In [114	<pre>import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns import warnings warnings.filterwarnings('ignore')</pre> <pre>Reading Data</pre>
In [115 In [116 Out[116]:	0 1 35 management married tertiary no 231 yes no unknown 5 may 139 1 -1 0 unknown 1 2 28 management single tertiary no 447 yes yes unknown 5 may 217 1 -1 0 unknown 2 3 42 entrepreneur divorced tertiary yes 2 yes no unknown 5 may 380 1 -1 0 unknown 3 4 58 retired married primary no 121 yes no unknown 5 may 50 1 -1 0 unknown 4 5 43 technician single secondary no 593 yes no unknown 5 may 55 1 -1 0 unknown 4 5 43 technician single secondary no 82 yes no unknown 9 may 654 1 -1 0 unknown 422 423 26 blue-collar single secondary no 82 yes no unknown 9 may 1692 2 -1 0 unknown 423 424 40 blue-collar married secondary no 10 yes no unknown 9 may 1692 2 -1 0 unknown 425 426 33 blue-collar single secondary no -197 yes no unknown 9 may 2016 2 -1 0 unknown 426 427 45 blue-collar married primary no 390 yes no unknown 12 may 460 5 -1 0 unknown 427 rows × 17 columns Data Cleaning Checking duplicate df. duplicated() 9 False 1 False
In [117	7 False 7 False 8 False 9 False 9 False 1 Fals
Out[117]: In [118 Out[118]:	riodsing, loan, contact, day, month, duration, campaign, 'pdays', 'previous', 'poutcome'], dtype='object') Rows in each columns df.count()
	df.shape (427, 17) checking null values df.isna().sum() Id 0 age 0 job 0 marital 0 education 0 default 0 balance 0
In [121	housing 0 loan 0 contact 0 day 0 month 0 duration 0 campaign 0 pdays 0 previous 0 previous 0 dtype: int64 droping unnecessary column df.drop(['contact', 'month', 'pdays', 'previous', 'poutcome', 'housing', 'campaign', 'day', 'duration'], axis=1, inplace=True)
In [122	<pre>df.drop(['Id'],axis=1,inplace=True)</pre> Top 5 rows df.head()
In [124 Out[124]:	3 58 retired married primary no 121 no 4 43 technician single secondary no 593 no Checking unique values in all columns df.nunique() age 40 job 12 marital 3 education 4 default 2 balance 310 loan 2
In [125 Out[125]: In [126 Out[126]:	<pre>checking unique values for particular column df['education'].unique() array(['tertiary', 'primary', 'secondary', 'unknown'], dtype=object) edu=df.education.value_counts() edu secondary 242 tertiary 87 primary 58</pre>
Out[127]: In [128 Out[128]: In [129	'services', 'blue-collar', 'self-employed', 'unemployed', 'housemaid', 'student', 'unknown'], dtype=object) df['marital'].unique() arroy(['married', 'single', 'diverced'], dtype=object)
In [130	<pre>df['loan'].unique() array(['no', 'yes'], dtype=object) checking types of columns df.dtypes</pre>
In [132 Out[132]:	balance object dtype: object Stastical information df.describe() age balance count 427,000000 427,000000 427,000000 mean 43,971897 640,302108 std 9.589908 2722.544765 min 22,000000 -674,0000000 25% 36,000000 23,500000 50% 43,000000 168,000000 75% 53,000000 45248,000000 45248,000000 45248,000000 45248,000000 45248,000000 45248,000000
In [136	Converting data types into integer df.job=df.job.replace({'management':1, 'entrepreneur':2, 'retired':3, 'technician':4, 'admin.':5, 'services':6, 'blue-collar':7, 'self-employed':9, 'hou df.marital=df.marital.replace({'married':1, 'single':2, 'divorced':3, 'unknown':4}) df.education=df.education.replace({'primary':1, 'secondary':2, 'tertiary':3, 'unknown':4}) df.loan=df.loan.replace({'yes':1, 'no':2}) df.default=df.default.replace({'yes':1, 'no':2})
Out[138]:	0 35 1 1 3 2 231 2 1 28 1 2 3 2 447 1 2 42 2 3 3 1 2 2 3 58 3 1 1 2 121 2 4 43 4 2 2 2 593 2 df.dtypes
Out[139]: In [140 Out[140]:	age int64 job int64 marital int64 education int64 default int64 balance int64 loan int64 dtype: object Heat Map df.corr() age job marital education default balance loan
In [141	age 1.000000 -0.119984 -0.085721 0.010033 -0.010437 0.004254 0.043394 job -0.119984 1.00000 -0.064153 -0.270193 -0.001241 -0.092210 0.056652 marital -0.085721 -0.064153 1.000000 -0.015460 -0.095898 0.058969 0.006856 education 0.010033 -0.270193 -0.015460 1.000000 0.000636 0.127697 0.078441 default -0.010437 -0.001241 -0.095898 0.00636 1.000000 0.022897 -0.037391 balance 0.004254 -0.092210 0.058969 0.127697 0.022897 1.000000 0.014398 loan 0.043394 0.056652 0.006856 0.078441 -0.037391 0.014398 1.000000 sns.heatmap(df.corr(),cmap='RdBu',annot=True,fmt='.1f') plt.show()
	age - 10
In [142 In [143	We can see there is a positive correlation between age & loan This makes sense since, The greater age results in a greater chance of loan. In addition, we see a negative correlation between job & our predictor i.e-loan. Seprating dependent and independent variable x=df.drop(['loan'], axis=1) y=df.loan
Out[143]: In [144 Out[144]:	0 35 1 1 3 2 231 1 28 1 2 3 2 447 2 42 2 3 3 1 2 3 58 3 1 1 2 121 4 43 4 2 2 2 593 422 26 7 2 2 2 82 423 40 7 1 2 2 10 424 28 7 2 2 2 -197 425 33 7 3 1 2 390 426 45 7 1 1 2 101 427 rows × 6 columns
In [145 Out[145]:	1
In [146 In [147 In [148	Name: loan, dtype: int64 Importing train-test split from sklearn.model_selection import train_test_split x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2) print(x_train.shape) print(y_train.shape) print(y_train.shape) print(y_test.shape) print(y_test.shape) (341, 6)
In [149 In [150 In [151 Out[151]:	(86, 6) (341,) (86,) 1-Logistic Regression from sklearn.linear_model import LogisticRegression model=LogisticRegression() model LogisticRegression()
In [152 Out[152]: In [153 Out[153]: In [154 Out[154]: In [155	model.score(x_train,y_train)*100 86.21700879765396 model.score(x_test,y_test)*100
Out[156]: In [157	<pre>y_predict array([2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2</pre>
In [159 In [160	94
In [161 Out[161]: In [162 In [163 In [164 Out[164]:	<pre>importing confusion matrix from sklearn.metrics import confusion_matrix performance=confusion_matrix(y_test,y_predict) performance performance</pre>
In [165 In [166	plot_confusion_matrix(model, x_test, y_test) plot_show()
	There are 8 type 1 error (False Positives)- You predicted positive and it's false. There are 0 type 2 error (False Negatives)- You predicted negative and it's false. Hence if we calculate the accuracy its # Correct Predicted/ # Total. In other words, where TP, FN, FP and TN represent the number of true positives, false negatives, false positives and true negatives. (TP + TN)/(TP + TN + FP + FN). (0 + 78)/(0+78+8+0) = 0.9069 = 90.69% accuracy Accuracy = 90.69% AS we are getting near by same accuracy in both training and testing so our model is balkanced fit.
In [167 In [168 In [169 In [170 In [171 In [172	Classification report from sklearn.metrics import classification_report performance=classification_report(y_test,y_predict) print(performance) precision recall f1-score support 1 0.00 0.00 0.00 8 2 0.01 1.00 0.95 78 accuracy 0.45 0.50 0.48 86 weighted avg 0.82 0.91 0.86 86 2-Decision Tree from sklearn.tree import DecisionTreeClassifier model_dt=DecisionTreeClassifier() model_dt=fit(x_train,y_train)
Out[172]: In [175 Out[175]: In [176 Out[176]: In [177 In [178 In [179 Out[179]:	DecisionTreeClassifier() model_dt.score(x_train,y_train)*100 100.0 model_dt.score(x_test,y_test)*100 75.5813953488372 from sklearn.metrics import accuracy_score y_predict=model_dt.predict(x_test) y_predict array([2, 2, 1, 2, 2, 2, 2, 1, 2, 2, 2, 1, 2, 2, 2, 2, 2, 2, 2, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
In [180 In [181 Out[181]: In [182 In [183 In [184	test_acc = accuracy_score(y_test,y_predict)*100 test_acc 75.5813953488372 3-Random Forest from sklearn.ensemble import RandomForestClassifier model_rf=RandomForestClassifier(n_estimators=100) model_rf.fit(x_train,y_train)
Out[185]: In [186 Out[186]: In [187 In [188 Out[189]:	<pre>model_rf.score(x_train,y_train)*100 100.0 model_rf.score(x_test,y_test)*100 90.69767441860465 from sklearn.metrics import accuracy_score y_predict=model_rf.predict(x_test) y_predict array([2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2</pre>
	test_acc 90.69767441860465 4-KNeighborsClassifier from sklearn.neighborsClassifier model_knn=KNeighborsClassifier(n_neighbors=10) model_knn.fit(x_train,y_train)
Out[194]: In [195 Out[195]: In [196 Out[196]: In [197	<pre>KNeighborsClassifier(n_neighbors=10) model_knn.score(x_train,y_train)*100 86.21700879765396 model_knn.score(x_test,y_test)*100 89.53488373003024</pre>
In [200 Out[200]:	<pre>y_predict array([2, 2, 2, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,</pre>
	Choservation From comparing the 4 models, we can conclude that Logistic regression and RandomForestClassifier Model has highest accuracy. With an accuracy of 90.69%. We have precision, recall, f1-score and support: Precision: be "how many are correctly classified among that class" Recall: "how many are correctly classified among that class." F1-score: harmonic mean of precision and recall values. F1 score reaches its best value at 1 and worst value at 0.F1 Score = 2 x ((precision x recall) / (precision + recall)) Support: of samples of the true response that lie in that class. Analysis of result We used precision. F1-score; recall and accuracy evaluation metrics for evaluating our models. False Positive(FP) is when a model incorrectly predicts a positive outcome. False Negative(FN) is when a model incorrectly predicts the negative outcome. True Positive(TP) is when model correctly predicts a negative outcome. Precision=FP/(TP+FP) = Recall = TP / (TP+FN) F1 score = 2 precision Recall / (precision + Recall) Machine Learning Models - Accuracy 1 - Logistic Regression - 74.50 2 - Decision Tree - 62.43 3 - Random Forest - 82.50 4 - KNaginborsClassifier - 82.50
In []:	Conclusions Our Random Forest algorithmand Logistic regression yields the highest accuracy, 90.69 %. Any accuracy above 80% is considered good. Thus, 90.69 % is the ideal accuracy! and Random forest model has and Logistic Regression Highest testing score i.e-90.69% and score of train data is 86%. Out of the 17 features we examined the 6 features i.e-age,job,marital,education,default,balance features that helped us classify between a yes and no (loan) As we are getting high accuracy in training as well as testing dataset,it is example of Right fit.