

cat_cols_2=pd.get_dummies(cat_cols.restecg, prefix='restecg').iloc[:, 1:] cat_cols_3=pd.get_dummies(cat_cols.slope, prefix='slope').iloc[:, 1:] cat_cols_4=pd.get_dummies(cat_cols.ca, prefix='ca').iloc[:, 1:] cat_cols_5=pd.get_dummies(cat_cols.thal, prefix="thal").iloc[:, 1:] #drop the original multiclass categorical columns In [38]: cat_cols=cat_cols.drop(['cp', 'restecg', 'slope', 'ca', 'thal'], axis=1) #cat_cols In [39]: #concat the numerical columns, binary columns and feature engineered columns df new= pd.concat([num cols, cat cols, cat cols 1, cat cols 2, cat cols 3, cat cols 4, cat c #display the first 5 rows of the new dataframe In [40]: df new.head() Out[40]: trestbps chol thalach oldpeak sex fbs exang cp_1 cp_2 ... slope_1 slope_2 ca_1 ca_2 0 0 0 0 63 233 150 2.3 0 0 0 0 0 0 145 37 130 250 187 3.5 0 0 0 0 0 0 204 0 0 0 0 0 2 41 172 0 0 0 130 1.4 56 236 8.0 0 0 0 0 0 0 3 120 178 0 ... 0 57 120 0 0 0 0 0 0 0 354 163 0.6 5 rows × 23 columns In [41]: #plot a heatmap in order to check the correlation between the feature variables sns.heatmap(df new.corr()) plt.show() -1.00age trestbps chol - 0.75 thalach oldpeak sex - 0.50 fbs exang cp_1 - 0.25 cp_2 cp_3 restecg_1 - 0.00 restecg_2 slope 1 slope_2 -0.25ca_1 ca_2 ca_3 -0.50ca_4 thal_1 thal_2 -0.75thal_3 target thal_1 thal 3 **Correlation Analysis:** The heatmap shows dark areas as negatively correlated and light areas as positively correlated features. From the heatmap we can see that 'thalach', 'cp_2', 'slope_2 and 'thal_2' are positively correlated with the target. 'slope_2' is positively correlated with 'thalach' and negatively correlated with 'oldpeak'. 'thalach' is negatively correlated with 'age' **Model Building** In [42]: #let's split the data into training and testing using train test split model from sklearn.model selection import train_test_split X=df_new.drop(['target'], axis=1) y=df new["target"] In [43]: print(X.shape) print(y.shape) (303, 22)(303,)In [44]: #set the testing size to 20% and training size to 80% X_train, X_test, y_train, y_test=train_test_split(X, y, test size=0.2, random state=1) In [45]: #fit the data using Logistic Regression from sklearn.linear model import LogisticRegression Logreg=LogisticRegression() Logreg.fit(X train, y train) LogisticRegression() Out[45]: In [46]: #check the accuracy score and f1 score accuracy=Logreg.score(X test, y test) print('Accuracy of Logistic Regression : ', accuracy*100) from sklearn.metrics import f1 score y pred=Logreg.predict(X test) f1=f1 score(y pred, y test) print('f1 score : ', f1) Accuracy of Logistic Regression: 77.04918032786885 fl_score : 0.7941176470588235 In [47]: #print the confusion matrix from sklearn.metrics import confusion matrix con_mat=confusion_matrix(y_test, y_pred) print(con_mat) [[20 10] [4 27]] From the confusion matrix, we can see that the model predicted 20 true positives (occurance of a Heart Attack) and 27 true negatives (Non occurance of a Heart Attack) In [48]: f, ax = plt.subplots(figsize = (6,5))sns.heatmap(con_mat) plt.xlabel('predicted value') plt.ylabel('actual value') plt.title('Confusion Matrix of Logistic Regression'); Confusion Matrix of Logistic Regression 0 - 20 value actual - 10 predicted value In [49]: # get the Precision, Recall and f1 score from sklearn.metrics import classification report print(classification_report(y_pred, y_test)) precision recall f1-score support

 0.67
 0.83
 0.74
 24

 0.87
 0.73
 0.79
 37

 0 1

 accuracy
 0.77
 61

 macro avg
 0.77
 0.78
 0.77
 61

 weighted avg
 0.79
 0.77
 0.77
 61

 Precision is the ratio of True Positive and Sum of True Positive and False Positive. Precision gives us the percentage of Positive Cases from Total Predicted cases. Precision = 0.87 Recall is the ratio of True Positive and Sum of True Positive and False Negative. Recall gives us the percentage of how many total Positive cases were Predicted correctly with our model. Recall = 0.73 f1-score gives the combined result of Precision and Recall, f1-score=0.79 In [50]: #import the required library from sklearn.ensemble import RandomForestClassifier from sklearn.metrics import accuracy score In [51]: # let us check for different n_estimators between 10 and 100 in steps of 10 to check the for i in range(10, 100, 10): rfc=RandomForestClassifier(n estimators=i, random state=0) rfc.fit(X train, y train) y predict=rfc.predict(X test) print('n estimators = ', i) accu=(accuracy_score(y_test, y_predict)*100) print('accuracy of Random Forest Classifier = ',accu) print('_ n = 10accuracy of Random Forest Classifier = 72.1311475409836 n = 100accuracy of Random Forest Classifier = 72.1311475409836 $n_{estimators} = 30$ accuracy of Random Forest Classifier = 73.77049180327869 $n_{estimators} = 40$ accuracy of Random Forest Classifier = 75.40983606557377 $n_{estimators} = 50$ accuracy of Random Forest Classifier = 75.40983606557377 n estimators = 60 accuracy of Random Forest Classifier = 78.68852459016394 $n_{estimators} = 70$ accuracy of Random Forest Classifier = 77.04918032786885 $n_{estimators} = 80$ accuracy of Random Forest Classifier = 77.04918032786885 n estimators = 90 accuracy of Random Forest Classifier = 77.04918032786885 The best fit is when $n_{estimators} = 60$, where the accuracy is maximum at 78.6885%. In [52]: #fit the data to a Random Forest Classifier to train the model rfc=RandomForestClassifier(n estimators=60, random state=0) rfc.fit(X_train, y_train) y predict=rfc.predict(X test) In [53]: from sklearn.metrics import accuracy_score, confusion_matrix, classification_report print('Accuracy with Random Forest: ',accuracy_score(y_test, y_predict)*100) print(confusion_matrix(y_test, y_predict)) print(classification report(y test, y predict)) Accuracy with Random Forest: 78.68852459016394 [[21 9] [4 27]] precision recall f1-score support 0.84 0.70 0.76 0.75 0.87 0.81 0 30 31

 accuracy
 0.79
 61

 macro avg
 0.79
 0.79
 0.78
 61

 weighted avg
 0.79
 0.79
 0.79
 61

 Accuracy Check: Accuracy of the Random Forest Classifier is slightly better at 78.6885% as compared to that of Logistic Regression at 77.0492% In [54]: f, ax = plt.subplots(figsize = (6,5)) sns.heatmap(confusion_matrix(y_test, y_predict)) plt.xlabel('predicted value') plt.ylabel('actual value') plt.title('Confusion Matrix of Random Forest Classifier'); Confusion Matrix of Random Forest Classifier 20 actual value - 15 predicted value **Statistical Analysis:** In [55]: #let us check the p-values of all the features using statsmodel for feature selection import statsmodels.api as sm log_reg=sm.Logit(y_train, X_train).fit() Optimization terminated successfully. Current function value: 0.272951 Iterations 8 In [56]: #print the result of the statsmodel Logit Regression print(log reg.summary()) Logit Regression Results ______ Dep. Variable: target No. Observations:
Model: Logit Df Residuals:
Method: MLE Df Model: 220
 Method:
 MLE
 Df Model:
 21

 Date:
 Tue, 02 Aug 2022
 Pseudo R-squ.:
 0.6029

 Time:
 06:52:04
 Log-Likelihood:
 -66.054

 converged:
 True
 LL-Null:
 -166.34

 Covariance Type:
 nonrobust
 LLR p-value:
 2.791e-31
 ______ coef std err z P>|z| [0.025 0.975] ______ For some of the features, the p-value is less than 0.05. There is a strong relationship between these variables and target (occurance or non-occurance of a heart attack). we can choose those features whose p-value is less than 0.05 and also standard error less than 1. In [57]: # compare the predicted values by the statts model against the actual values yhat = log reg.predict(X test) prediction = list(map(round, yhat)) print('Actual values', list(y test.values)) print('Predictions :', prediction) Actual values [0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 0, 0, 1] 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 0, 1] In [58]: print(confusion_matrix(y_test, prediction)) print('Accuracy = ', accuracy score(y test, prediction)*100) [[21 9] [4 27]] Accuracy = 78.68852459016394 In [59]: #Let's select those features whose p-value<0.05 and also standard error<1 and train X reduced=df new[['thalach','trestbps','sex', 'cp 2', 'cp 3', 'ca 1', 'ca 2']] y=df_new[['target']] In [60]: from sklearn.model_selection import train_test_split X reduced train, X reduced test, y train, y test = train test split(X reduced, y, te In [61]: **from** sklearn.linear_model **import** LogisticRegression #fit training data into Logis Logreg=LogisticRegression() Logreg.fit(X reduced train, y train) y predt=Logreg.predict(X reduced test) In [62]: from sklearn.metrics import confusion_matrix, accuracy_score, classification_report print('Accuracy of Logistic Regression = ', accuracy score(y test,y predt)*100) print(confusion_matrix(y_test, y_predt)) print(classification_report(y_test, y_predt)) Accuracy of Logistic Regression = 81.9672131147541 [[22 5] [6 28]] precision recall f1-score support 0 0.79 0.81 0.80 27 1 0.85 0.82 0.84 34

 accuracy
 0.82
 61

 macro avg
 0.82
 0.82
 0.82
 61

 weighted avg
 0.82
 0.82
 0.82
 61

 In [63]: rfc=RandomForestClassifier(n_estimators=550, random_state=0) # fit the training dat # the target values rfc.fit(X_reduced_train, y_train) y_predicted=rfc.predict(X_reduced_test) In [64]: print('Accuracy of Random Forest Classifier = ',accuracy_score(y test, y predicted)* print(confusion_matrix(y_test, y_predicted)) print(classification_report(y_test, y_predicted)) Accuracy of Random Forest Classifier = 75.40983606557377 [[19 8] [7 27]] precision recall f1-score support 0 0.73 0.70 0.72 27 1 0.77 0.79 0.78 34

 accuracy
 0.75
 61

 macro avg
 0.75
 0.75
 0.75

 weighted avg
 0.75
 0.75
 0.75
 61

 After feature selection using statsmodel (p-value < 0.05) the accuracy of the Logistic Regression Model has improved from 77.0492% to 81.9672%. Whereas, the accuracy of Random Forest Classifier Model has reduced from 78.6885% to 75.4098% In [65]: #let's check for the important features using a Random Forest Classifier from sklearn.ensemble import RandomForestClassifier params={'random_state':0, 'n_jobs':-1, 'n_estimators':500, 'max_depth':8} classifier=RandomForestClassifier(**params) classifier=classifier.fit(X,y) important_features = pd.Series(data=classifier.feature_importances_, index=X.column plt.figure(figsize=(10,12)) plt.title("Feature importance") ax = sns.barplot(y=important_features.index, x=important_features.values, palette=" plt.show() Feature importance thalach thal_2 oldpeak thal_3 age chol trestbps exang slope_2 cp 2 ca_1 slope_1 sex ca_2 cp_3 restecg_1 ca_3 cp_1 fbs thal 1 ca_4 restecg 2 0.00 0.02 0.04 0.06 0.08 0.10 0.12 From the above graph, we see that 'thalach' or Maximum Heart Rate Achieved is the most important feature in predicting a Heart Attack, followed by thal_2 (Thalassemia_2, oldpeak(ST depression induced by exercise relative to rest), Thalassemia_3, age etc. #Let's train the model using the first five important features from the Random Fox In [66]: X reduce=df new[['thalach', 'thal 2', 'oldpeak', 'thal 3', 'age']] y=df_new[['target']] In [67]: from sklearn.model_selection import train test split X reduce train, X reduce test, y train, y test = train test split(X reduce, y, test Logreg=LogisticRegression() In [68]: Logreg.fit(X_reduce_train, y_train) y_pred1=Logreg.predict(X_reduce_test) print('Accuracy of Logistic Regression = ', accuracy_score(y_test,y_pred1)*100) print(confusion_matrix(y_test, y_pred1)) print(classification_report(y_test, y_pred1)) Accuracy of Logistic Regression = 83.60655737704919 [[21 6] [4 30]] precision recall f1-score support 0.78 0 0.84 0.81 27 0.83 0.88 0.86 34 0.84 61 accuracy 0.84 macro avg 0.83 0.83 61 weighted avg 0.84 0.84 In [69]: rfc=RandomForestClassifier(n_estimators=100, random_state=0) rfc.fit(X_reduce_train, y_train) y pred2=rfc.predict(X reduce test) print('Accuracy of Random Forest Classifier = ',accuracy_score(y_test, y_pred2)*1 print(confusion_matrix(y_test, y_pred2)) print(classification_report(y_test, y_pred2)) Accuracy of Random Forest Classifier = 81.9672131147541 [[23 4] [7 27]] precision recall f1-score support 0 0.77 0.85 0.81 27 1 0.87 0.79 0.83 34

 accuracy
 0.82
 61

 macro avg
 0.82
 0.82
 61

 weighted avg
 0.82
 0.82
 0.82
 61

 After using features from Random Forest Classifier Feature Selection, the accuracy of Logistic Regression has slightly improved to 83.6066% and the accuracy of Random Forest Classifier to 81.9672% Inference: After selecting all the features: Accuracy of Logistic Regression with all the features: 77.0492% f1_score : 0.79 Accuracy of Random Forest Classifier with all the features: 78.6885% f1_score : 0.81 After selecting features based on p-value from statts model: Accuracy of Logistic Regression with 7 features: 81.9672% f1_score : 0.84 Accuracy of Random Forest Classifier with 7 features: 75.4098% f1_score: 0.78 Using Feature Selection by Random Forest Classifier: Accuracy of Logistic Regression using 5 important features: 83.6066% f1_score: 0.86 Accuracy of Random Forest Classifier using 5 important features: 81.9672% f1_score:0.83 In []:

In [37]: cat_cols_1=pd.get_dummies(cat_cols.cp, prefix='cp').iloc[:, 1:]