

# Heart Disease Prediction

In this machine learning project, I will predict whether any person is suffering from heart disease

Data will be downloaded from Kaggle as indicated below;  
(<https://www.kaggle.com/onitf/heart-disease-uci>)

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from matplotlib import rcParams
from matplotlib.cm import rainbow
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: df = pd.read_csv('/Users/asumankabugo/Desktop/heart.csv')
```

```
In [3]: 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.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
```

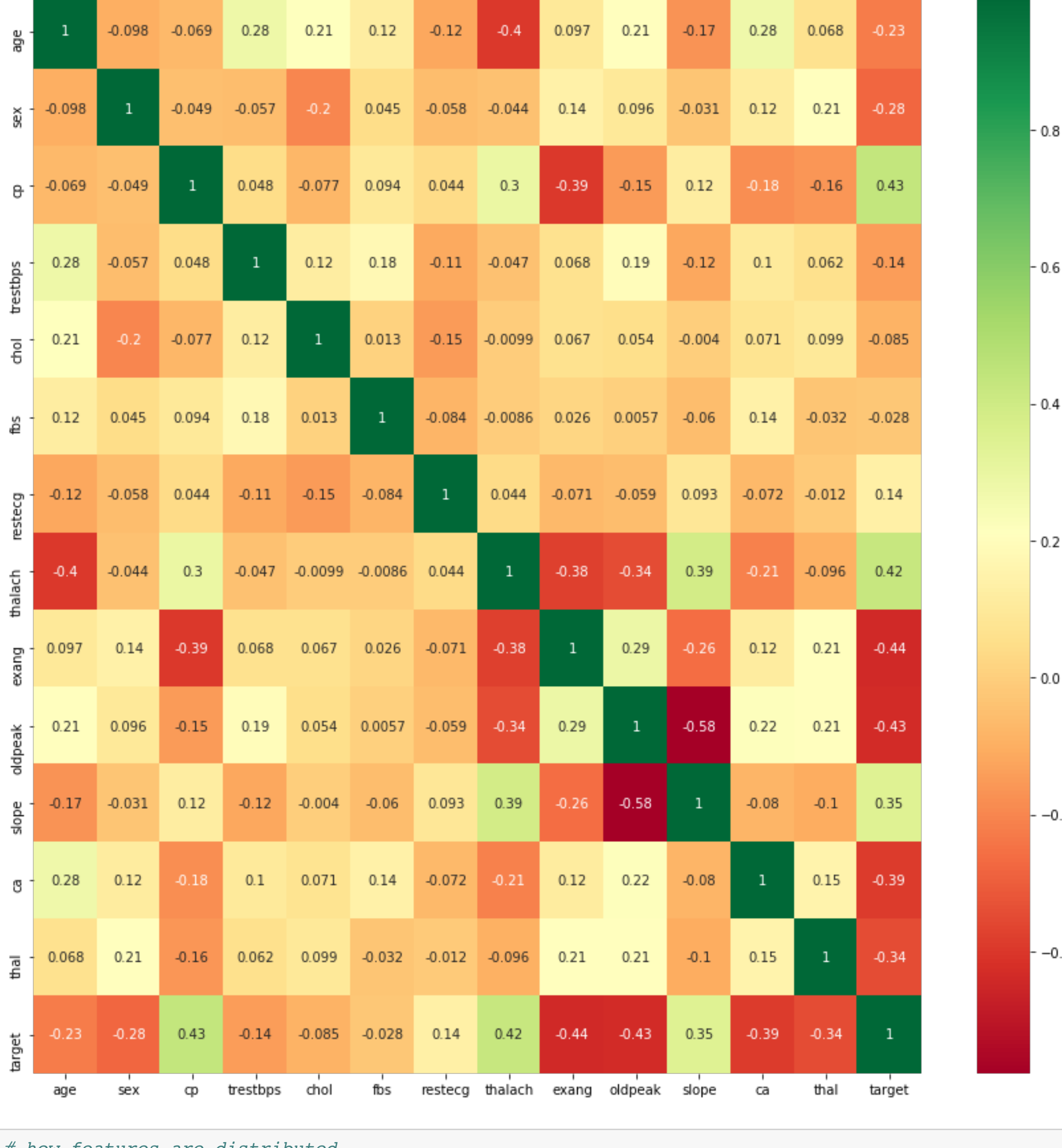
```
In [5]: df.describe()
```

```
Out[5]:
```

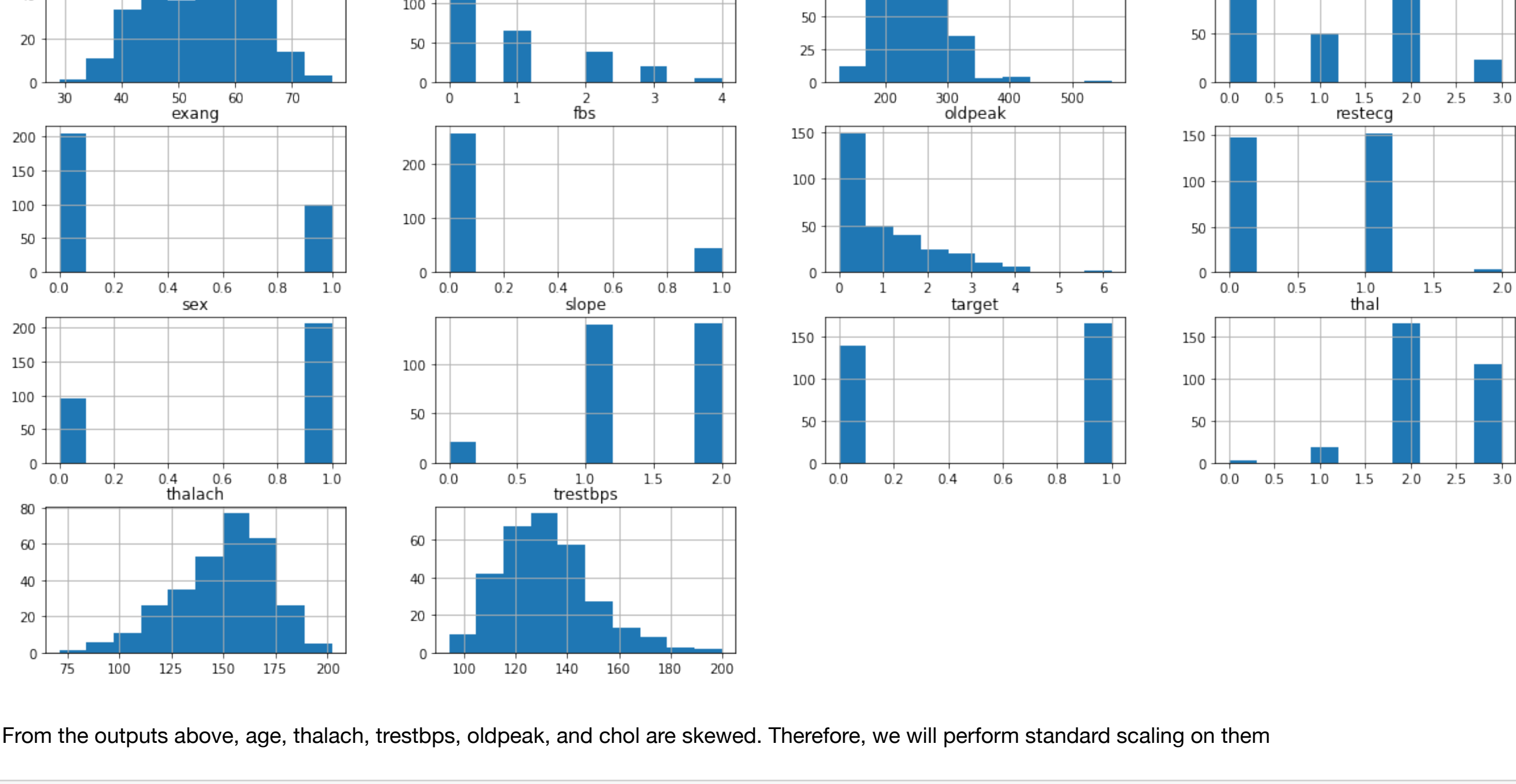
	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca		
count	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.00	303.00
mean	54.366337	0.683188	0.966997	131.623762	246.264026	0.148515	0.528053	149.646865	0.326733	1.039604	1.399340	0.729373	2.31	
std	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525860	22.905161	0.469794	1.161075	0.616226	1.022606	0.61	
min	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000	71.000000	0.000000	0.000000	0.000000	0.000000	0.00	
25%	47.500000	0.000000	0.000000	120.000000	211.000000	0.000000	0.000000	133.500000	0.000000	0.000000	1.000000	0.000000	2.00	
50%	55.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000	153.000000	0.000000	0.800000	1.000000	0.000000	2.00	
75%	61.000000	1.000000	2.000000	140.000000	274.500000	0.000000	1.000000	166.000000	1.000000	1.600000	2.000000	1.000000	3.00	
max	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000	202.000000	1.000000	6.200000	2.000000	4.000000	3.00	

## Feature Selection

```
In [6]: #get correlations of each features in dataset
corrmat = df.corr()
top_corr_features = corrmat.index
plt.figure(figsize=(15,15))
#plot heat map
g=sns.heatmap(df[top_corr_features].corr(),annot=True,cmap="RdYlGn")
```

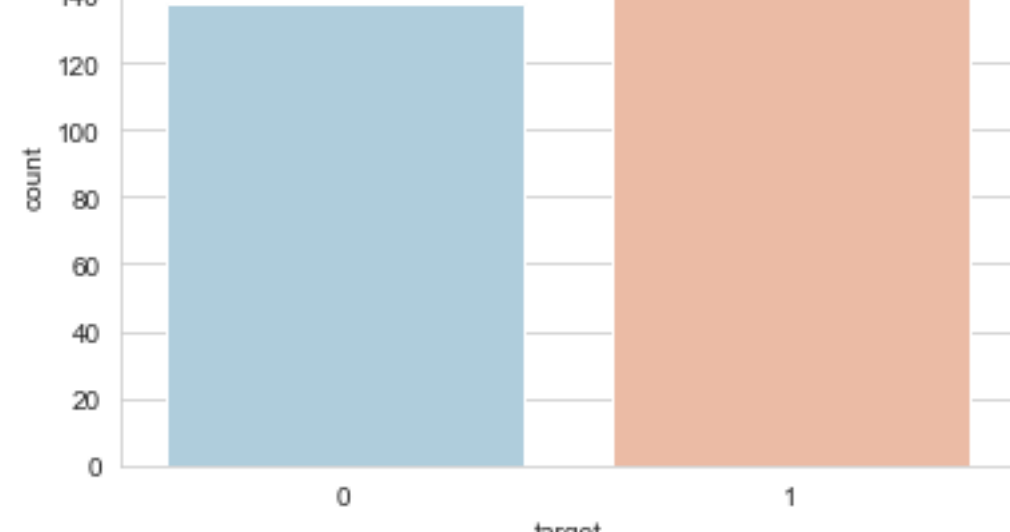


```
In [7]: # how features are distributed
df.hist(figsize=(20,10))
plt.show()
```



From the outputs above, age, thalach, trestbps, oldpeak, and chol are skewed. Therefore, we will perform standard scaling on them

```
In [8]: # understand whether the dataset is balanced or not using the target variable
sns.set_style('whitegrid')
sns.countplot(x='target',data=df,palette='RdBu_r')
plt.show()
```



The above output above indicate individuals with heart disease are approx 140 and those with heart disease approx 168. Therefore , this dataset is considered balanced.

## Data Processing

After exploring the dataset, I observed that I need to convert some categorical variables into dummy variables and scale all the values before training the Machine Learning models. First, I'll use the get\_dummies method to create dummy columns for categorical variables.

```
In [9]: dataset = pd.get_dummies(df, columns = ['sex', 'cp', 'fbs', 'restecg', 'exang', 'slope', 'ca', 'thal'])
```

## Here we will be experimenting with 4 algorithms

1. KNeighborsClassifier
2. DecisionTreeClassifier
3. RandomForestClassifier
4. Gaussian Naive Bayes

```
In [10]: from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.naive_bayes import GaussianNB
```

```
In [11]: from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
standardScaler = StandardScaler()
columns_to_scale = ['age', 'trestbps', 'chol', 'thalach', 'oldpeak']
dataset[columns_to_scale] = standardScaler.fit_transform(dataset[columns_to_scale])
```

```
In [12]: dataset.head()
```

```
Out[12]:
```

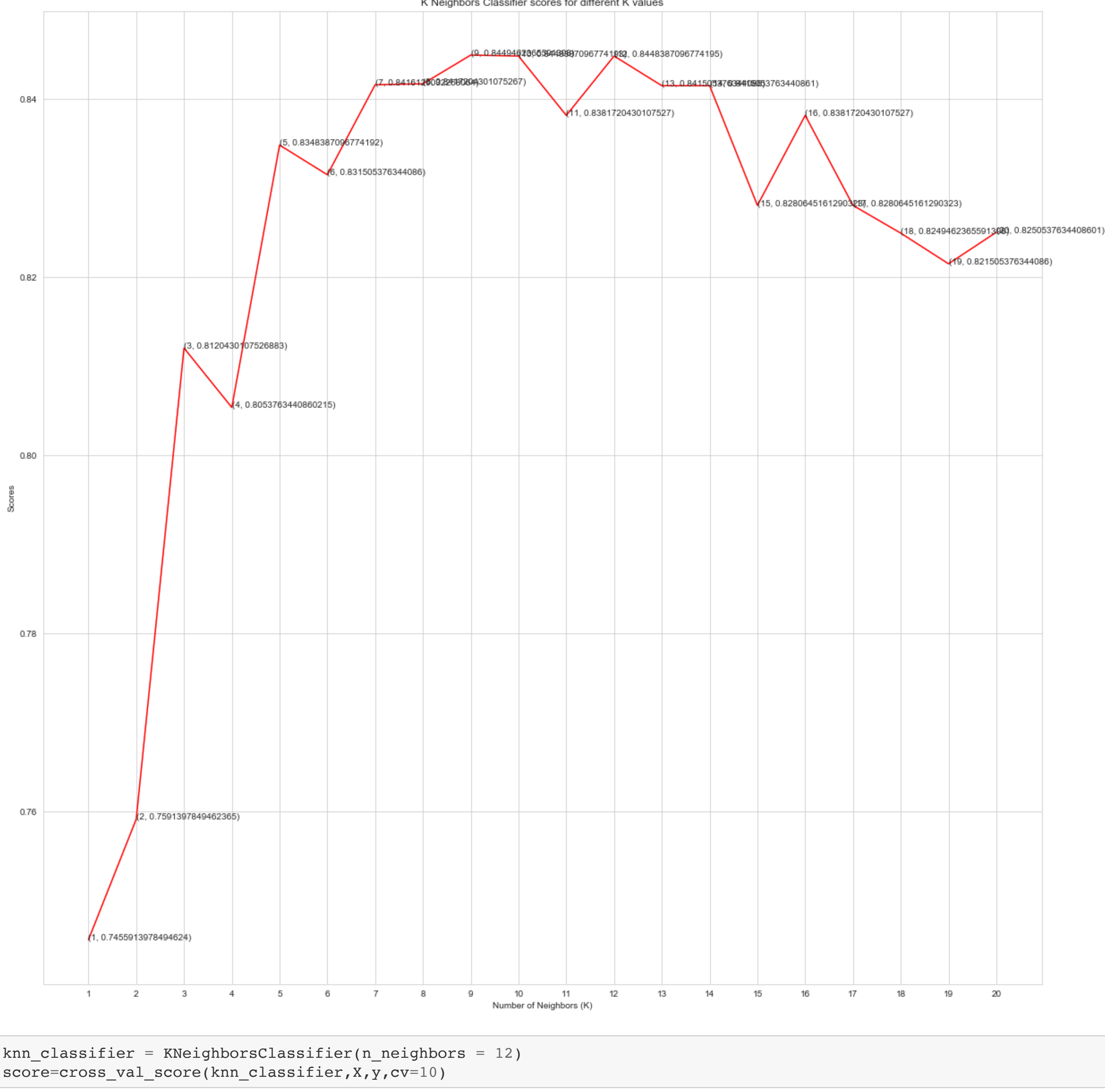
	age	trestbps	chol	thalach	oldpeak	target	sex_0	sex_1	cp_0	cp_1	...	slope_2	ca_0	ca_1	ca_2	ca_3	ca_4	thal_0	thal_1	thal_2	tha
0	0.952197	0.763956	-0.256334	0.015443	1.087338	1	0	1	0	0	...	0	1	0	0	0	0	0	1	0	
1	-1.915313	-0.092738	0.072199	1.833471	2.122573	1	0	1	0	0	...	0	1	0	0	0	0	0	0	1	
2	-1.474158	-0.092738	-0.816773	0.977514	0.310912	1	1	0	0	1	...	1	1	0	0	0	0	0	0	1	
3	0.180175	-0.663867	-0.198357	1.239897	-0.206705	0	0	1	0	1	...	1	1	0	0	0	0	0	0	1	
4	0.290464	-0.663867	2.082050	0.583939	-0.379244	1	1	0	1	0	...	1	1	0	0	0	0	0	0	1	

5 rows x 31 columns

```
In [13]: y = dataset['target']
X = dataset.drop(['target'], axis = 1)
```

```
In [14]: from sklearn.model_selection import cross_val_score
knn_scores = []
for k in range(1,21):
    knn_classifier = KNeighborsClassifier(n_neighbors = k)
    score=cross_val_score(knn_classifier,X,y,cv=10)
    knn_scores.append(score.mean())
```

```
In [15]: plt.figure(figsize=(20,20))
plt.plot([k for k in range(1, 21)], knn_scores,color = 'red')
for i in range(1,21):
    plt.text(i, knn_scores[i-1], (i, knn_scores[i-1]))
plt.xticks([i for i in range(1, 21)])
plt.xlabel('Number of Neighbors (K)')
plt.ylabel('Scores')
plt.title('K Neighbors Classifier scores for different K values')
plt.show()
```



```
In [16]: knn_classifier = KNeighborsClassifier(n_neighbors = 12)
score=cross_val_score(knn_classifier,X,y,cv=10)
```

```
In [17]: score.mean()
```

```
Out[17]: 0.8448387096774195
```

## Random Forest Classifier

```
In [18]: randomforest_classifier= RandomForestClassifier(n_estimators=10)
score=cross_val_score(randomforest_classifier,X,y,cv=10)
```

```
In [19]: score.mean()
```

```
Out[19]: 0.8186021505376344
```

## Gaussian Naive Bayes

```
In [20]: gaussian_NB = GaussianNB()
score=cross_val_score(gaussian_NB,X,y,cv=10)
```

```
In [21]: score.mean()
```

```
Out[21]: 0.8018279569892475
```

## Decision Tree Classifier

```
In [22]: decision_tree= DecisionTreeClassifier(random_state=42, max_depth=5)
score=cross_val_score(decision_tree,X,y,cv=10)
```

```
In [23]: score.mean()
```

```
Out[23]: 0.7450537634408603
```

```
In [ ]:
```