

### PIMPRI CHINCHWAD COLLEGE OF ENGINEERING



DESCRIPTION

TITLE- BANK CUSTOMERS EXITING PREDICTIVE MODEL

**GUIDE: PROF. PRIYA SURANA** 

CLASSIFICATION

VARUN GADDE - BECOB203

VEDANT UPGANLAWAR - BECOB<sub>270</sub>

PREPROCESSING

1. DESCRIPTION

**1.A PROBLEM STATEMENT** 

Consider a labeled dataset belonging to an application domain. Apply suitable data preprocessing steps such as handling of null values, data reduction, discretization. For prediction of class labels of given data instances, build classifier models using different techniques (minimum 3), analyze the confusion matrix and compare these models. Also apply cross validation while preparing the training and testing datasets.

1.B DATASET DESCRIPTION A bank is investigating a very high rate of customer leaving the bank. Here is a 10,000 records dataset to investigate and predict which of

the customers are more likely to leave the bank soon.

We will be guiding you with the implementation part simultaneosuly. Current project can be run on any Python, Jupyter notebook, Google Colab etc.

The project requires the following imports:

2. Pandas - for Data Preprocessing and CSV I/O

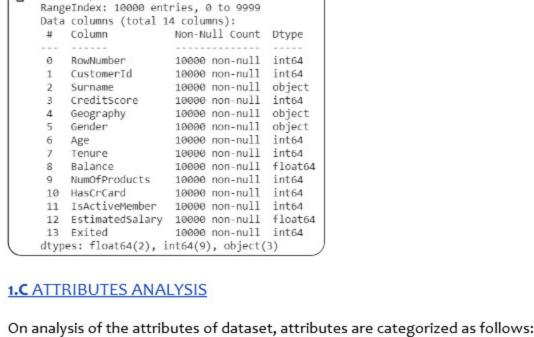
1. Numpy - for Linear Algebra

- 3. Matplotlib Data Visualization 4. sklearn.modelselection - for Modelling
- 5. sklearn.linear\_model for Logistic Regression Classifier 6. sklearn.svm - for Support Vector Classifier
- 7. sklearn.naive\_bayes for Gaussian Naive Bayes Classifier 8. sklearn.neighbors - for KNeighbors Classifier
- sklearn.ensemble for Random Forest Classifier 10. sklearn.metrics - for Accuracy Score, Confusion Matrix and Classification Report

As and when the need aries the imports are executed in the project.

Title of the dataset that is to be anaysed throughout the project is 'Churn\_Modelling.csv' which is manually uploaded to the execution environment. We loaded the dataset into a variable 'data'. Following is the description of all the attributes of the dataframe -

data.info() <class 'pandas.core.frame.DataFrame'>



#### Categorical Variables Geography

Gender HasCrCard

```
IsActiveMember
          Exited
          NumOfProducts
  2. Numerical Variables
          CreditScore
          Age
          Tenure
          Balance
         EstimatedSalary
From the following features, 'Exited' is a categorical feature which determines whether the customer left the bank or not.
```

2.A MISSING VALUES

# We can check missing value also by isnull() method as below:

2. PREPROCESSING

data.isnull().any()

RowNumber False E. CustomerId False Surname False

```
CreditScore
                      False
    Geography
                       False
    Gender
                      False
                       False
    Tenure
                       False
    Balance
                      False
    NumOfProducts
                       False
    HasCrCard
                      False
    IsActiveMember
                      False
    EstimatedSalary
                      False
    Exited
                      False
    dtype: bool
From above, it can be concluded that there are no missing value in any of the columns. There are 10000 items and each column has 10000
non-null items.
2.B DROPPING UNNECESSARY COLUMNS
```

Balance NumOfProducts HasCrCard IsActiveMember EstimatedSalary Exited CreditScore Geography Gender Age Tenure

```
619
                        France Female
                                                         0.00
                                                                                                            101348.88
                                                                                                            112542.58
                608
                         Spain Female
                                                    83807.86
                                                                                                                           0
                         France Female
                502
                                                 8 159660.80
                                                                                                            113931.57
                        France Female
                                                                                                             93826.63
                         Spain Female
                                                 2 125510.82
Here we can clearly observe that columns'Rownumber', 'CustomerId' and 'Surname' does not provide any meaningful information. Hence
we are dropping the above mentioned columns.
2.C DISCRETIZATION
```

 Geography Geography\_France Geography\_Germany Geography\_Spain

1

0

0

0

0

0

## 0

0 0 0

0

Now we will convert the following attributes into categorical values.

data.drop(["RowNumber","CustomerId","Surname"], axis=1, inplace=True)

	'	U	0
	0	0	1
2.	Gender Gender_Female Gen	der_Male	
	1	0	
	1	0	
	1	0	
	1	0	

0 0

0

0

0

0

1

1

0

0

0

0

3. Tenure

NumOfProducts_1	NumOfProducts_2	NumOfProducts_3	NumOfProducts_4
1	0	0	0
1	0	0	0
0	0	1	0
0	1	0	0
1	0	0	0

4. Splitting the X\_train and y\_train as X\_Train,y\_Train,X\_Test,y\_test. The split ratio varies according to the discretion of the user. Here we

Tenure\_2 Tenure\_3 Tenure\_4 Tenure\_5 Tenure\_6 Tenure\_7 Tenure\_8 Tenure\_9 Tenure\_10

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

0

### have used the splitting ratio as 0.33. Here we have executed some instructions for better understanding of the splitted data shape.

3. CLASSIFICATION

3.A MODELLING

#creating X\_train and y\_train X\_train = df.drop("Exited", axis = 1) y\_train = df["Exited"] #Features (CreditScore, Age, Balance, EstimatedSalary) to be normalized:

Modelling includes creation of train, test data by splitting the dataset. The steps followed below are:

 Create X\_train by dropping the dependent attribute "Exited" from the dataframe df. 2. Create y\_train by considering the dependent attribute "Exited" from the dataframe df.

Normalizing necessary attributes (CreditScore, Age, Balance, EstimatedSalary)

for each in ["CreditScore", "Age", "Balance", "EstimatedSalary"]:

normalize\_feature(each, X\_train)

#X\_train and Y\_train data shape summary print("Length of X\_train: ",len(X\_train)) print("Shape of X\_train: ", X\_train.shape) print("Length of Y\_tain: ", len(y\_train)) print("Shape of Y\_train: ", y\_train.shape)

Length of X\_train: 10000

3.B LOGISTIC REGRESSION

3.C SVM

3.E KNN

3.D NAIVE BAYES

```
Shape of X_train: (10000, 26)
    Length of Y_tain: 10000
    Shape of Y_train: (10000,)
from sklearn.model_selection import train_test_split
    X_train, X_test, y_train, y_test = train_test_split(
        X train,
        y_train,
        test_size = 0.33,
        random_state = 42
    print("Length of X_train: ",len(X_train))
    print("Length of X_test: ",len(X_test))
    print("Length of y_train: ",len(y_train))
    print("Length of y_test: ",len(y_test))
    Length of X_train: 6700
    Length of X test: 3300
    Length of y_train: 6700
    Length of y_test: 3300
```

3.F RANDOM FOREST <a>
▼</a> CLASSIFY

cm=confusion\_matrix(y\_test,y\_pred)

print(classification\_report(y\_test, y\_pred))

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import classification\_report

print("RF accuracy with test data :", rf.score(X\_test, y\_test)\*100)

from sklearn.metrics import accuracy\_score from sklearn.metrics import confusion\_matrix

#### rf = RandomForestClassifier(n\_estimators = 24, random\_state = 42) rf.fit(X train, y train) y\_pred=rf.predict(X\_test)

print() print(cm) print()

RF accuracy with test data : 86.51515151515152 [[2554 103] [ 342 301]] precision 0 0.88 1 0.75 accuracy macro avg 0.81 weighted avg 0.86

print("accuracies :",accuracies)

print("mean accuracy :", accuracies.mean())

recall f1-score support 0.96 0.92 2657 0.47 0.57 643 0.87 3300 0.71 0.75 3300 0.87 0.85 3300 4. METRICS 4.A K-FOLD CROSS VALIDATION In this case of cross validation, we will divide the complete dataset into k parts. Each section will act as the test set and rest of the k-1 sections will be used for training. Hence we get k different accuracy scores. For this dataset, we have k=10. Finally from these k scores, we obtain a mean accuracy which helps us compare the classifiers with better precision. Below is implementation of k-fold cross validation on Random Forest Classifier from sklearn.model selection import cross val score

accuracies = cross\_val\_score(estimator =rf, X = X\_train, y = y\_train, cv = 10)

eir respective mean

urther classifier's summary is given in the below table (Following table is arranged in descending order curacies) -										
Classifier	Accuracy ( Test Data )	K-Fold Cross Validation Accuracy	Confusion Matrix	Classification Report						
Random Forest ( RF )	86.51	85.31	[[2554 103] [ 342 301]]	0	precision 0.88 0.75	0.96 0.47	f1-score 0.92 0.57	support 2657 643		
Logistic Regression (LR)	85.0	83.61	[[2552 105] [ 399 244]]	0	precision 0.86 0.70	0.96 0.38	f1-score 0.91 0.49	support 2657 643		
Support Vector ( SVM )	84.66	83.46	[[2602 55] [ 451 192]]	0	precision 0.85 0.78	0.98 0.30	f1-score 0.91 0.43	support 2657 643		
Naive Bayes ( NB )	82.78	81.56	[[2638 19] [ 549 94]]	0	precision 0.83 0.83	0.99 0.15	f1-score 0.90 0.25	support 2657 643		
K Nearest Neighbour ( KNN )	82.0	80.20	[[2559 98] [ 502 141]]	0	precision 0.84 0.59	0.96 0.22	f1-score 0.90 0.32	support 2657 643		



**METRICS**