

```
[1]: import numpy as np
import pandas as pd

In [2]: train_data = pd.read_csv('Churn_Modelling.csv')

In [3]: train_data.shape
Out[3]: (10086, 14)

In [4]: train_data
Out[4]:
   RowNumber  CustomerId  Surname  CreditScore  Geography  Gender  Age  Tenure  Balance  NumOfProducts  HasCrCard  IsActiveMember  EstimatedSalary  Exited
0           0           1  15634002  Hargrave         619      France  Female  42     2     0.00             1             1             1     101348.86      1
1           1           2  15647311         H8         608      Spain  Female  41     1     83807.86            1             0             1     112542.56      0
2           2           3  15618304         Ono         502      France  Female  42     8     159660.80            3             1             0     113931.57      1
3           3           4  15701254         Bore         699      France  Female  39     1     0.00             2             0             0      93826.63      0
4           4           5  15737888         Mitchell         850      Spain  Female  43     2     125510.82            1             1             1      79084.10      0
...
9995          9995  15603220         Ohjalu         771      France  Male   39     5     0.00             2             1             0      96270.64      0
9996          9997  15608823         Johnson         516      France  Male   35    10     57369.61            1             1             1     101099.77      0
9997          9998  15614032         Ili         709      France  Female  36     7     0.00             1             0             1      42085.56      1
9998          9999  15683295         Substane         772      Germany  Male  42     3     75075.31            2             1             0      92888.52      1
9999         10000  15628319         Walker         792      France  Female  28     4     130142.79            1             1             0      38190.78      0
10000 rows x 14 columns

In [5]: train_data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10086 entries, 0 to 9999
Data columns (total 14 columns):
#   Column                Non-Null Count  Dtype
---  --
0   RowNumber              10086 non-null    int64
1   CustomerId             10086 non-null    int64
2   Surname                10086 non-null    object
3   CreditScore             10086 non-null    int64
4   Geography              10086 non-null    object
5   Gender                 10086 non-null    object
6   Age                    10086 non-null    int64
7   Tenure                 10086 non-null    int64
8   Balance                 10086 non-null    float64
9   NumOfProducts          10086 non-null    int64
10  HasCrCard               10086 non-null    int64
11  IsActiveMember          10086 non-null    int64
12  EstimatedSalary         10086 non-null    float64
13  Exited                  10086 non-null    int64
dtypes: float64(2), int64(8), object(5)
memory usage: 1.1+ MB

In [6]: train_data.describe()
Out[6]:
   RowNumber  CustomerId  CreditScore  Age  Tenure  Balance  NumOfProducts  HasCrCard  IsActiveMember  EstimatedSalary  Exited
count  10000.000000  1.000000e+04  10000.000000  10000.000000  10000.000000  10000.000000  10000.000000  10000.000000  10000.000000  10000.000000
mean      2992.000000  1.560904e+04  650.528800  38.921800  5.012800  76485.889288  1.532000  0.70550  0.515100  100090.239881  0.203700
std      2886.895658  7.193619e+04  86.653290  10.487806  2.892174  62397.405202  0.581654  0.45584  0.499797  57510.492018  0.402769
min       1.000000  1.556570e+07  350.000000  18.000000  0.000000  0.000000  1.000000  0.00000  0.000000  1158000.000000  0.000000
25%      2500.750000  1.560234e+07  584.000000  32.000000  3.000000  0.000000  1.000000  0.00000  0.000000  51052.110000  0.000000
50%      5000.500000  1.560724e+07  652.000000  37.000000  5.000000  97196.540000  1.000000  1.00000  1.000000  100193.910000  0.000000
75%      7500.250000  1.575322e+07  718.000000  44.000000  7.000000  127644.240000  2.000000  1.00000  1.000000  149388.247000  0.000000
max      10000.000000  1.581569e+07  850.000000  10.000000  250895.090000  4.000000  1.00000  1.00000  1.000000  199992.480000  1.000000

In [7]: train_data.describe(include = 'object')
Out[7]:
   Surname  Geography  Gender
count  10000  10000  10000
unique    2992         5         2
freq      5014  5014  5457

In [8]: train_data.isnull().sum()
Out[8]:
RowNumber      0
CustomerId      0
Surname         0
CreditScore     0
Geography       0
Gender          0
Age            0
Tenure         0
Balance        0
NumOfProducts  0
HasCrCard      0
IsActiveMember  0
EstimatedSalary 0
Exited         0
dtypes: int64

In [9]: train_data.dropna(ignore_index=True, inplace=True)

In [10]: data = train_data.drop(columns=['RowNumber', 'CustomerId', 'Surname'])
Out[10]: data
Out[11]:
   CreditScore  Geography  Gender  Age  Tenure  Balance  NumOfProducts  HasCrCard  IsActiveMember  EstimatedSalary  Exited
0           619      France  Female  42     2     0.00             1             1             1     101348.86      1
1           608      Spain  Female  41     1     83807.86            1             0             1     112542.56      0
2           502      France  Female  42     8     159660.80            3             1             0     113931.57      1
3           699      France  Female  39     1     0.00             2             0             0      93826.63      0
4           850      Spain  Female  43     2     125510.82            1             1             1      79084.10      0
...
9995          771      France  Male   39     5     0.00             2             1             0      96270.64      0
9996          516      France  Male   35    10     57369.61            1             1             1     101099.77      0
9997          709      France  Female  36     7     0.00             1             0             1      42085.56      1
9998          772      Germany  Male  42     3     75075.31            2             1             0      92888.52      1
9999          792      France  Female  28     4     130142.79            1             1             0      38190.78      0
10000 rows x 11 columns

In [12]: data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10086 entries, 0 to 9999
Data columns (total 11 columns):
#   Column                Non-Null Count  Dtype
---  --
0   CreditScore            10086 non-null    int64
1   Geography              10086 non-null    object
2   Gender                 10086 non-null    object
3   Age                    10086 non-null    int64
4   Tenure                 10086 non-null    int64
5   Balance                 10086 non-null    float64
6   NumOfProducts          10086 non-null    int64
7   HasCrCard              10086 non-null    int64
8   IsActiveMember          10086 non-null    int64
9   EstimatedSalary         10086 non-null    float64
10  Exited                  10086 non-null    int64
dtypes: float64(2), int64(7), object(2)
memory usage: 859.3+ KB

In [13]: data
Out[13]:
   CreditScore  Geography  Gender  Age  Tenure  Balance  NumOfProducts  HasCrCard  IsActiveMember  EstimatedSalary  Exited
0           619      France  Female  42     2     0.00             1             1             1     101348.86      1
1           608      Spain  Female  41     1     83807.86            1             0             1     112542.56      0
2           502      France  Female  42     8     159660.80            3             1             0     113931.57      1
3           699      France  Female  39     1     0.00             2             0             0      93826.63      0
4           850      Spain  Female  43     2     125510.82            1             1             1      79084.10      0
...
9995          771      France  Male   39     5     0.00             2             1             0      96270.64      0
9996          516      France  Male   35    10     57369.61            1             1             1     101099.77      0
9997          709      France  Female  36     7     0.00             1             0             1      42085.56      1
9998          772      Germany  Male  42     3     75075.31            2             1             0      92888.52      1
9999          792      France  Female  28     4     130142.79            1             1             0      38190.78      0
10000 rows x 11 columns

In [14]: ax = data['Exited'].value_counts().to_frame().plot(kind='bar', legend=False)
ax.set_xlabel('Exited')
ax.set_ylabel('Count')
Out[14]: Text(0, 0.5, 'Count')


In [15]: import seaborn as sns
c = data['Geography'].value_counts().to_frame()
c.reset_index()
Out[15]:
   Geography  count
0      France   5014
1         Spain   2509
2         Spain   2477

In [16]: c = data.groupby(['Exited'])['Geography'].value_counts().reset_index()
Out[16]:
   Exited  Geography  count
0         0      France  4204
1         0         Spain  2064
2         0      Germany  1695
3         1      Germany   814
4         1         France   810
5         1         Spain   413

In [17]: import seaborn as sns
sns.countplot(x='Geography', data=data, hue='Exited')
Out[17]:
<Axes: xlabel='Geography', ylabel='count'>


In [18]: sns.barplot(x='Geography', y='count', data=data, hue='Exited')
Out[18]:
<Axes: xlabel='Geography', ylabel='count'>


In [19]: sns.countplot(x='Gender', data=data, hue='Exited')
Out[19]:
<Axes: xlabel='Gender', ylabel='count'>


In [20]: sns.countplot(x='NumOfProducts', data=data, hue='Exited')
Out[20]:
<Axes: xlabel='NumOfProducts', ylabel='count'>


In [21]: sns.countplot(x='HasCrCard', data=data, hue='Exited')
Out[21]:
<Axes: xlabel='HasCrCard', ylabel='count'>


In [22]: sns.countplot(x='IsActiveMember', data=data, hue='Exited')
Out[22]:
<Axes: xlabel='IsActiveMember', ylabel='count'>


In [23]: data['Age'].value_counts()
Out[23]:
Age
37     478
38     477
35     474
36     456
34     447
...
92         2
82         1
88         1
85         1
83         1
Name: count, Length: 70, dtype: int64

In [24]: sns.kdeplot(x=data['Age'], bins=30, data=data, kde=True, hue='Exited')
Out[24]:
<Axes: xlabel='Age', ylabel='count'>


In [25]: sns.histplot(x=data['Balance'], bins=30, data=data, kde=True, hue='Exited')
Out[25]:
<Axes: xlabel='Balance', ylabel='Count'>


In [26]: data['Tenure'].value_counts()
Out[26]:
Tenure
2     1848
1     1635
7     1528
8     1525
5     1499
...
9         98
6         984
4         989
3         994
2         999
Name: count, dtype: int64

In [27]: sns.countplot(x='Tenure', data=data, hue='Exited')
Out[27]:
<Axes: xlabel='Tenure', ylabel='count'>


In [28]: data.dtypes
Out[28]:
CreditScore    int64
Geography      object
Gender         object
Age            int64
Tenure         int64
Balance        float64
NumOfProducts  int64
HasCrCard      int64
IsActiveMember int64
EstimatedSalary float64
Exited         object
dtype: object

In [29]: data['HasCrCard'] = data['HasCrCard'].astype('object')

In [30]: data['IsActiveMember'] = data['IsActiveMember'].astype('object')
data.dtypes
Out[30]:
CreditScore    int64
Geography      object
Gender         object
Age            int64
Tenure         int64
Balance        float64
NumOfProducts  int64
HasCrCard      object
IsActiveMember object
EstimatedSalary float64
Exited         object
dtype: object

In [31]: data['Exited'] = data['Exited'].astype('object')
data.dtypes
Out[31]:
CreditScore    int64
Geography      object
Gender         object
Age            int64
Tenure         int64
Balance        float64
NumOfProducts  int64
HasCrCard      object
IsActiveMember object
EstimatedSalary float64
Exited         object
dtype: object

In [32]: from sklearn.preprocessing import StandardScaler
data['CreditScore'] = data['CreditScore'].astype('int')
data['Geography_Spain'] = le.fit_transform(data['Geography_Spain'])
data['Gender_Male'] = le.fit_transform(data['Gender_Male'])
data
Out[32]:
   CreditScore  Geography  Gender  Age  Tenure  Balance  NumOfProducts  HasCrCard  IsActiveMember  EstimatedSalary  Exited  Geography_Germany  Geography_Spain  Gender_Male
0   -0.326221  0.295117  2   -1.225848            1             1             0.021886      1             0.021886      0             0             0             0
1   -0.440036  0.196164            1  0.117350            1             0             0.216534      0             0.216534      0             0             1             0
2   -1.536794  0.295117  8   1.330553            3             1             0.240687      1             0.240687      1             0             0             0
3   0.501521  0.007457  1   -1.225848            2             0             -0.108918      0             -0.108918      0             0             0             0
4   2.063884  0.388871  2   0.785728            1             1             -0.365276            0             -0.365276            0             0             1             0
...
9995  1.245488  0.007457  5   -1.225848            2             1             -0.065419            0             -0.065419            0             0             0             1
9996  -1.391339  -0.373958    10  -0.306379            1             1             0.027988            0             0.027988            0             0             0             1
9997  0.604988  -0.278054  7   -1.225848            1             0             -1.009643            1             -1.009643            1             0             1             0
9998  1.256835  0.295117  3   -0.022608            2             1             0.012523            1             0.012523            1             1             0             1
9999  1.463771  1.041433  4   0.859965            1             1             -0.107670            0             -0.107670            0             0             0             0
10000 rows x 12 columns

In [33]: data = pd.get_dummies(data, columns=['Geography', 'Gender'], drop_first=True)
data
Out[33]:
   CreditScore  Age  Tenure  Balance  NumOfProducts  HasCrCard  IsActiveMember  EstimatedSalary  Exited  Geography_Germany  Geography_Spain  Gender_Male
0   -0.326221  0.295117  2   -1.225848            1             1             0.021886      1             0.021886      0             0             0
1   -0.440036  0.196164            1  0.117350            1             0             0.216534      0             0.216534      0             1             0
2   -1.536794  0.295117  8   1.330553            3             1             0.240687      1             0.240687      1             0             0
3   0.501521  0.007457  1   -1.225848            2             0             -0.108918      0             -0.108918      0             0             0
4   2.063884  0.388871  2   0.785728            1             1             -0.365276            0             -0.365276            0             1             0
...
9995  1.245488  0.007457  5   -1.225848            2             1             -0.065419            0             -0.065419      0             0             1
9996  -1.391339  -0.373958    10  -0.306379            1             1             0.027988            0             0.027988      0             0             1
9997  0.604988  -0.278054  7   -1.225848            1             0             -1.009643            1             -1.009643      1             0             0
9998  1.256835  0.295117  3   -0.022608            2             1             0.012523            1             0.012523            1             1             0
9999  1.463771  1.041433  4   0.859965            1             1             -0.107670            0             -0.107670      0             0             0
10000 rows x 12 columns

In [34]: from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
data['Geography_Germany'] = le.fit_transform(data['Geography_Germany'])
data['Geography_Spain'] = le.fit_transform(data['Geography_Spain'])
data['Gender_Male'] = le.fit_transform(data['Gender_Male'])
data
Out[34]:
   CreditScore  Age  Tenure  Balance  NumOfProducts  HasCrCard  IsActiveMember  EstimatedSalary  Exited  Geography_Germany  Geography_Spain  Gender_Male
0   -0.326221  0.295117  2   -1.225848            1             1             0.021886      1             0.021886      0             0             0
1   -0.440036  0.196164            1  0.117350            1             0             0.216534      0             0.216534      0             1             0
2   -1.536794  0.295117  8   1.330553            3             1             0.240687      1             0.240687      1             0             0
3   0.501521  0.007457  1   -1.225848            2             0             -0.108918      0             -0.108918      0             0             0
4   2.063884  0.388871  2   0.785728            1             1             -0.365276            0             -0.365276      0             1             0
...
9995  1.245488  0.007457  5   -1.225848            2             1             -0.065419            0             -0.065419      0             0             1
9996  -1.391339  -0.373958    10  -0.306379            1             1             0.027988            0             0.027988      0             0             1
9997  0.604988  -0.278054  7   -1.225848            1             0             -1.009643            1             -1.009643      1             0             0
9998  1.256835  0.295117  3   -0.022608            2             1             0.012523            1             0.012523            1             1             0
9999  1.463771  1.041433  4   0.859965            1             1             -0.107670            0             -0.107670      0             0             0
10000 rows x 12 columns

In [35]: data.columns
Out[35]:
Index(['CreditScore', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'HasCrCard', 'IsActiveMember', 'EstimatedSalary', 'Exited', 'Geography_Germany', 'Geography_Spain', 'Gender_Male',
      'Geography_Spain', 'Gender_Male'],
      dtype='object')

In [36]: data['HasCrCard'] = data['HasCrCard'].astype('int')
data['IsActiveMember'] = data['IsActiveMember'].astype('int')
data['Exited'] = data['Exited'].astype('int')
data
Out[36]:
   CreditScore  Age  Tenure  Balance  NumOfProducts  HasCrCard  IsActiveMember  EstimatedSalary  Exited  Geography_Germany  Geography_Spain  Gender_Male
0   -0.326221  0.295117  2   -1.225848            1             1             0.021886      1             0.021886      0             0             0
1   -0.440036  0.196164            1  0.117350            1             0             0.216534      0             0.216534      0             1             0
2   -1.536794  0.295117  8   1.330553            3             1             0.240687      1             0.240687      1             0             0
3   0.501521  0.007457  1   -1.225848            2             0             -0.108918      0             -0.108918      0             0             0
4   2.063884  0.388871  2   0.785728            1             1             -0.365276            0             -0.365276      0             1             0
...
9995  1.245488  0.007457  5   -1.225848            2             1             -0.065419            0             -0.065419      0             0             1
9996  -1.391339  -0.373958    10  -0.306379            1             1             0.027988            0             0.027988      0             0             1
9997  0.604988  -0.278054  7   -1.225848            1             0             -1.009643            1             -1.009643      1             0             0
9998  1.256835  0.295117  3   -0.022608            2             1             0.012523            1             0.012523            1             1             0
9999  1.463771  1.041433  4   0.859965            1             1             -0.107670            0             -0.107670      0             0             0
10000 rows x 12 columns

In [37]: data.columns
Out[37]:
Index(['CreditScore', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'HasCrCard', 'IsActiveMember', 'EstimatedSalary', 'Exited', 'Geography_Germany', 'Geography_Spain', 'Gender_Male',
      'Geography_Spain', 'Gender_Male'],
      dtype='object')

In [38]: data['HasCrCard'] = data['HasCrCard'].astype('int')
data['IsActiveMember'] = data['IsActiveMember'].astype('int')
data['Exited'] = data['Exited'].astype('int')
data
Out[38]:
   CreditScore  Age  Tenure  Balance  NumOfProducts  HasCrCard  IsActiveMember  EstimatedSalary  Exited  Geography_Germany  Geography_Spain  Gender_Male
0   -0.326221  0.295117  2   -1.225848            1             1             0.021886      1             0.021886      0             0             0
1   -0.440036  0.196164            1  0.117350            1             0             0.216534      0             0.216534      0             1             0
2   -1.536794  0.295117  8   1.330553            3             1             0.240687      1             0.240687      1             0             0
3   0.501521  0.007457  1   -1.225848            2             0             -0.108918      0             -0.108918      0             0             0
4   2.063884  0.388871  2   0.785728            1             1             -0.365276            0             -0.365276      0             1             0
...
9995  1.245488  0.007457  5   -1.225848            2             1             -0.065419            0             -0.065419      0             0             1
9996  -1.391339  -0.373958    10  -0.306379            1             1             0.027988            0             0.027988      0             0             1
9997  0.604988  -0.278054  7   -1.225848            1             0             -1.009643            1             -1.009643      1             0             0
9998  1.256835  0.295117  3   -0.022608            2             1             0.012523            1             0.012523            1             1             0
9999  1.463771  1.041433  4   0.859965            1             1             -0.107670            0             -0.107670      0             0             0
10000 rows x 12 columns

In [39]: from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier

In [40]: lr = LogisticRegression(random_state=42)
lr.fit(X=train, y=train)
y_pred_train = lr.predict(X=train)
y_pred_test = lr.predict(X=test)
array([0, 0, 0, ..., 0, 0, 0])

In [41]: from sklearn.metrics import accuracy_score
acc_train = accuracy_score(y_train, y_pred_train)
acc_test = accuracy_score(y_test, y_pred_test)
0.81

In [42]: c1f1 = RandomForestClassifier(random_state=42)
c1f1.fit(X=train, y=train)
y_pred_train = c1f1.predict(X=train)
y_pred_test = c1f1.predict(X=test)
array([0, 0, 1, ..., 1, 1, 0])

In [43]: acc_train1 = accuracy_score(y_train, y_pred_train1)
acc_test1 = accuracy_score(y_test, y_pred_test1)
0.67

In [44]:
```