

```
1 import pandas as pd
```

```
1 import warnings
```

```
2 warnings.filterwarnings("ignore")
```

```
1 data = pd.read_csv("/content/ML Case Study - Data.csv")
```

```
1 import pandas as pd
```

```
2 import numpy as np
```

```
3 import matplotlib.pyplot as plt
```

```
4 %matplotlib inline
```

```
5 import seaborn as sns
```

```
6 pd.options.display.float_format = '{:,.2f}'.format
```

```
7 from IPython.core.display import display, HTML
```

```
8 display(HTML("<style>.container {width : 98% !important; }</style>"))
```

```
9 plt.style.use('ggplot')
```

```
10 plt.rcParams['figure.figsize'] = [10,10]
```

```
11 pd.set_option('display.max_columns', 500)
```

```
12 pd.set_option('display.max_rows', 500)
```

```
↳
```

```
1 from scipy import stats
```

```
2 from sklearn import metrics
```

```
1 import yellowbrick as yb
```

▼ Checking Data Types

```
1 data.info()
```

```
↳
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5000 entries, 0 to 4999
Data columns (total 14 columns):
#   Column              Non-Null Count  Dtype
---  -
0   ID                   5000 non-null   int64
1   Age                  5000 non-null   int64
```

▼ Statistical Summary

```
6   CCAvg                5000 non-null   float64
```

```
1 data.describe()
```

	ID	Age	Experience	Income	ZIP Code	Family	CCAvg	Education	Mo
count	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5
mean	2,500.50	45.34	20.10	73.77	93,152.50	2.40	1.94	1.88	
std	1,443.52	11.46	11.47	46.03	2,121.85	1.15	1.75	0.84	
min	1.00	23.00	-3.00	8.00	9,307.00	1.00	0.00	1.00	
25%	1,250.75	35.00	10.00	39.00	91,911.00	1.00	0.70	1.00	
50%	2,500.50	45.00	20.00	64.00	93,437.00	2.00	1.50	2.00	
75%	3,750.25	55.00	30.00	98.00	94,608.00	3.00	2.50	3.00	
max	9,999.00	99.00	40.00	999.00	99,999.00	4.00	4.00	4.00	

```
1 data.describe()
```

	ID	Age	Experience	Income	ZIP Code	Family	CCAvg	Education	Mo
count	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5
mean	2,500.50	45.34	20.10	73.77	93,152.50	2.40	1.94	1.88	
std	1,443.52	11.46	11.47	46.03	2,121.85	1.15	1.75	0.84	
min	1.00	23.00	-3.00	8.00	9,307.00	1.00	0.00	1.00	
25%	1,250.75	35.00	10.00	39.00	91,911.00	1.00	0.70	1.00	
50%	2,500.50	45.00	20.00	64.00	93,437.00	2.00	1.50	2.00	
75%	3,750.25	55.00	30.00	98.00	94,608.00	3.00	2.50	3.00	
max	9,999.00	99.00	40.00	999.00	99,999.00	4.00	4.00	4.00	

```
1 # from pandas_profiling import ProfileReport
2 # prof = ProfileReport(data)
3 # prof.to_file("profile.html")
```

▼ Checking for Null Values

There were none

```
1 data.isnull().sum()
```

```

↳ ID          0
   Age         0
   Experience   0
   Income       0
   ZIP Code     0
   Family       0
   CCAvg        0
   Education    0
   Mortgage     0
   Personal Loan 0
   Securities Account 0
   CD Account    0
   Online        0
   CreditCard    0
   dtype: int64

```

▼ Checking For Incorrect Values

```
1 data.describe()
```

```

↳

```

	ID	Age	Experience	Income	ZIP Code	Family	CCAvg	Education	Mo
count	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5
mean	2,500.50	45.34	20.10	73.77	93,152.50	2.40	1.94	1.88	
std	1,443.52	11.46	11.47	46.03	2,121.85	1.15	1.75	0.84	
min	1.00	23.00	-3.00	8.00	9,307.00	1.00	0.00	1.00	
25%	1,250.75	35.00	10.00	39.00	91,911.00	1.00	0.70	1.00	
50%	2,500.50	45.00	20.00	64.00	93,437.00	2.00	1.50	2.00	
75%	3,750.25	55.00	30.00	98.00	94,608.00	3.00	2.50	3.00	
max	5,000.00	67.00	40.00	99.00	99,951.00	4.00	4.00	4.00	

▼ Incorrect values are in experience as it cannot be negative

```

1 for i,k in enumerate(data["Experience"]):
2     if k < 0:

```

```
3 data["Experience"][i] = 0
```

```
1 data.columns
```

```
↳ Index(['ID', 'Age', 'Experience', 'Income', 'ZIP Code', 'Family', 'CCAvg',
        'Education', 'Mortgage', 'Personal Loan', 'Securities Account',
        'CD Account', 'Online', 'CreditCard'],
        dtype='object')
```

▼ Corrected the Experience Values

```
1 data.describe()
```

```
↳
```

	ID	Age	Experience	Income	ZIP Code	Family	CCAvg	Education	Mo
count	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5,000.00	5
mean	2,500.50	45.34	20.12	73.77	93,152.50	2.40	1.94	1.88	
std	1,443.52	11.46	11.44	46.03	2,121.85	1.15	1.75	0.84	
min	1.00	23.00	0.00	8.00	9,307.00	1.00	0.00	1.00	
25%	1,250.75	35.00	10.00	39.00	91,911.00	1.00	0.70	1.00	
50%	2,500.50	45.00	20.00	64.00	93,437.00	2.00	1.50	2.00	
75%	3,750.25	55.00	30.00	98.00	94,608.00	3.00	2.50	3.00	
max	5,000.00	67.00	40.00	224.00	99,954.00	4.00	4.00	4.00	

▼ Unique value in each column

```
1 for i in list( data.columns):
2     print(i,data[i].nunique())
```

```
↳
```

```
ID 5000
```

```
1
```

```
income 102
```

➤ Number Of People with 0 Mortgage

```
Mortgage 347
```

```
1 (data["Mortgage"] == 0).sum()
```

```
↳ 3462
```

```
CreditCard 2
```

➤ Number of people 0 credit card spending

```
1 (data["CCAvg"] == 0).sum()
```

```
↳ 106
```

```
1 data.drop("ID", axis = 1, inplace= True)
```

```
1 data.drop("ZIP Code", axis = 1, inplace= True)
```

➤ Value Counts of different categorical variable

```
1 for i in list( ['Family',
2               'Education', 'Personal Loan', 'Securities Account',
3               'CD Account', 'Online', 'CreditCard']):
4     print(data[i].value_counts())
5     print()
```

```
↳
```

```

1    1472
2    1296
4    1222
3    1010
Name: Family, dtype: int64

```

```

1    2096
3    1501
2    1403
Name: Education, dtype: int64

```

```

0    4520
1     480
Name: Personal Loan, dtype: int64

```

```

0    4478
1     522

```

```
1 pd.pivot_table(data , values="Personal Loan", index = "Family",columns="Education", aggfun
```

```

↳ Education  1    2    3
      Family
1          9  40  58
2          4  50  52
3         40  44  49
4         40  48  46

```

Can be seen than that as the family and and education grows the number of people that go for personal loans have increased

▼ Data Preproceesing

```

1 x = data.loc[:,['Age', 'Experience', 'Income', 'Family', 'CCAvg', 'Education',
2               'Mortgage', 'Securities Account', 'CD Account',
3               'Online', 'CreditCard']]
4 y = data["Personal Loan"]
5 print(x.shape,y.shape)

```

```
↳ (5000, 11) (5000,)
```

```

1 dic = {0:0,1:1}
2 for k in y:
3     dic[k] +=1
4 # temp = dic[0]
5 # dic[0] = dic[1]

```

```
6 # dic[1] = temp
```

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, Y_train, Y_test = train_test_split(x,y,test_size = 0.3,random_state=7)
```

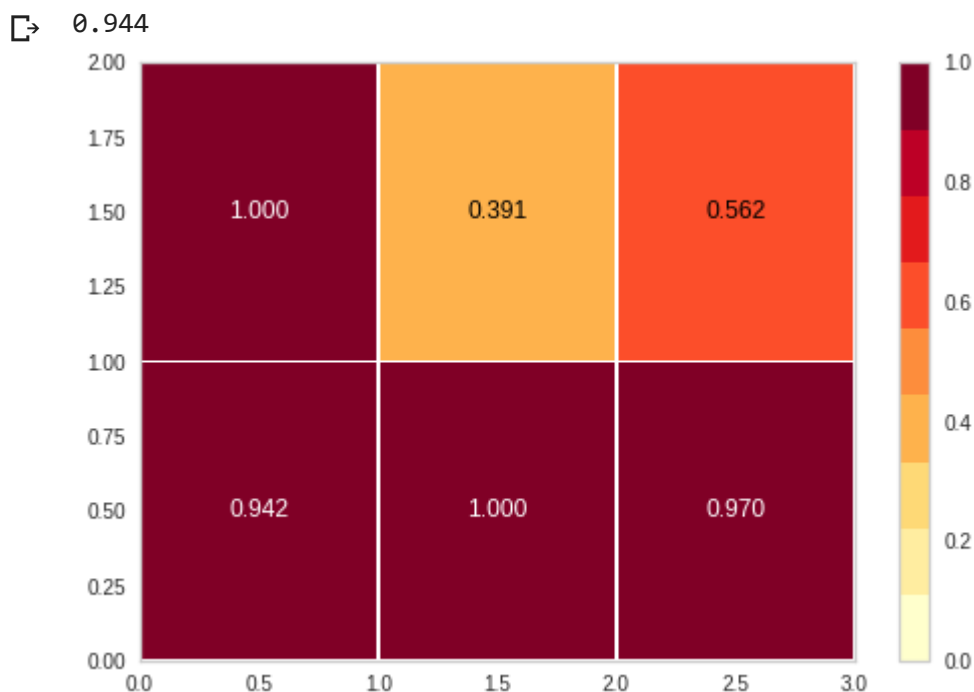
▼ First Model Logistic Regression

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import accuracy_score
3 lg_reg = LogisticRegression(max_iter = 100000, class_weight=dic)
4 lg_reg.fit(X_train, Y_train)
5 print("test_accuracy",lg_reg.score(X_test,Y_test))
6 print("train_accuracy",lg_reg.score(X_train,Y_train))
```

```
↳ test_accuracy 0.944
   train_accuracy 0.9382857142857143
```

```
1 y_pred = lg_reg.predict(X_test)
```

```
1 from yellowbrick.classifier import ClassificationReport
2 viz = ClassificationReport(lg_reg)
3 viz.fit(X_train, Y_train)
4 viz.score(X_test, Y_test)
5
```



```
1 lg_reg.get_params()
```

```

{
  'C': 1.0,
  'class_weight': {0: 4520, 1: 481},
  'dual': False,
  'fit_intercept': True,
  'intercept_scaling': 1,
  'l1_ratio': None,
  'max_iter': 100000,
  'multi_class': 'auto',
  'n_jobs': None,
  'penalty': 'l2',
  'random_state': None,
  'solver': 'lbfgs',
  'tol': 0.0001,
  'verbose': 0,
  'warm_start': False}

```

```

1 from sklearn.metrics import confusion_matrix
2 confusion_matrix(Y_test,y_pred)

```

```

array([[1362,    0],
       [  84,   54]])

```

▼ Performance Metrics

We choose ***f1_score*** as the evaluation metric due to the class size imbalance. 1 are 481 and 0 are 4520.

Metrics from <http://onlineconfusionmatrix.com/>

Measure	Value
Sensitivity	0.8365
Specificity	0.9635
Precision	0.6304
Negative Predictive Value	0.9875
False Positive Rate	0.0365
False Discovery Rate	0.3696
False Negative Rate	0.1635
Accuracy	0.9547
F1 Score	0.7190
Matthews Correlation Coefficient	0.7081

▼ Tweaking Hyperparameters

Since the most there is a imbalance in the number of dataset of each class.

```

1 from sklearn.metrics import f1_score
2 Lr_rate = [0.001,0.01,0.1,0.3,0.9,1.0,3,10.0]
3 for i in Lr_rate:
4     lg_reg = LogisticRegression(C = i,max_iter =1000,class_weight="balanced")
5     lg_reg.fit(X_train, Y_train)
6     Y_pred = lg_reg.predict(X_test)
7     print(f1_score(Y_test,Y_pred),"C =" + str(i))

```

```

☞ 0.5617977528089888 C =0.001
   0.5794392523364487 C =0.01
   0.5975903614457831 C =0.1
   0.5980861244019139 C =0.3
   0.6024096385542169 C =0.9
   0.6038647342995169 C =1.0
   0.6038647342995169 C =3
   0.6024096385542169 C =10.0

```

Ways To Improving our model

Since there are 5 categorical variables in the data using **Random Forest** and **decision trees** classifiers make more sense. But we will work our way upto it.

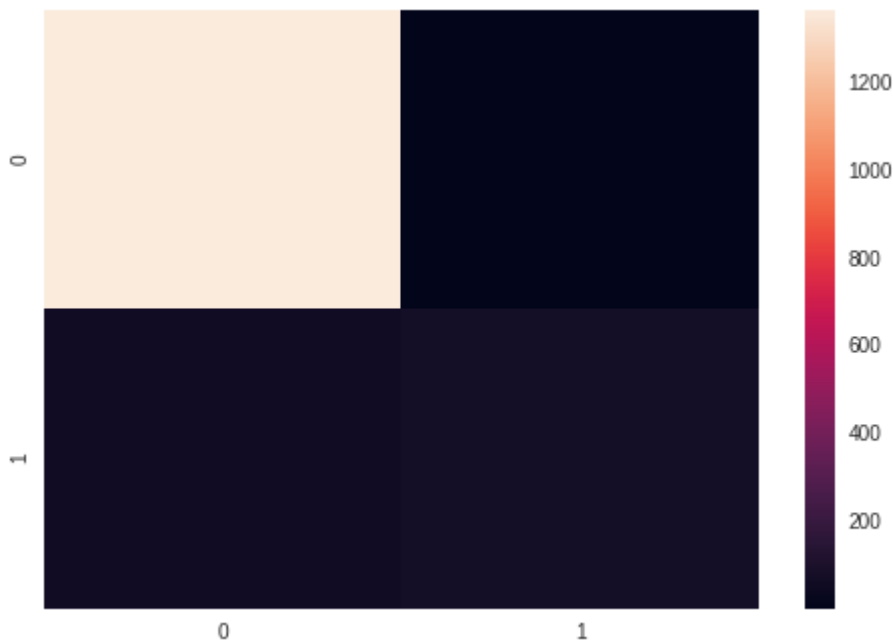
▼ SVM Based Classifier

```
1 from sklearn.svm import SVC
2 Spvm = SVC(kernel = "poly",class_weight=dic)
3 Spvm.fit(X_train,Y_train)
4 Spvm.score(X_test,Y_test)
```

☞ 0.9566666666666667

```
1 Y_pred = Spvm.predict(X_test)
2 from sklearn.metrics import confusion_matrix
3 sns.heatmap(confusion_matrix(Y_test,Y_pred))
4 f1_score(Y_test,Y_pred)
```

☞ 0.6948356807511737



▼ Decision Tree Classifier

As predicted we get a lot better prediction scores.

```
1 from sklearn.tree import DecisionTreeClassifier
2 Dt = DecisionTreeClassifier()
3 Dt.fit(X_train,Y_train)
```

```

3 print(Dt.score(X_test,Y_test))
4 print(Dt.score(X_test,Y_test))
5 Y_pred = Dt.predict(X_test)
6 f1_score(Y_test,Y_pred)

```

```

0.9833333333333333
0.9110320284697508

```

```

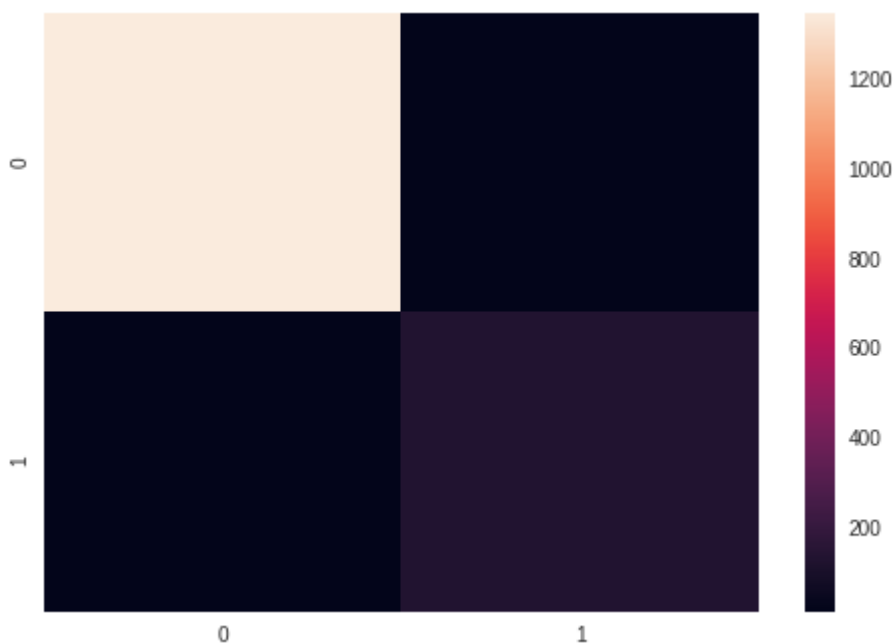
1 Y_pred = Dt.predict(X_test)
2 from sklearn.metrics import confusion_matrix
3 print(confusion_matrix(Y_test,Y_pred))
4 sns.heatmap(confusion_matrix(Y_test,Y_pred))

```

```

[[1347  15]
 [ 10 128]]
<matplotlib.axes._subplots.AxesSubplot at 0x7f3253c10438>

```



```

1 from sklearn import tree
2 import graphviz
3 dt_data = tree.export_graphviz(Dt, out_file=None)
4 graph = graphviz.Source(dt_data)
5 graph.render("iris")
6 tree.export_graphviz(Dt, out_file=None,
7                       feature_names=x.columns,
8                       class_names = "Personal Loan",
9                       special_characters=True)

```

```

'digraph Tree {\nnode [shape=box, style="filled, rounded", color="black", fontname=helvetica, fontname=helvetica] ;\n0 [label=<Income &le; 100.5<br/>gini = 0.176<br/>samples = 3500<br/>158, 342<br/>class = P>, fillcolor="#e88f4e"] ;\n1 [label=<CCAvg &le; 2.95<br/>gini = 0.0<br/>samples = 2638<br/>value = [2607, 31]<br/>class = P>, fillcolor="#e5823b"] ;\n0 -> 1 [label=<belong to class = P>, fillcolor="#e58139"] ;\n1 -> 2 ;\n3 [label=<CD Account &le; 0.5<br/>gini = 0.0<br/>samples = 166<br/>value = [135, 31]<br/>class = P>, fillcolor="#e58139"] ;\n1 -> 3 ;\n4 [label=<

```

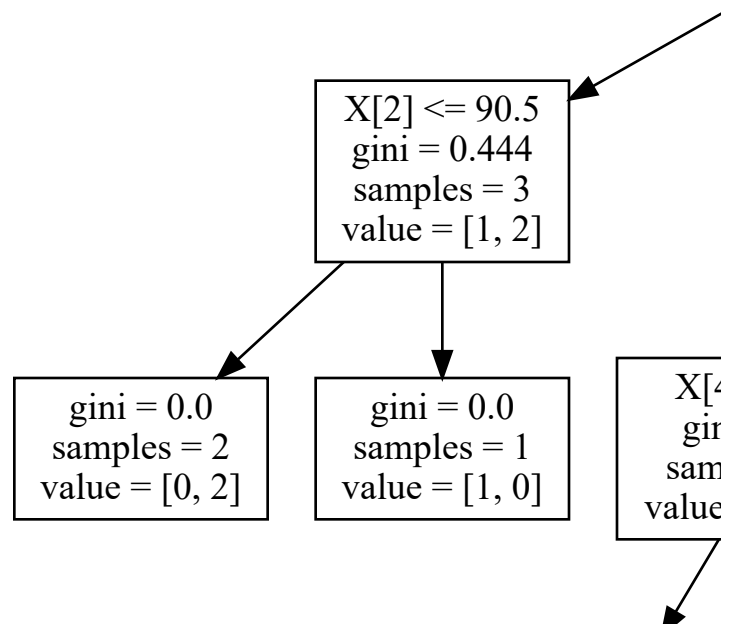
```
1 x.columns
```

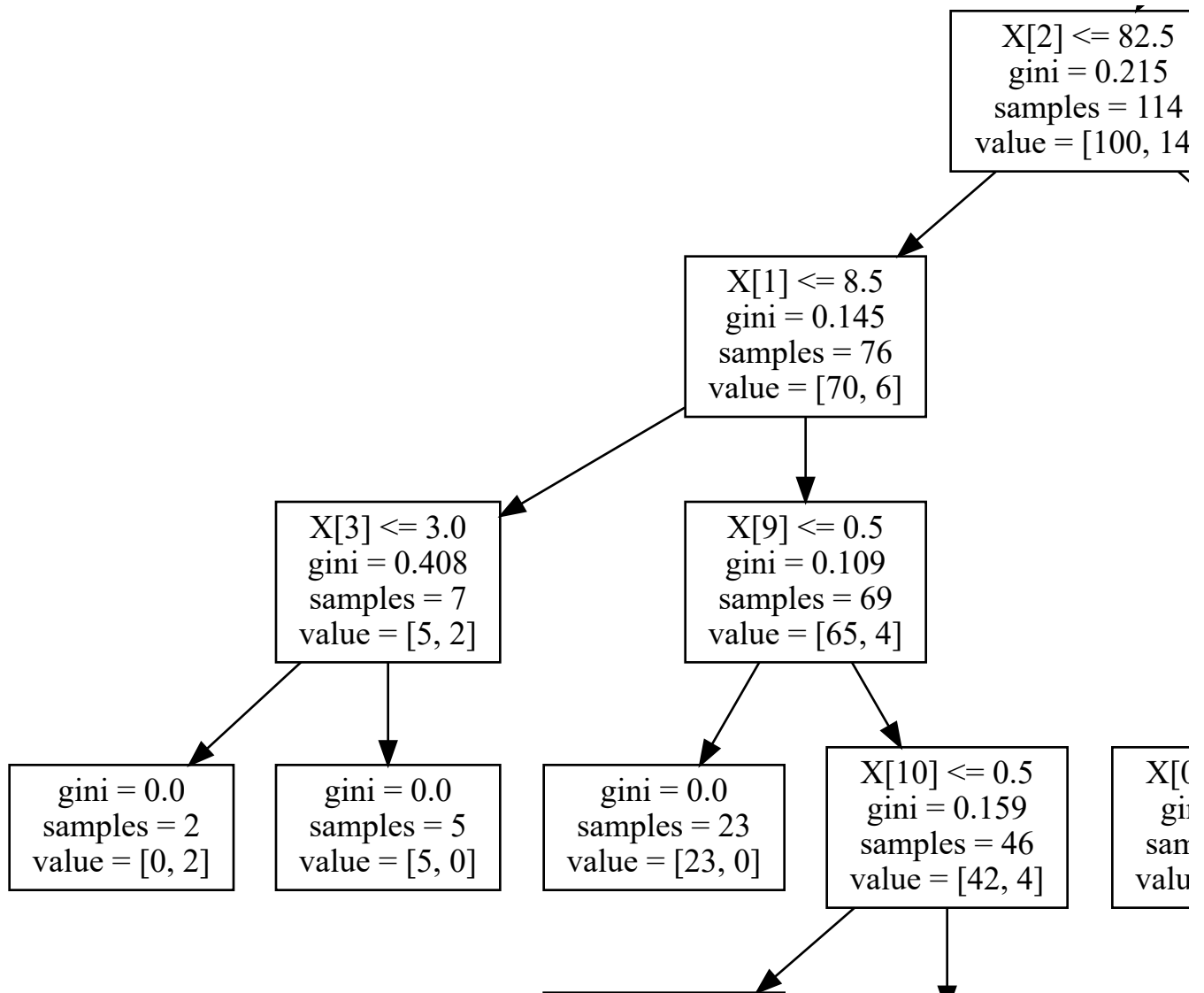
```
↳ Index(['Age', 'Experience', 'Income', 'Family', 'CCAvg', 'Education',  
         'Mortgage', 'Securities Account', 'CD Account', 'Online', 'CreditCard'],  
        dtype='object')
```

```
1 graph = graphviz.Source(dt_data)
```

```
2 graph
```

```
↳
```





RandomForestClassifier

```

1 from sklearn.ensemble import RandomForestClassifier
2 n_est = [25,50,75,100,125]
3 max_depth = [None,3,4,5,6]
4 acc = []
5 for i in n_est:
6     for k in max_depth:
7         rbf_clf = RandomForestClassifier(n_estimators=i,max_depth=k)
8         rbf_clf.fit(X_train, Y_train)
9         Y_pred = rbf_clf.predict(X_test)
10        print(rbf_clf.score(X_test, Y_test),f1_score(Y_test,Y_pred), "N_est = " +str(i),"m

```



```

0.9886666666666667 0.9368029739776952 N_est = 25 max_depth = None
0.9606666666666667 0.730593607305936 N_est = 25 max_depth = 3
0.9793333333333333 0.8734693877551021 N_est = 25 max_depth = 4
0.984 0.9069767441860466 N_est = 25 max_depth = 5
0.982 0.8957528957528957 N_est = 25 max_depth = 6
0.9886666666666667 0.9368029739776952 N_est = 50 max_depth = None
0.9573333333333334 0.7009345794392524 N_est = 50 max_depth = 3
0.972 0.8220338983050847 N_est = 50 max_depth = 4
0.9833333333333333 0.9034749034749036 N_est = 50 max_depth = 5
0.9853333333333333 0.9166666666666666 N_est = 50 max_depth = 6
0.986 0.9207547169811321 N_est = 75 max_depth = None
0.9433333333333334 0.5595854922279793 N_est = 75 max_depth = 3
0.964 0.7589285714285714 N_est = 75 max_depth = 4
0.9793333333333333 0.8774703557312253 N_est = 75 max_depth = 5
0.9846666666666667 0.9118773946360155 N_est = 75 max_depth = 6
0.9866666666666667 0.9259259259259259 N_est = 100 max_depth = None
0.944 0.5670103092783505 N_est = 100 max_depth = 3
0.9733333333333334 0.8347107438016529 N_est = 100 max_depth = 4
0.9906666666666667 0.9477611940298507 N_est = 100 max_depth = 5

```

```

1 rbf_clf = RandomForestClassifier(n_estimators=100,max_depth = None)
2 rbf_clf.fit(X_train, Y_train)
3 Y_pred = rbf_clf.predict(X_test)
4 print(rbf_clf.score(X_test, Y_test),f1_score(Y_test,Y_pred))

```

```
0.9906666666666667 0.9477611940298507
```

▼ Best F1 Score 94 % Accuracy 99

Since the accuracy varies a lot so let us boost it to its maximum accuracy.(I'll use XGBoost)

```

1 import pickle
2 pickle.dump(rbf_clf,open("rbf_clf.pkl","wb"))

```

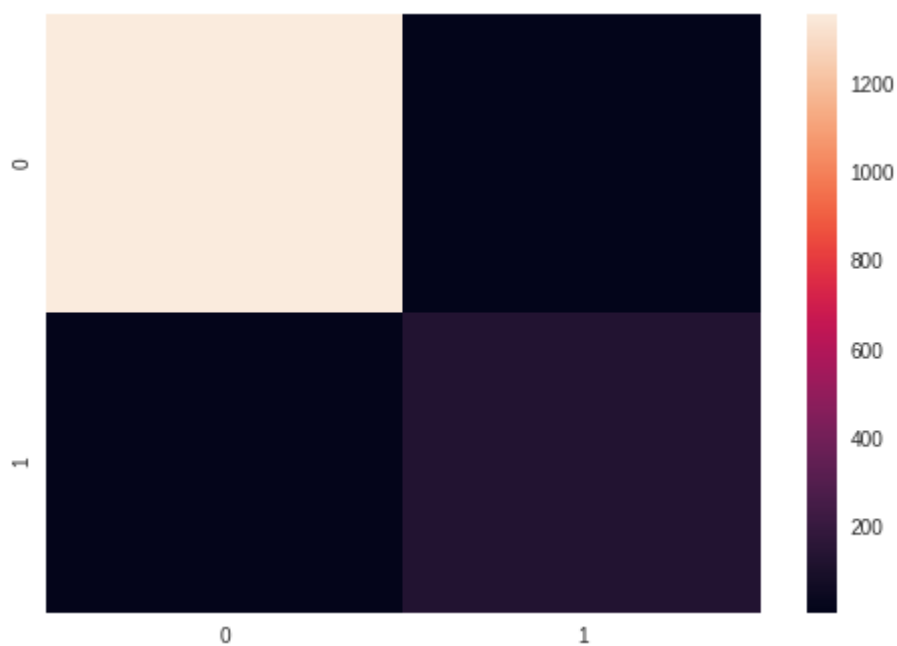
```
1 rbf_clf.get_params()
```

```
0.9906666666666667 0.9477611940298507
```

```
{'bootstrap': True,  
  'ccp_alpha': 0.0,  
  'class_weight': None,
```

```
1 Y_pred = rbf_clf.predict(X_test)  
2 print(confusion_matrix(Y_test,Y_pred))  
3 sns.heatmap(confusion_matrix(Y_test,Y_pred))  
4 f1_score(Y_test,Y_pred)
```

```
[[1359   3]  
 [  11 127]]  
0.9477611940298507
```



Measure	Value
<hr/>	
1 X_train.columns	
↳ Index(['Age', 'Experience', 'Income', 'Family', 'CCAvg', 'Education', 'Mortgage', 'Securities Account', 'CD Account', 'Online', 'CreditCard'], dtype='object')	
Negative Predictive Value	0.9930
1 # data.groupby(['Loan_Status', 'Gender', 'Property_Area'])['Married'].value_counts()	
2 # from yellowbrick.classifier import ClassificationReport	
3 # viz = ClassificationReport(LR_model)	
4 # viz.fit(XT, YT)	
5 # viz.score(Xt, Yt)	
6 # viz.show()	
7 # import statsmodels.api as sm	
8 # X2 = sm.add_constant(XT)	
9 # regressor2 = sm.Logit(YT, X2).fit()	
10 # print("p-Values for each column: ")	
11 # print()	
12 # print(regressor2.pvalues)	
1 from xgboost import XGBClassifier	
2 model = XGBClassifier(max_depth=3,n_estimators=100)	
3 model.fit(X_train, Y_train)	
↳ XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1, colsample_bynode=1, colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=3, min_child_weight=1, missing=None, n_estimators=100, n_jobs=1, nthread=None, objective='binary:logistic', random_state=0, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None, silent=None, subsample=1, verbosity=1)	
1 print(model.score(X_test,Y_test),model.score(X_train,Y_train))	
↳ 0.9893333333333333 0.9908571428571429	
1 Y_pred = model.predict(X_test)	
1 print(confusion_matrix(Y_test,Y_pred))	
2 sns.heatmap(confusion_matrix(Y_test,Y_pred))	
3 f1_score(Y_test,Y_pred)	
↳	