## Import libraries In [1]: import pandas as pd ##data processing, dataframe (read,drop,replace..etc) import numpy as np ##linear algebra import matplotlib.pyplot as plt ##..Data visualization %matplotlib inline import seaborn as sns ##Statistical Data Visualization #import warning import warnings warnings.filterwarnings('ignore') ##ignore unwanted warning In [2]: ## import the dataset hc data=pd.read csv('health care.csv') # Read the data In [3]: hc data trestbps chol fbs restecg thalach exang oldpeak slope thal Out[3]: sex ca target age 0 63 1 3 145 233 0 150 0 2.3 0 0 1 37 130 250 187 3.5 0 2 41 0 1 130 204 0 172 0 1.4 2 0 1 3 56 120 236 0 178 8.0 2 0 4 57 0 0 120 354 0 1 163 0.6 2 0 2 1 57 0 298 0 140 241 0 123 0.2 1 0 3 0 299 45 110 264 132 1.2 300 68 1 0 193 1 1 141 0 3.4 1 2 3 0 144 301 57 130 131 115 0 1.2 302 57 0 130 236 0 174 0.0 1 2 0 303 rows × 14 columns hc\_data1=hc\_data.copy() In [4]: #check shape of dataset In [5]: hc data.shape (303, 14) Out[5]: #check top 5 rows by using head In [6]: hc data.head() Out[6]: 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 1 130 250 0 187 3.5 0 0 2 2 41 0 1 130 204 0 172 1.4 2 0 1 3 56 1 1 120 236 0 178 8.0 0 1 57 120 354 163 0.6 hc data.tail() trestbps chol fbs restecg thalach exang oldpeak slope ca thal target Out[7]: age sex cp 3 0 298 57 0 0 140 241 0 123 0.2 1 0 0 299 110 0 132 264 1.2 193 300 68 141 301 131 115 302 57 0 130 236 0 0 174 0.0 0 ## CHeck null values In [8]: hc data.isnull().sum().any() False Out[8]: hc\_data.isnull().sum() 0 Out[9]: 0 0 ср 0 trestbps 0 chol fbs 0 restecg 0 thalach 0 exang 0 oldpeak 0 slope 0 0 са 0 thal 0 target dtype: int64 In [10]: #check info hc data.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 303 entries, 0 to 302 Data columns (total 14 columns): # Column Non-Null Count Dtype int64 303 non-null age 1 303 non-null sex int64 303 non-null int64 ср trestbps 303 non-null int64 chol 303 non-null int64 303 non-null int64 fbs restecg 303 non-null int64 7 thalach 303 non-null int64 303 non-null int64 exang oldpeak 303 non-null float64 10 slope 303 non-null int64 11 303 non-null int64 са 303 non-null 12 int64 thal 13 target 303 non-null int64 dtypes: float64(1), int64(13)memory usage: 33.3 KB In [11]: #check all description from dataset hc data.describe() Out[11]: trestbps chol fbs thalach oldpeak slope restecg exang age sex **count** 303.000000 303.000000 303.000000 303.000000 303.000000 303.000000 303.000000 303.000000 303.000000 303.000000 303 303.000000 0.966997 131.623762 246.264026 149.646865 54.366337 0.683168 0.148515 0.528053 0.326733 1.039604 1.399340 0 mean 9.082101 17.538143 51.830751 22.905161 0.466011 1.032052 0.356198 0.525860 0.469794 1.161075 0.616226 1 std 29.000000 0.000000 0.000000 94.000000 126.000000 0.000000 0.000000 71.000000 0.000000 0.000000 0.000000 0 min 120.000000 211.000000 0.000000 1.000000 0 25% 47.500000 0.000000 0.000000 0.000000 0.000000 133.500000 0.000000 50% 55.000000 130.000000 240.000000 1.000000 153.000000 0.000000 0.800000 0 1.000000 1.000000 0.000000 1.000000 61.000000 140.000000 1.000000 166.000000 2.000000 **75**% 1.000000 2.000000 274.500000 0.000000 1.000000 1.600000 1 2.000000 202.000000 3.000000 200.000000 564.000000 6.200000 2.000000 77.000000 1.000000 1.000000 1.000000 max # Histogram of the Heart Dataset fig = plt.figure(figsize = (40,30)) hc data.hist(ax = fig.gca()); From the above histogram plots, we can see that the features are skewed and not normally distributed. Also, the scales are different between one and another. **Understanding the Data** we can see the correlation between different features in below heat mat # Creating a correlation heatmap In [13]: sns.heatmap(hc data.corr(),annot=True, cmap='terrain', linewidths=0.1) fig=plt.gcf() fig.set size inches (20,20) plt.show() -0.098 -0.069 0.28 0.21 0.12 0.097 0.21 0.28 0.068 -0.098 -0.049 -0.057 0.045 -0.058 -0.044 -0.031 0.14 0.096 0.12 0.21 -0.077 0.094 -0.16 -0.069 trestbps -0.057 0.12 0.18 -0.11 -0.047 0.068 0.6 -0.077 0.12 1 0.013 -0.15 -0.0099 0.067 0.054 -0.004 0.071 -0.085 0.045 0.013 -0.084 0.0057 0.12 0.094 0.18 1 -0.0086 0.026 -0.06 0.14 -0.032 -0.028 -0.058 -0.11 -0.084 -0.071 -0.12 0.044 -0.15 1 0.044 -0.059 0.093 -0.072 -0.012 0.14 restecq -0.0099 -0.0086 -0.096 thalach -0.071 -0.15 oldpeak -0.096 0.19 0.054 0.0057 -0.059 0.22 0.21 -0.031 -0.17 -0.12 -0.004 -0.06 -0.58 -0.08 -0.1 0.12 0.093 0.39 0.35 slope -0.2 0.28 0.12 -0.18 0.1 0.071 0.14 -0.072 0.12 0.22 0.15 0.21 -0.16 -0.032 -0.012 -0.096 0.21 target --0.14 -0.085 -0.028 0.14 1 9 fbs From the above HeatMap, we can see that cp and thalach are the features with highest positive correlation whereas exang, oldpeak and ca are negatively correlated. While other features do not hold much correlation with the response variable "target" **Outlier Detection** Since the dataset is not large, we cannot discard the outliers. We will treat the outliers as potential observations. # Boxplots In [14]: fig dims = (15,8)fig, ax = plt.subplots(figsize=fig\_dims) sns.boxplot(data=hc\_data, ax=ax); 500 400 300 200 100 trestbps oldpeak chol fbs restecg thalach exang slope thal target Handling Imbalance Imbalance in a dataset leads to inaccuracy and high precision, recall scores. There are certain resampling techniques such as undersampling and oversampling to handle these issues. Considering our dataset, the response variable target has two outcomes "Patients with Heart Disease" and "Patients without Heart Disease". Let us now observe their distribution in the dataset. hc data["target"].value counts() In [15]: 165 Out[15]: 138 Name: target, dtype: int64 From the above chart, we can conclude even when the distribution is not exactly 50:50, but still the data is good enough to use on machine learning algorithms and to predict standard metrics like Accuracy and AUC scores. So, we do not need to resample this dataset. **Train-Test Split** Divide the data into training and test datasets using the train\_test\_split() function #Extracting Variable In [16]: X = hc data.drop("target",axis=1) y = hc data["target"] X.shape (303, 13)Out[17]: y.shape In [18]: (303,)Out[18]: In [19]: # Splitting the dataset into training and test set. from sklearn.model selection import train test split X train, X test, y train, y test = train test split(X, y, test size=0.20, stratify=y, random state=7) **Logistic Regression** from sklearn.linear\_model import LogisticRegression classifier = LogisticRegression() In [21]: classifier.fit(X train, y train) LogisticRegression() Out[21]: In [22]: pred = classifier.predict(X test) In [23]: pred array([1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, Out[23]: 1, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1], dtype=int64) Test Accuracy of the result In [24]: #Creating the Confusion matrix from sklearn.metrics import accuracy score, confusion matrix, classification report In [25]: # Accuracy on Test data accuracy\_score(y\_test, pred) 0.8032786885245902 Out[25]: In [26]: # Accuracy on Train data accuracy score(y train, classifier.predict(X train)) 0.8471074380165289 Out[26]: In [27]: # Building a predictive system import warnings in data = (57,0,0,140,241,0,1,123,1,0.2,1,0,3)# Changing the input data into a numpy array in\_data\_as\_numpy\_array = np.array(in\_data) # Reshaping the numpy array as we predict it in data reshape = in data as numpy array.reshape(1,-1) pred = classifier.predict(in data reshape) print (pred) **if**(pred[0] == 0): print('The person does not have heart disease.') print('The person has heart disease.') The person does not have heart disease. \*\*End\*\*.. #Thank you!!... In [ ]: