

ML-MINOR-DEC

SUBMITTED BY : JOSWIN V JAISON

DATE OF EXERCISE GIVEN : 22/1/2021

DATE OF SUBMISSION : 22/1/2021

EMAIL:joswinvjaison@karunya.edu.in

```
import pandas as pd
df=pd.read_csv('heart.csv')
df.head()
```

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

submitted by : joswin v jaison email:[joswinvjaison@karunya.edu.in](mailto:joswinvjaison@karunya.edu.in)

Double-click (or enter) to edit

ML-MINOR-DEC

```
df.tail()
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
<b>298</b>	57	0	0	140	241	0	1	123	1	0.2	1	0	
<b>299</b>	45	1	3	110	264	0	1	132	0	1.2	1	0	
<b>300</b>	68	1	0	144	193	1	1	141	0	3.4	1	2	
<b>301</b>	57	1	0	130	131	0	1	115	1	1.2	1	1	
<b>302</b>	57	0	1	130	236	0	0	174	0	0.0	1	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.describe()
```

	age	sex	cp	trestbps	chol	fbs	restecg
<b>count</b>	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000
<b>mean</b>	54.366337	0.683168	0.966997	131.623762	246.264026	0.148515	0.528053
<b>std</b>	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525860
<b>min</b>	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000
<b>25%</b>	47.500000	0.000000	0.000000	120.000000	211.000000	0.000000	0.000000
<b>50%</b>	55.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000
<b>75%</b>	61.000000	1.000000	2.000000	140.000000	274.500000	0.000000	1.000000
<b>max</b>	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000

```
df.isnull().sum()
```

```

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

```

```

from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import chi2
data = df.copy()

```

```

x = data.iloc[:,0:13] #independent columns
y = data.iloc[:, -1]   #target column i.e price range
#apply SelectKBest class to extract top 10 best features
bestfeatures = SelectKBest(score_func=chi2, k=10)
fit = bestfeatures.fit(X,y)
dfscores = pd.DataFrame(fit.scores_)
dfcolumns = pd.DataFrame(X.columns)
#concat two dataframes for better visualization
featureScores = pd.concat([dfcolumns,dfscores],axis=1)
featureScores.columns = ['Specs','Score'] #naming the dataframe columns
print(featureScores.nlargest(12,'Score')) #print 10 best features

```

	Specs	Score
7	thalach	188.320472
9	oldpeak	72.644253
11	ca	66.440765
2	cp	62.598098
8	exang	38.914377
4	chol	23.936394
0	age	23.286624
3	trestbps	14.823925
10	slope	9.804095
1	sex	7.576835
12	thal	5.791853
6	restecg	2.978271

```

from sklearn.ensemble import ExtraTreesClassifier
import matplotlib.pyplot as plt

```

```

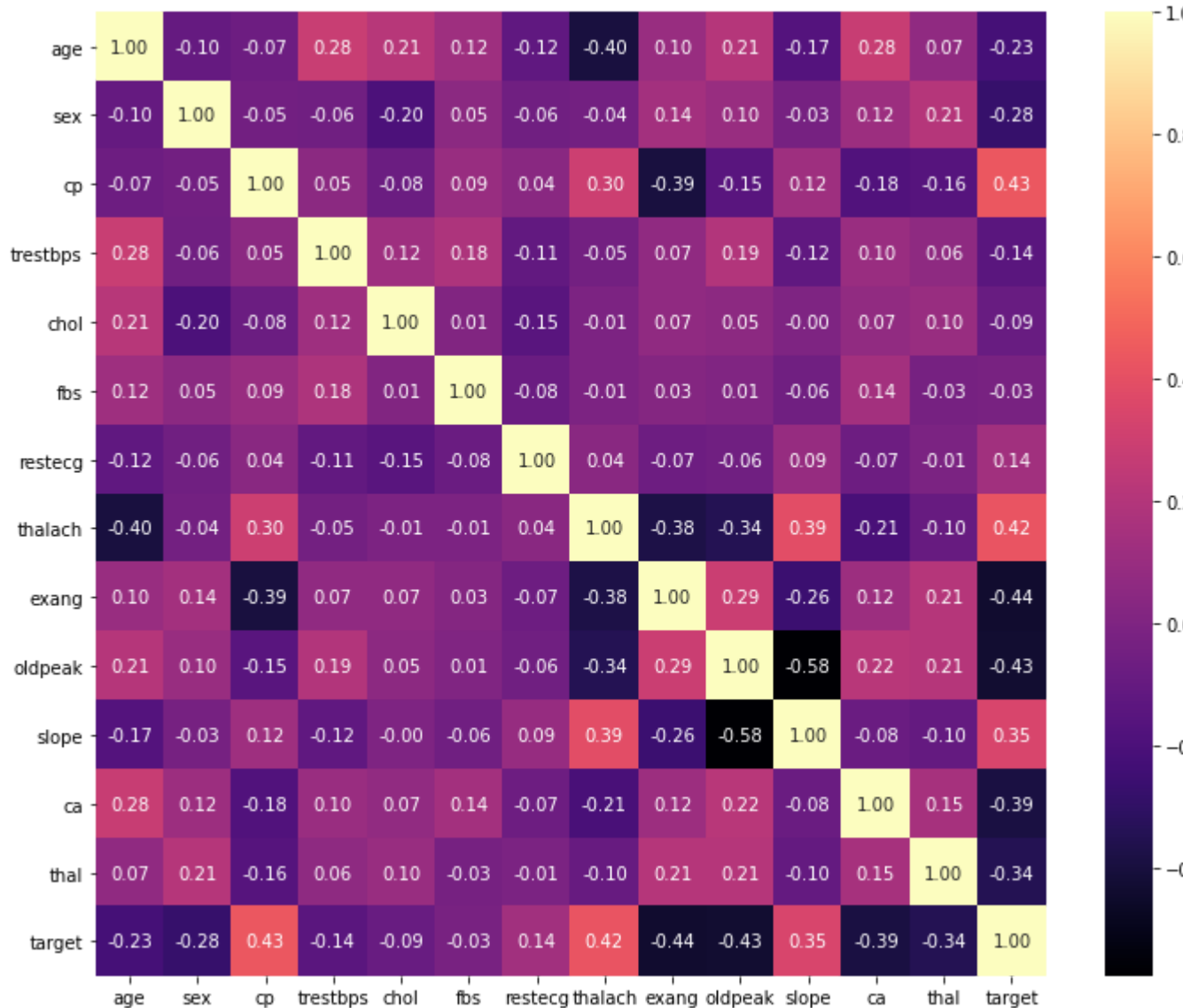
model = ExtraTreesClassifier()
model.fit(X,y)
print(model.feature_importances_) #use inbuilt class feature_importances of tree based classi
#plot graph of feature importances for better visualization
feat_importances = pd.Series(model.feature_importances_, index=X.columns)
feat_importances.nlargest(13).plot(kind='barh')
plt.show()

```

```
[0.06829859 0.05365263 0.12177194 0.06383964 0.05981136 0.01973602
```

```
import seaborn as sns
plt.figure(figsize=(12,10))
sns.heatmap(df.corr(),annot=True,cmap="magma",fmt='.2f')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f48ffcee160>
```



```
for i in df.columns:
    print(i,len(df[i].unique()))
```

```
age 41
sex 2
cp 4
trestbps 49
chol 152
fbs 2
restecg 3
thalach 91
exang 2
oldpeak 40
slope 3
```

```
ca 5
thal 4
target 2
```

```
sns.set_style('darkgrid')
sns.set_palette('Set2')
```

```
df2=df.copy()
df2.head()
```

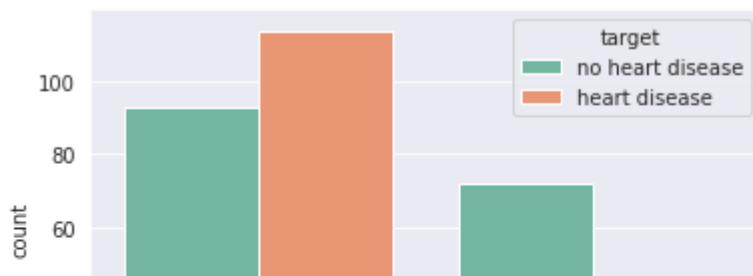
	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
<b>0</b>	63	1	3	145	233	1	0	150	0	2.3	0	0	1
<b>1</b>	37	1	2	130	250	0	1	187	0	3.5	0	0	2
<b>2</b>	41	0	1	130	204	0	0	172	0	1.4	2	0	2
<b>3</b>	56	1	1	120	236	0	1	178	0	0.8	2	0	2
<b>4</b>	57	0	0	120	354	0	1	163	1	0.6	2	0	2

```
def chng(sex):
    if sex == 0:
        return 'female'
    else:
        return 'male'
df2['sex'] = df2['sex'].apply(chng)
```

```
def chng2(prob):
    if prob==0:
        return 'heart disease'
    else:
        return 'no heart disease'
df2['target']=df2['target'].apply(chng2)
```

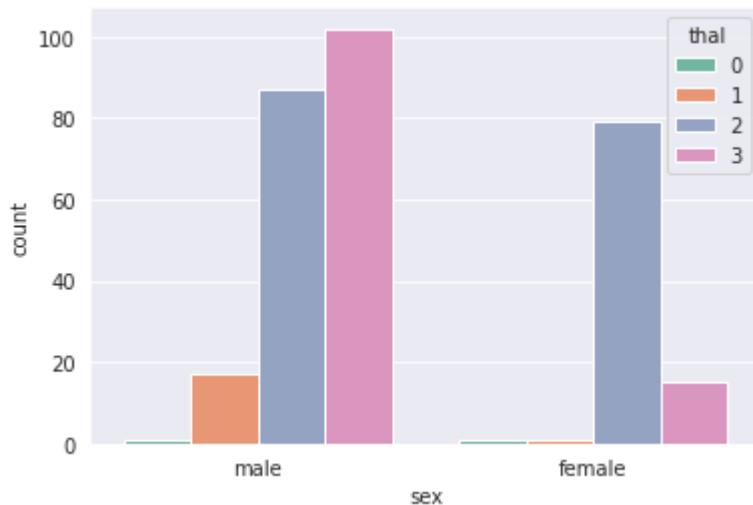
```
sns.countplot(data=df2,x='sex',hue='target')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f48f6ca7b70>
```



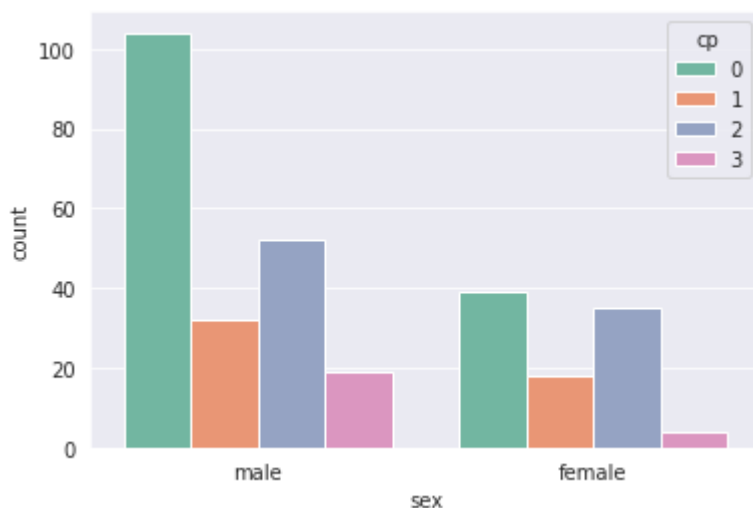
```
sns.countplot(data=df2,x='sex',hue='thal')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f48f516b2b0>
```



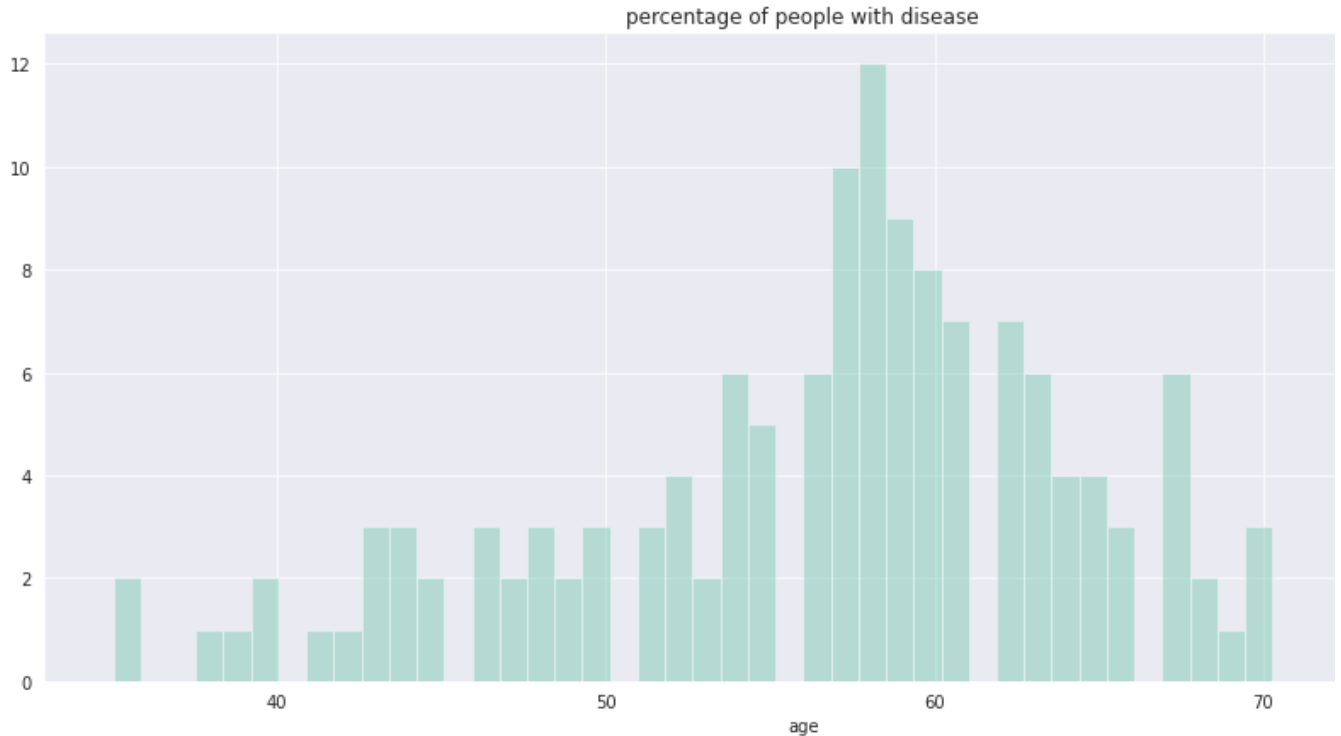
```
sns.countplot(data=df2,x='sex',hue='cp')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f48f6c68c88>
```



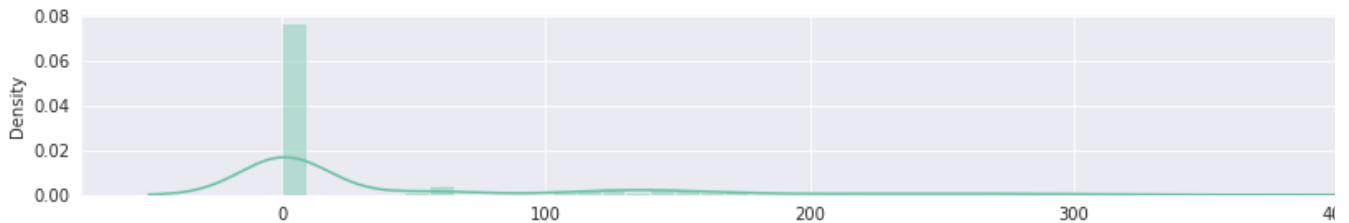
```
plt.figure(figsize=(16,7))
sns.distplot(df[df['target']==0]['age'],kde=False,bins=50)
plt.title("percentage of people with disease")
```

```
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2557: FutureWarning: `d
warnings.warn(msg, FutureWarning)
Text(0.5, 1.0, 'percentage of people with disease')
```



```
#cholesterol of heart diseased patients
plt.figure(figsize=(16,2))
sns.distplot(df[df['target']==0 ])
```

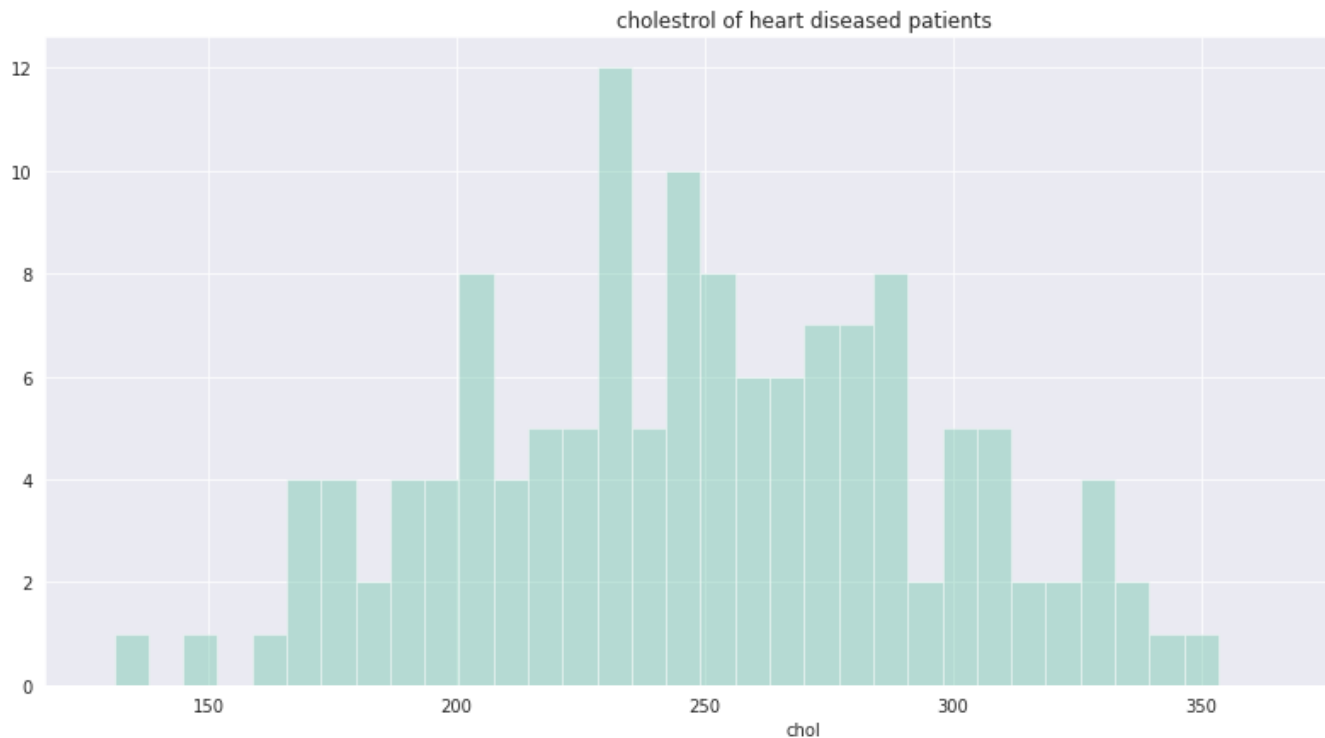
```
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2557: FutureWarning: `d
warnings.warn(msg, FutureWarning)
<matplotlib.axes._subplots.AxesSubplot at 0x7f48f4f6fb70>
```



```
plt.figure(figsize=(16,7))
sns.distplot(df[df['target']==0][ 'chol'],kde=False,bins=40)
plt.title('cholesterol of heart diseased patients')
```

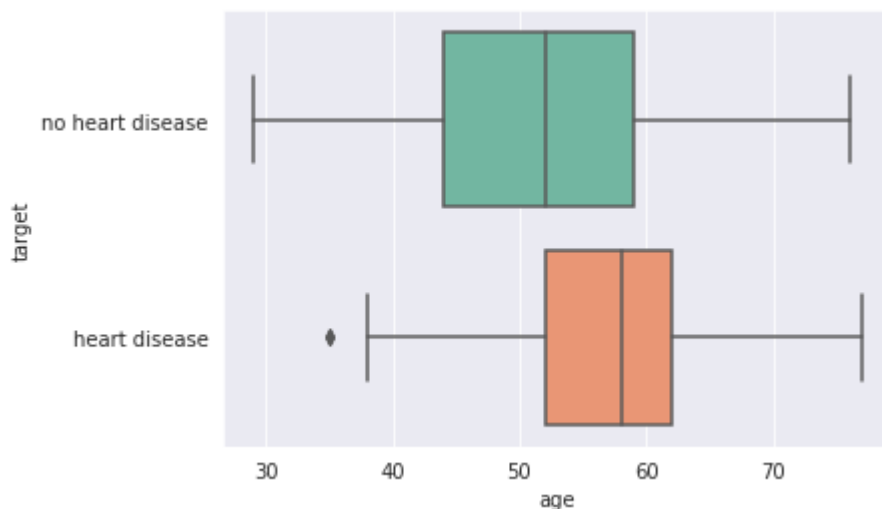


```
/usr/local/lib/python3.6/dist-packages/seaborn/distributions.py:2557: FutureWarning: `d
warnings.warn(msg, FutureWarning)
Text(0.5, 1.0, 'cholesterol of heart diseased patients')
```



```
#boxplot to determine heart disease with age
sns.boxplot(data=df2,x='age',y='target')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f48f4e70ac8>
```



## CLASSIFICATION MODEL BUILDING

```
x=df.drop('target',axis=1)
```

```
y=df['target']
```

```
x.head()
```

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

```
x.isnull().sum()
```

```
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
dtype: int64
```

```
y.isnull().sum()
```

```
0
```

```
y.head()
```

```
0    1
1    1
2    1
3    1
4    1
Name: target, dtype: int64
```

```
x.dtypes
```

```
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
dtype: object

```

```
x.shape
```

```
(303, 13)
```

```

from sklearn.model_selection import train_test_split
xtrain,xtest,ytrain,ytest=train_test_split(x,y,test_size=0.3,stratify=y,random_state=42)

```

```
x.describe()
```

	age	sex	cp	trestbps	chol	fbs	restecg
<b>count</b>	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000
<b>mean</b>	54.366337	0.683168	0.966997	131.623762	246.264026	0.148515	0.528053
<b>std</b>	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525860
<b>min</b>	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000
<b>25%</b>	47.500000	0.000000	0.000000	120.000000	211.000000	0.000000	0.000000
<b>50%</b>	55.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000
<b>75%</b>	61.000000	1.000000	2.000000	140.000000	274.500000	0.000000	1.000000
<b>max</b>	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000

```

from sklearn.preprocessing import MinMaxScaler
scaler=MinMaxScaler()
xtrain=scaler.fit_transform(xtrain)
xtest=scaler.transform(xtest)
from sklearn.neighbors import KNeighborsClassifier
model=KNeighborsClassifier(n_neighbors=8)

```

```
model.fit(xtrain,ytrain)
```

```

KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                    metric_params=None, n_jobs=None, n_neighbors=8, p=2,

```

```
weights='uniform')
```

```
model.score(xtest,ytest)
```

```
0.7912087912087912
```

```
import numpy as np
```

```
np.sqrt(303)
```

```
17.406895185529212
```

```
x.head()
```

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

```
data=pd.DataFrame({0:[64,25],1:[1,0],2:[3,2],3:[120,130],4:[200,300],5:[1,0],6:[0,1],7:[100,2
```

```
data=scaler.transform(data)
```

```
model.predict(data)
```

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

FROM THE CLASSIFICATION MODEL IT IS ASSUMED THAT A has no heart disease even though the age is 64, but most of the features are tend to be normal. B has the chance of getting heart disease because of high cholestrol and high heart rate.

Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.