	import matplotlib.pyplot as plt %matplotlib inline import seaborn as sns import warnings warnings, filterwarnings('ignore') from sklearn.model_selection import train_test_split from sklearn.metrics import * from sklearn.preprocessing import LabelEncoder, StandardScaler from sklearn.linear_model import LogisticRegression from statsmodels.stats.outliers_influence import variance_inflation_factor from sklearn.ensemble import RandomForestClassifier import graphviz from sklearn.tree import export_graphviz from IPython.display import Image L. Data Analysis:
7	a. Import the dataset substance is b. Get information about the dataset (mean, max, min, quarties etc.) # a. Import the dataset substance is visual fields # a. Import the d
2 3 3 3	3 56 1 1 1 120 236 0 1 178 0 0.8 2 0 2 1 4 57 0 0 0 120 354 0 1 163 1 0.6 2 0 2 1 298 57 0 0 0 140 241 0 1 123 1 0.2 1 0 3 0 299 45 1 3 110 264 0 1 132 0 1.2 1 0 3 0 300 68 1 0 144 193 1 1 141 0 3.4 1 2 3 0 301 57 1 0 130 131 0 1 115 1 1.2 1 1 3 0 302 57 0 1 1 30 236 0 0 174 0 0.0 1 1 2 0 303 rows × 14 columns
< R D	df.info() **Cclass '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 2 cp 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 6 restecg 303 non-null int64 7 thalach 303 non-null int64
d m	8 exang
a	mem 54.366337 0.683168 0.966997 131.623762 246.264026 0.148515 0.528053 149.646865 0.326733 1.039604 1.399340 0.729373 2.313531 0.544554 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 0.612277 0.498835 min 29.000000 0.000000 0.000000 120.00000 126.000000 0.000000 0.000000 0.000000 0.000000
c t c f r t e o s c t t d	sex 0 trestbps 0 thoi 0 fbs 0 restecg 0 thalach 0 exang 0 poldpeak 0 slope 0 target 0 dtype: int64 # c. Find the correlation between all fields df.corr()
	4 5 5 5 1
2	oldpeak 0.210013 0.096093 -0.149230 0.193216 0.05952 0.005477 -0.344187 0.288223 1.00000 -0.577537 0.222682 0.210244 -0.43696 slope -0.168814 -0.030711 0.119717 -0.121475 -0.004038 -0.059894 0.093045 0.386784 -0.257748 -0.577537 1.00000 -0.08155 -0.104764 0.345877 ca 0.276326 0.118261 -0.181053 0.101389 0.07511 0.13779 -0.072042 -0.21377 0.115739 0.222682 -0.080155 1.00000 -0.151832 -0.391724 target -0.225439 -0.280937 0.434931 -0.085239 -0.028046 0.137230 0.421741 -0.436757 -0.436656 0.345877 -0.391724 -0.344029 1.000000 2. Data Visualization: * a. Visualization has a patient has disease and not having a heart disease and not having a h
; ;	* a. Visualize the number of patients having a heart disease and not having a heart disease sms.countplot(data=df, x='target') plt.title("Chance of heart disease") plt.xticks(ticks=[0, 1], labels=['No', 'Yes']) plt.show() Chance of heart disease 160 - 140 -
	120 - 100 - 80 - 40 - 20 -
: !	# b. Visualize the age and whether a patient has disease or not sns.boxplot(x='target', y='age', data=df) plt.xlabel('Has Disease') plt.ylabel('ge') plt.title('Distribution of Age for Patients with and without Disease') plt.show() Distribution of Age for Patients with and without Disease
	70 - 60 - 9 50 - 40 -
	# c. Visualize correlation between all features using a heat map plt.subplots(figsize=(12,6)) sns.heatmap(data=df.corr(), annot=True) plt.show() age - 1
	sex - 0.098
	oldpeak - 0.21 0.096 -0.15 0.19 0.054 0.0057 -0.059 -0.34 0.29 1 -0.58 0.22 0.21 -0.43 slope - 0.17 -0.031 0.12 -0.12 -0.004 -0.06 0.093 0.39 -0.26 -0.58 1 -0.08 -0.1 0.35 ca - 0.28 0.12 -0.18 0.1 0.071 0.14 -0.072 -0.21 0.12 0.22 -0.08 1 0.15 -0.39 thal - 0.068 0.21 -0.16 0.062 0.099 -0.032 -0.012 -0.096 0.21 0.21 -0.1 0.15 1 -0.34 target - 0.23 -0.28 0.43 -0.14 -0.085 -0.028 0.14 0.42 -0.44 -0.43 0.35 -0.39 -0.34 1 age sex cp trestbps chol fbs restecg thalach exang oldpeak slope ca thal target
V	<pre>Creating models //F col_list = [] for col in df.columns: iff((df[col].dtype!='object')&(col!='target')): col_list.append(col) x = df[col_list] vif_data = pd.obatFrame() vif_data['feature'] = X.columns vif_data['VIF'] = [variance_inflation_factor(X.values, i) for i in range(len(X.columns))] print(vif_data)</pre>
1	feature VIF 1
	<pre>col_list = [] for col in df.columns: if((df[col].dtype!='object')&(col!='target')): col_list.append(col) X = df[col_list] vif_data = pd.DataFrame() vif_data['feature'] = X.columns vif_data['VIF'] = [variance_inflation_factor(X.values, i) for i in range(len(X.columns))] print(vif_data) feature</pre>
1	age 2.39990 1 sex 3.522376 2 cp 2.403644 3 chol 25.923434 4 fbs 1.250256 5 restecg 2.058206 6 thalach 34.026519 7 exang 1.990878 8 oldpeak 2.973430 9 slope 10.067190 10 ca 1.808102 11 thal 17.005949 df.drop(columns=['thalach'], inplace=True)
	<pre>for col in df.columns: if((df[col].dtype!='object')&(col!='target')): col_list.append(col) X = df[col_list] vif_data = pd.DataFrame() vif_data['feature'] = X.columns vif_data['VIF'] = [variance_inflation_factor(X.values, i) for i in range(len(X.columns))] print(vif_data) feature</pre>
	<pre>5 restecg 2.022210 6 exang 1.955987 7 oldpeak 2.965697 8 slope 8.372679 9</pre>
,	fbs 1.232428 4 restecg 2.006017 5 exang 1.947640 6 oldpeak 2.827322 7 slope 7.983150 8 ca 1.710828
	<pre>df.drop(columns=['chol'], inplace=True) col_list = [] for col in df.columns: iff((df[col].dtype!='object')&(col!='target')): col_list.append(col) X = df[col_list] vif_data = pd.DataFrame() vif_data['vIF'] = [variance_inflation_factor(X.values, i) for i in range(len(X.columns))] print(vif_data)</pre>
	feature VIF sex 3.348022 1
	<pre>col_ist = [] for col in df.columns: if((df[col].dtype!='object')&(col!='target')): col_list.append(col) X = df[col_list] vif_data = pd.DataFrame() vif_data['feature'] = X.columns vif_data['VIF'] = [variance_inflation_factor(X.values, i) for i in range(len(X.columns))] print(vif_data) feature</pre>
a	siope 3.756703 Ca 1.686170 3. Logistic Regression: a. Build a simple logistic regression model: i. Divide the dataset in 70:30 ratio ii. Build the model on train set and predict the values on test set iii. Build the confusion matrix and get the accuracy score
	<pre>X = df.drop('target', axis=1) y = df['target'] # i. Divide the dataset in 70:30 ratio X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.7, random_state=0) # ii. Build the model on train set and predict the values on test set LR = LogisticRegression() LR.fit(X_train, y_train) * LogisticRegression</pre> <pre>conititieRegression()</pre>
i i	<pre>pred_lr = LR.predict(X_test) # iii. Build the confusion matrix and get the accuracy score print('Confusion Matrix:\n',confusion_matrix(y_test, y_pred_lr)) print('Accuracy Score ', accuracy_score(y_test, y_pred_lr)) Confusion Matrix: [[30 14] [7 40]] Accuracy Score 0.7692307692307693</pre>
a	5. Random Forest: a. Build a Random Forest model: i. Divide the dataset in 70:30 ratio ii. Build the model on train set and predict the values on test set iii. Build the confusion matrix and calculate the accuracy iv. Visualize the model using the Graphviz package # i. Divide the dataset in 70:30 ratio x_train, x_test, y_train, y_test = train_test_split(x, y, train_size=0.7, random_state=0)
R	<pre># ii. Build the model on train set and predict the values on test set rf = RandomForestClassifier() rf.fit(X_train, y_train) v RandomForestClassifier RandomForestClassifier() y_pred_rf = rf.predict(X_test) # iii. Build the confusion matrix and calculate the accuracy</pre>
C A	<pre>print('Confusion Matrix:\n',confusion_matrix(y_test, y_pred_rf)) print('Accuracy Score ', accuracy_score(y_test, y_pred_rf)) Confusion Matrix: [[31 13] [6 41]]</pre>
9	# Use Graphviz to render the decision tree graph = graphviz.Source(dot_data, format='png') graph.render('decision_tree', format='png', cleanup=True, quiet=True) display(graph) examg < 0.5
	color colo
	Cas 0.5 Simples = 1 Simp
6	Select the best model Select the select model Select the best model Select the best model Select the select model Select model Select the select model Select the select model Select mode
C	 b. Print the classification report of all classifiers c. Calculate Recall Precision and F1 score of all the models d. Visualize confusion matrix using heatmaps e. Select the best model based on the best accuracies # a. Print the confusion matrix of all classifiers print('Confusion Matrix of Logistic Regression:\n',confusion_matrix(y_test, y_pred_lr)) print('Confusion Matrix of Random Forest:\n',confusion_matrix(y_test, y_pred_rf))
; 	# b. Print the classification report of all classifiers print('Classification Report of Logistic Regression:\n', classification_report(y_test, y_pred_lr)) print('Classification Report of Random Forest:\n', classification_report(y_test, y_pred_rf)) Classification Report of Logistic Regression: precision recall f1-score support 0 0.81 0.68 0.74 44 1 0.74 0.85 0.79 47 accuracy 0.77 91
W	macro avg
	<pre>print('F1 Score of Logistic Regression:', f1_score(y_test, y_pred_lr)) print('Pecall Score of Random Forest:', recall_score(y_test, y_pred_rf)) print('F1 Score of Random Forest:', f1_score(y_test, y_pred_rf)) Recall Score of Logistic Regression: 0.851063829787234 F1 Score of Logistic Regression: 0.7920792079207921 Recall Score of Random Forest: 0.8723404255319149 F1 Score of Random Forest: 0.8723404255319149 F1 Score of Random Forest: 0.8118811881188118 # d. Visualize confusion matrix using heatmaps plt.figure(figsize=(16,6)) plt.subplot(1,2,1) sns.heatmap(confusion_matrix(y_test, y_pred_lr), annot=True)</pre>
	sns.heatmap(confusion_matrix(y_test, y_pred_lr), annot=True) plt.subplot(1,2,2) sns.heatmap(confusion_matrix(y_test, y_pred_rf), annot=True) plt.show() Heatmap of Logistic Regression -40 -35 Heatmap of Random Forest -40 -35
ı	o - 30 14 o - 31 13 - 30 - 30