
thal         0
target       0
dtype: int64
```

In [10]:

```
df.dtypes
```

Out[10]:

```
age          int64
sex          int64
cp           int64
trestbps     int64
chol         int64
fbs          int64
restecg      int64
thalach      int64
exang        int64
oldpeak      float64
slope        int64
ca           int64
thal         int64
target       int64
dtype: object
```

In [ ]:

In [12]:

```
df.info()
```

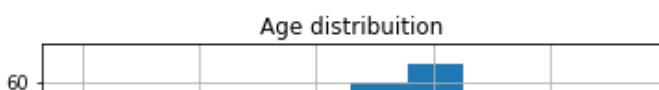
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
 #   Column        Non-Null Count  Dtype  
---  -
 0   age           303 non-null   int64  
 1   sex           303 non-null   int64  
 2   cp            303 non-null   int64  
 3   trestbps      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
```

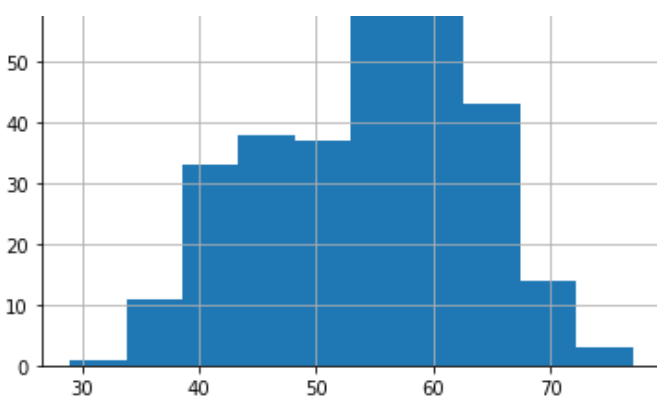
In [11]:

```
df['age'].hist(grid=True, bins=10);
plt.title('Age distribution')
```

Out[11]:

```
Text(0.5, 1.0, 'Age distribution')
```





In the above graph, we can analyse the distribution of Age column, and we can say that there are 60+ people who are having age between 57 to 63.

In [16]:

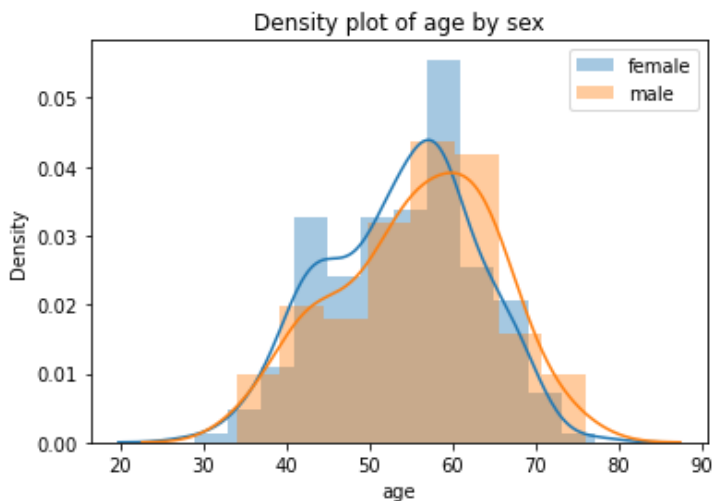
```
sns.distplot(df[df['sex']==1]['age'], label='female')
sns.distplot(df[df['sex']==0]['age'], label='male')
plt.legend()
plt.title('Density plot of age by sex')
plt.show()
```

C:\Users\lenovo\anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)

C:\Users\lenovo\anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

warnings.warn(msg, FutureWarning)



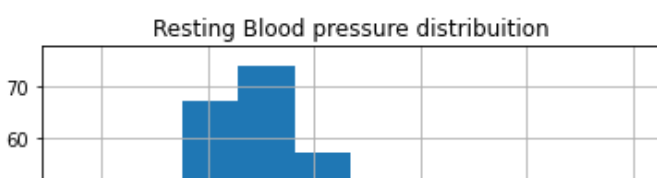
Density graph shows the smoothed distribution of points along the numerical axis. The density peaks where there is the highest concentration of points. In sum, density graphs can be considered smoothed histograms.

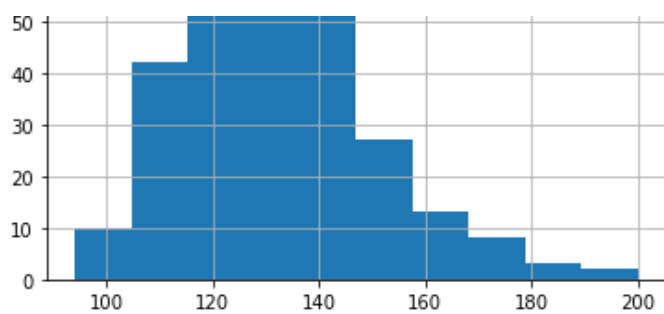
In [18]:

```
df['trestbps'].hist()
plt.title('Resting Blood pressure distribution')
```

Out[18]:

Text(0.5, 1.0, 'Resting Blood pressure distribution')

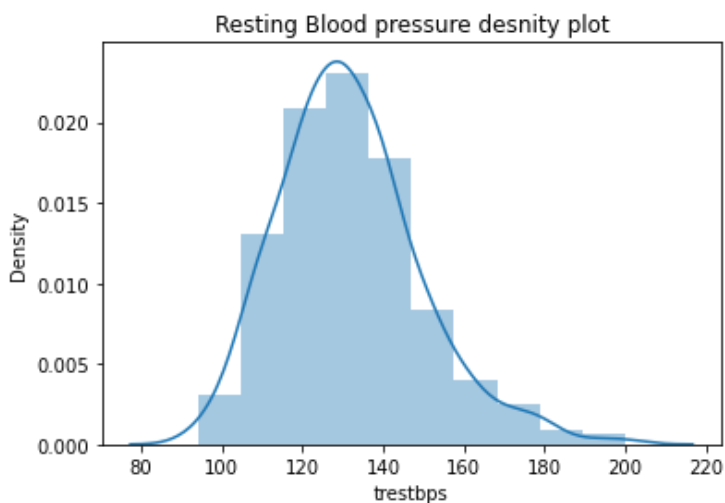




In [19]:

```
sns.distplot(df['trestbps'], bins=10)
plt.title('Resting Blood pressure desnity plot');
```

C:\Users\lenovo\anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)



**In the above graphh, we are having a normal distribution**

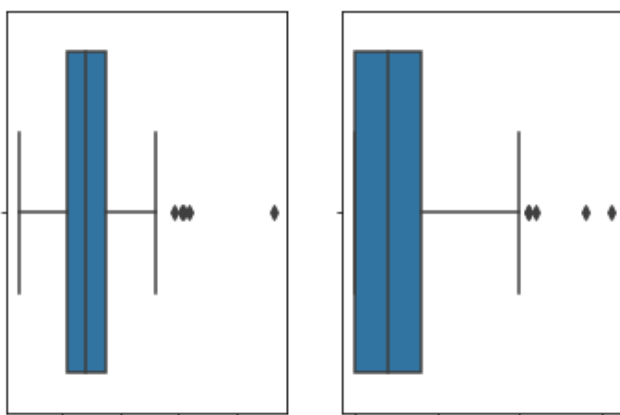
In [20]:

```
fig, axes = plt.subplots(nrows = 1, ncols=2)
sns.boxplot(x='chol', data=df, orient='v', ax=axes[0])
sns.boxplot(x='oldpeak', data=df, orient='v', ax=axes[1])
```

C:\Users\lenovo\anaconda3\lib\site-packages\seaborn\\_core.py:1319: UserWarning: Vertical orientation ignored with only `x` specified.  
warnings.warn(single\_var\_warning.format("Vertical", "x"))  
C:\Users\lenovo\anaconda3\lib\site-packages\seaborn\\_core.py:1319: UserWarning: Vertical orientation ignored with only `x` specified.  
warnings.warn(single\_var\_warning.format("Vertical", "x"))

Out[20]:

<AxesSubplot:xlabel='oldpeak'>

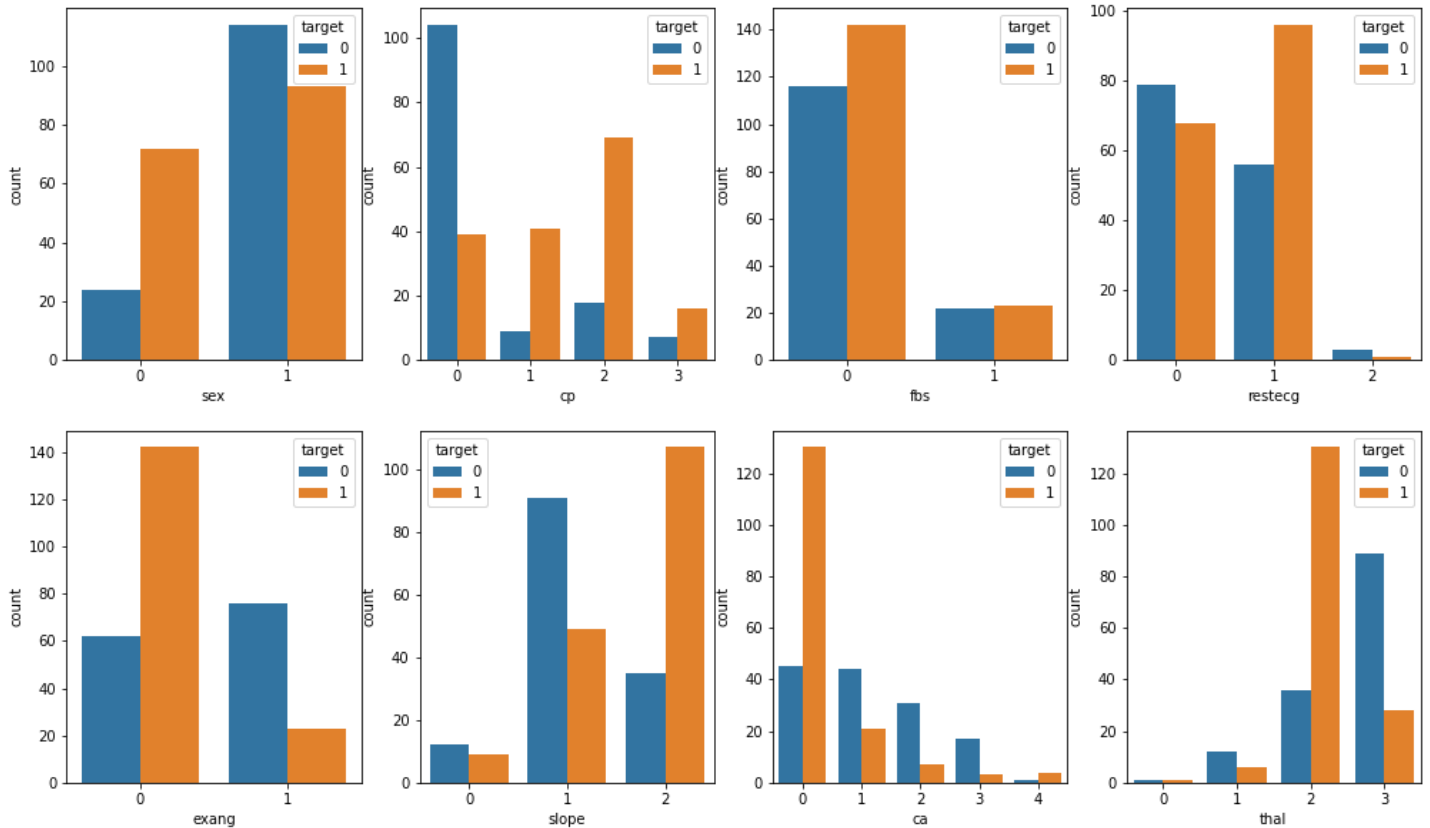


200 300 400 500 0 2 4 6  
chol oldpeak

In [26]:

```
fig, axes = plt.subplots(nrows=2, ncols=4, figsize=(17,10))
cat_feat = ['sex', 'cp', 'fbs', 'restecg', 'exang', 'slope', 'ca', 'thal', 'target']

for idx, feature in enumerate(cat_feat):
    if feature != 'target':
        ax = axes[int(idx/4), idx%4]
        sns.countplot(x=feature, hue='target', data=df, ax=ax)
```



Let's get some insights from this chart:

**Chest pain:** the heart disease diagnosis is greater among the patients that feel any chest pain.

**Restecg - Eletrocardiograph results:** the rate of heart disease diagnoses higher for patients with a ST-T wave abnormality .

**Slope:** The ratio of patients diagnosed with heart disease is higher for slope = 2

**Ca:** The diagnosed ratio decreases for ca between 1 and 3.

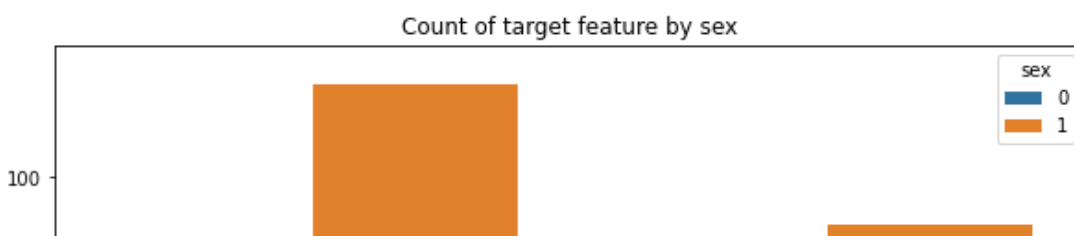
**Thal:** the diagnosed ratio is higher for thal = 2.

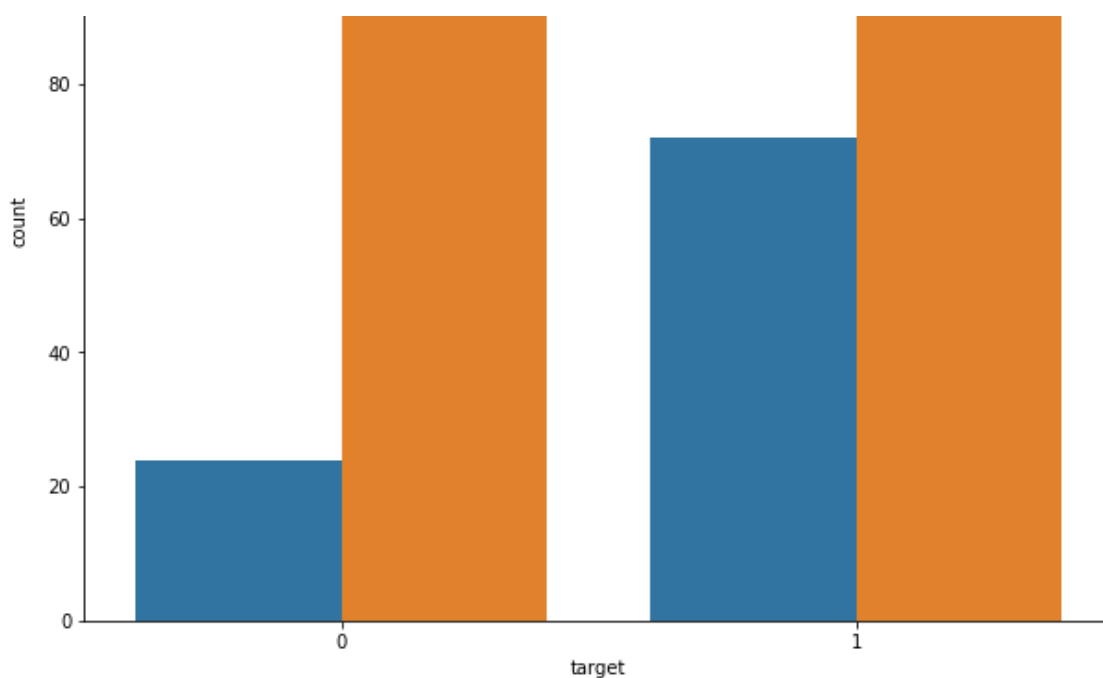
In [27]:

```
plt.rcParams['figure.figsize'] = (10,8)
sns.countplot(x='target', hue='sex', data=df);
plt.title('Count of target feature by sex')
```

Out[27]:

Text(0.5, 1.0, 'Count of target feature by sex')





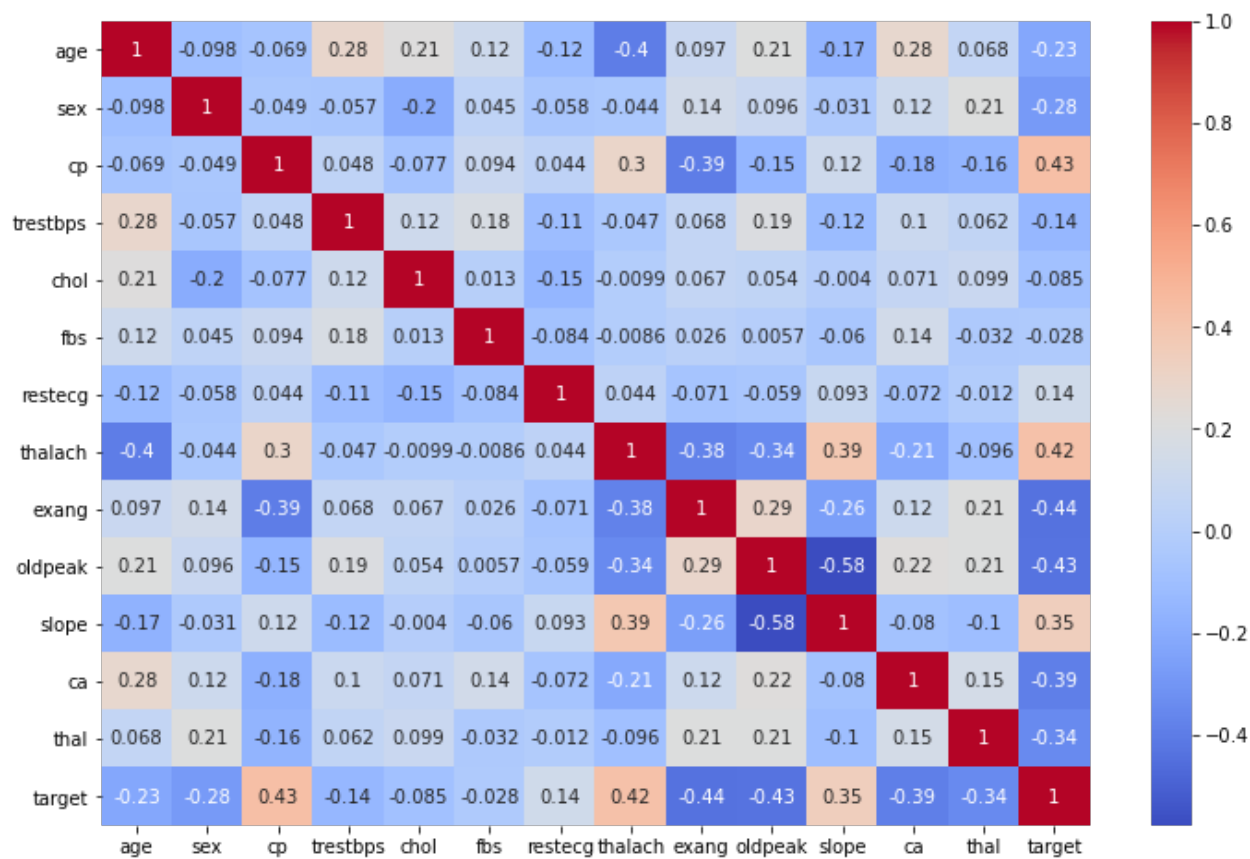
The amount of healthy male people is greater than the amount of unhealthy. For women, the number of unhealthy women is higher.

In [29]:

```
plt.figure(figsize=(12,8))
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
```

Out[29]:

<AxesSubplot:>



Apparently there are no features with a pretty strong correlation (above |0.7|)

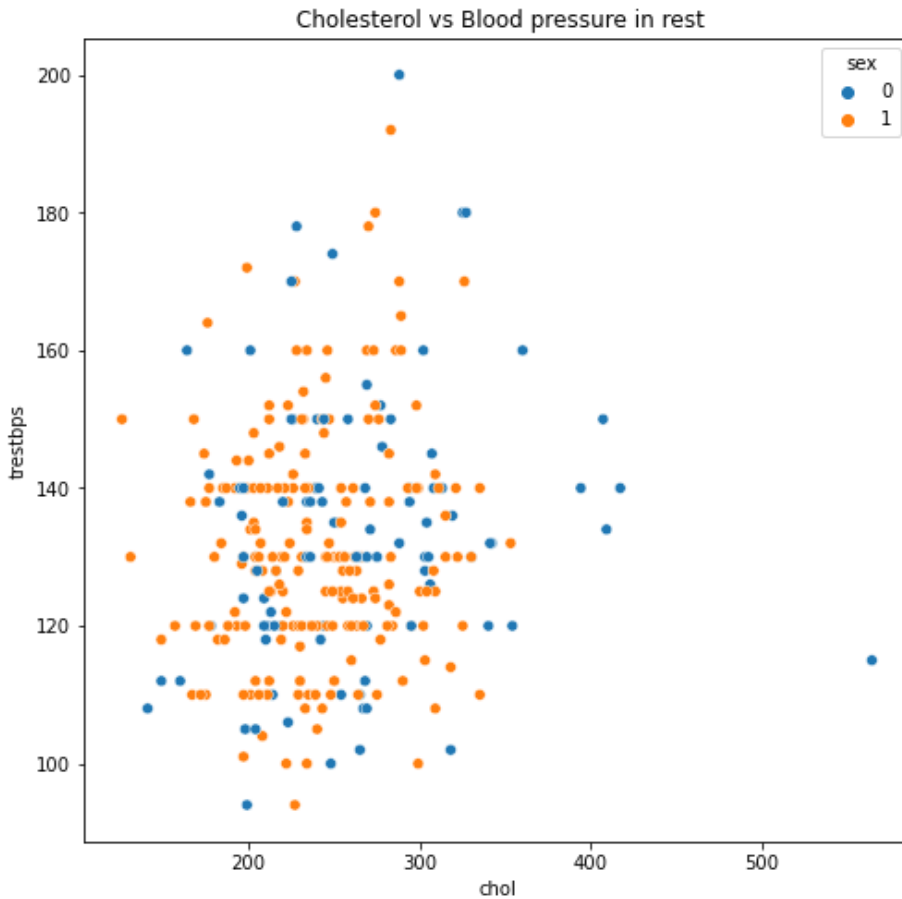
In [30]:

```
plt.rcParams['figure.figsize'] = (8,8)
sns.scatterplot(x='chol', y='trestbps', hue='sex', size=None, data=df)
```

```
plt.title(' Cholesterol vs Blood pressure in rest')
```

Out[30]:

```
Text(0.5, 1.0, ' Cholesterol vs Blood pressure in rest')
```



As can be seen there is a patient with high cholesterol. But, there's not a specific division between those that feel pain during exercise practice and those of not feel pain. We can use hue to filter by sex. It's also possible to filter using size = 'label\_to\_filer'.

In [18]:

```
X = df.drop(columns=['target'])#independent variable
y = df['target']#dependnet or target value
print(X.shape)
print(y.shape)
```

```
(303, 13)
```

```
(303,)
```

In [20]:

```
x_train,x_test,y_train,y_test = train_test_split(X,y,random_state=0, test_size=0.3)
print(x_train.shape)
print(x_test.shape)
```

```
(212, 13)
```

```
(91, 13)
```

In [21]:

```
clf = tree.DecisionTreeClassifier()
clf.fit(x_train,y_train)
y_train_pred = clf.predict(x_train)
y_test_pred = clf.predict(x_test)
```

In [22]:

```
y_train_pred
```

Out[22]:

```
array([1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0, 0, 0,
       0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1,
       0, 1, 1, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0,
       0, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0,
       0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1,
       0, 1, 0, 0, 0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0,
       0, 0, 1, 1, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0,
       0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0,
       1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1, 0,
       1, 1, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0], dtype=int64)
```

In [23]:

```
y_test_pred
```

Out[23]:

```
array([0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 0,
       0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0,
       0, 0, 0, 1, 0, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1,
       1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1,
       0, 1, 0], dtype=int64)
```

In [28]:

```
a=90
print('The value of a is', a)
```

The value of a is 90

In [29]:

```
print(f'The value of a is {a}')
```

The value of a is 90

In [1]:

```
# helper function
def plot_confusionmatrix(y_train_pred, y_train, dom):
    print(f'{dom} Confusion matrix')
    cf = confusion_matrix(y_train_pred, y_train)
    sns.heatmap(cf, annot=True, cmap='Blues', fmt='g') #For g and G , the maximum number of
    significant digits
    plt.tight_layout()
    plt.show()
```

In [ ]:

In [31]:

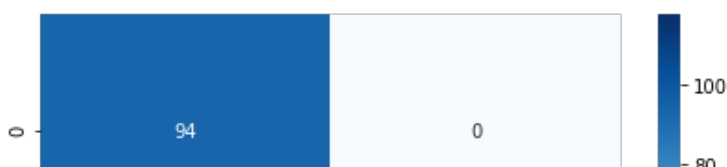
```
print(f'Train score {accuracy_score(y_train_pred, y_train)}')
print(f'Test score {accuracy_score(y_test_pred, y_test)}')
```

Train score 1.0  
Test score 0.7252747252747253

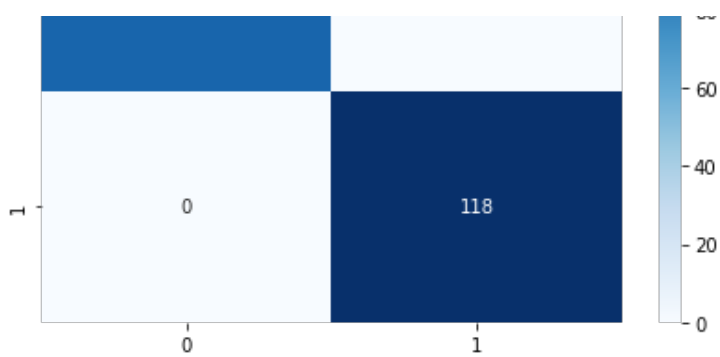
In [32]:

```
plot_confusionmatrix(y_train_pred, y_train, dom='Train')
plot_confusionmatrix(y_test_pred, y_test, dom='Test')
```

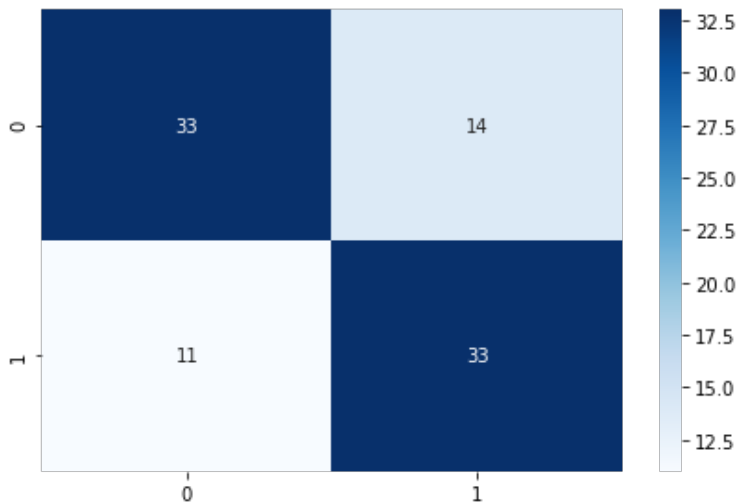
Train Confusion matrix







Test Confusion matrix



In [ ]:

```
#           Actual Values
#predicted    1    0
#           1  TP  FP
#           0  FN  TN
```

In [37]:

```
c_parameter_name = 'max_depth'
c_parameter_values = [1,2,3,4,5,6,7,8,9,10,11,12,13,14,15]
df = pd.DataFrame(columns=[c_parameter_name, 'accuracy'])

for input_parameter in c_parameter_values:
    model = tree.DecisionTreeClassifier(max_depth=input_parameter,splitter='best')
    model.fit(x_train, y_train)
    y_pred = model.predict(x_test)
    acc_score = accuracy_score(y_test,y_pred)*100
    df = df.append({c_parameter_name : input_parameter , 'accuracy' : acc_score}, ignore_index=True)
print(df)
print("")
```

	max_depth	accuracy
0	1.0	76.923077
1	2.0	73.626374
2	3.0	81.318681
3	4.0	76.923077
4	5.0	73.626374
5	6.0	71.428571
6	7.0	74.725275
7	8.0	72.527473
8	9.0	74.725275
9	10.0	71.428571
10	11.0	71.428571
11	12.0	71.428571
12	13.0	74.725275
13	14.0	71.428571
14	15.0	72.527473

In [ ]:

```
for input_para in c_para_values:
    model = DecisionTreeClassifier(max_depth=input_para, splitter='best')
    model.fit(X_train, y_train)
    y_pred1 = model.predict(X_test)
    acc_score = accuracy_score(y_test, y_pred)*100
    df = df.append({c_para_name: input_para, 'accuracy':acc_score}, ignore_index=True)
```

In [34]:

```
from sklearn.metrics import classification_report
```

In [35]:

```
print(classification_report(y_test_pred,y_test))
```

	precision	recall	f1-score	support
0	0.75	0.70	0.73	47
1	0.70	0.75	0.73	44
accuracy			0.73	91
macro avg	0.73	0.73	0.73	91
weighted avg	0.73	0.73	0.73	91

In [ ]:

```
#recall->tp / (tp + fn)
#The recall is the measure of our model correctly identifying True Positives.
#Thus, for all the customers who actually have heart disease, recall tells us how many we
correctly identified as a heart patient.

#precision of class 0 = TP of class 0/total number of object
#What is the Precision for our model? Yes, it is 0.843 or, when it predicts that a patient
has heart disease, it is correct around 84% of the time.
#precision of class 1 = TP of class 1/total number of object

#macro average = (precision of class 0 + precision of class 1)/2

#weighted average is precision of all classes merge together
#weighted average = (TP of class 0 + TP of class 1)/(total number of class 0 + total number
of class 1)

#F1-score is a measure of a model's accuracy on a dataset
#a good F1 score means that you have low false positives and low false negatives,
#Accuracy is used when the True Positives and True negatives are more important while
#F1-score is used when the False Negatives and False Positives are crucial.
#Support is the number of actual occurrences of the class in the specified dataset.
```

In [ ]: