# ▾ Bank Customer Churn Prediction

In this project, we use supervised learning models to identify customers who are likely to churn in the future. Furthermore, we will analyze top factors that influence user retention. [Dataset information](.).

## ▾ Contents

- [Part 1: Data Exploration](.)
- [Part 2: Feature Preprocessing](.)
- [Part 3: Model Training and Results Evaluation](.)

## ▾ Part 0: Setup Google Drive Environment / Data Collection

check this [link](.) for more info

```
1 # install pydrive to load data
2 !pip install -U -q PyDrive
3
4 from pydrive.auth import GoogleAuth
5 from pydrive.drive import GoogleDrive
6 from google.colab import auth
7 from oauth2client.client import GoogleCredentials
8
9 auth.authenticate_user()
10 gauth = GoogleAuth()
11 gauth.credentials = GoogleCredentials.get_application_default()
12 drive = GoogleDrive(gauth)
```

```
1 # the same way we get id from last class
2 #https://drive.google.com/file/d/1szdCZ98EK59cfJ4jG03g1HOv_OhC1oyN/view?usp=sharing
3 id = "1szdCZ98EK59cfJ4jG03g1HOv_OhC1oyN"
4 file = drive.CreateFile({'id':id})
5 file.GetContentFile('bank_churn.csv')
```

```
1 import pandas as pd
2
3 df = pd.read_csv('bank_churn.csv')
4 df.head()
```

| | RowNumber | CustomerId | Surname | CreditScore | Geography | Gender | Age | Tenure | Balanc |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 15634602 | Hargrave | 619 | France | Female | 42 | 2 | 0.0 |
| **1** | 2 | 15647311 | Hill | 608 | Spain | Female | 41 | 1 | 83807.8 |
| **2** | 3 | 15619304 | Onio | 502 | France | Female | 42 | 8 | 159660.8 |
| **3** | 4 | 15701354 | Boni | 699 | France | Female | 39 | 1 | 0.0 |
| **4** | 5 | 15737888 | Mitchell | 850 | Spain | Female | 43 | 2 | 125510.8 |

# Part 1: Data Exploration

## Part 1.1: Understand the Raw Dataset

```
1 import pandas as pd
2 import numpy as np
3
4 churn_df = pd.read_csv('bank_churn.csv')
```

```
1 churn_df.head()
```

| | RowNumber | CustomerId | Surname | CreditScore | Geography | Gender | Age | Tenure | Balanc |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 15634602 | Hargrave | 619 | France | Female | 42 | 2 | 0.0 |
| **1** | 2 | 15647311 | Hill | 608 | Spain | Female | 41 | 1 | 83807.8 |
| **2** | 3 | 15619304 | Onio | 502 | France | Female | 42 | 8 | 159660.8 |
| **3** | 4 | 15701354 | Boni | 699 | France | Female | 39 | 1 | 0.0 |
| **4** | 5 | 15737888 | Mitchell | 850 | Spain | Female | 43 | 2 | 125510.8 |

```
1 # check data info
2 churn_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 14 columns):
 #   Column          Non-Null Count  Dtype
---  ------          --------------  -----
 0   RowNumber       10000 non-null  int64
 1   CustomerId      10000 non-null  int64
 2   Surname         10000 non-null  object
 3   CreditScore     10000 non-null  int64
```

```
 4   Geography        10000 non-null   object
 5   Gender           10000 non-null   object
 6   Age              10000 non-null   int64
 7   Tenure           10000 non-null   int64
 8   Balance          10000 non-null   float64
 9   NumOfProducts    10000 non-null   int64
 10  HasCrCard        10000 non-null   int64
 11  IsActiveMember   10000 non-null   int64
 12  EstimatedSalary  10000 non-null   float64
 13  Exited           10000 non-null   int64
dtypes: float64(2), int64(9), object(3)
memory usage: 1.1+ MB
```

```
1 # check the unique values for each column
2 churn_df.nunique()
```

```
RowNumber          10000
CustomerId         10000
Surname             2932
CreditScore          460
Geography              3
Gender                 2
Age                   70
Tenure                11
Balance             6382
NumOfProducts          4
HasCrCard              2
IsActiveMember         2
EstimatedSalary     9999
Exited                 2
dtype: int64
```

```
1 # Get target variable
2 y = churn_df['Exited']
```

## Part 1.2: Understand the features

```
1 # check missing values
2 churn_df.isnull().sum()
```

```
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
      dtype: int64
```
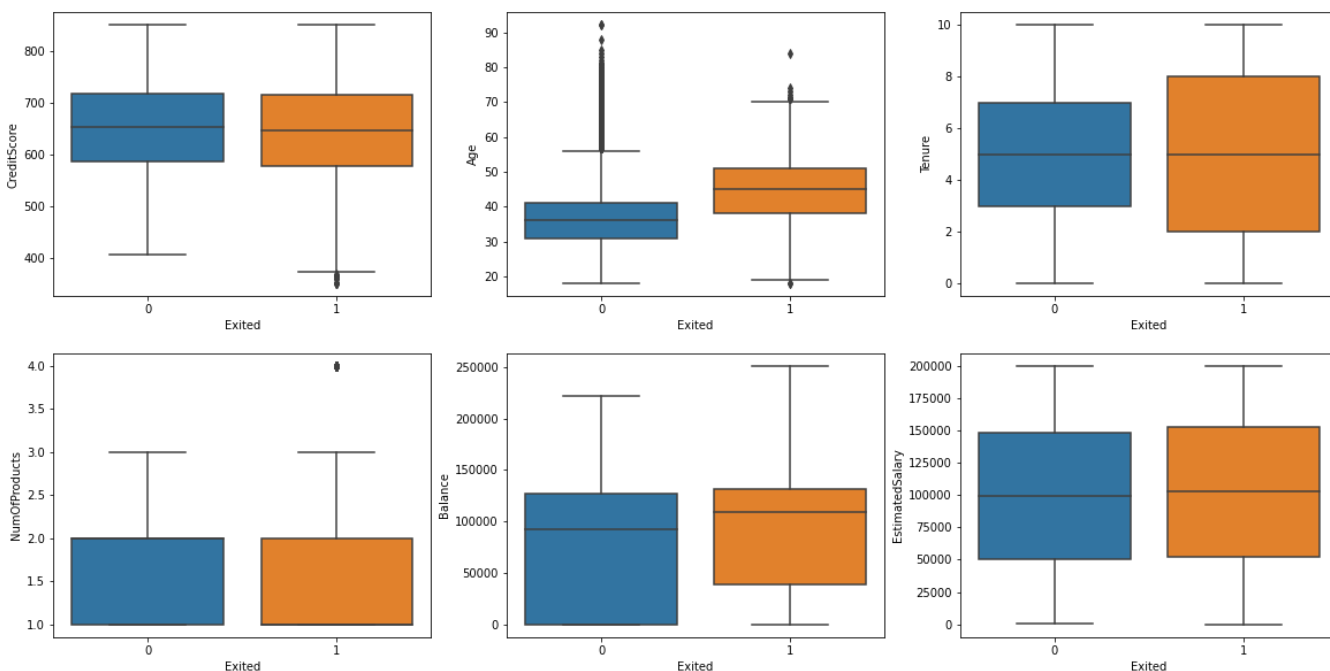
```
1 # understand Numerical feature
2 # discrete/continuous
3 # 'CreditScore', 'Age', 'Tenure', 'NumberOfProducts'
4 # 'Balance', 'EstimatedSalary'
5 churn_df[['CreditScore', 'Age', 'Tenure', 'NumOfProducts','Balance', 'EstimatedSalary']].d
```

|       | CreditScore | Age | Tenure | NumOfProducts | Balance | EstimatedS |
|-------|-------------|-----|--------|---------------|---------|-----------|
| count | 10000.000000 | 10000.000000 | 10000.000000 | 10000.000000 | 10000.000000 | 10000.0 |
| mean | 650.528800 | 38.921800 | 5.012800 | 1.530200 | 76485.889288 | 100090.2 |
| std | 96.653299 | 10.487806 | 2.892174 | 0.581654 | 62397.405202 | 57510.4 |
| min | 350.000000 | 18.000000 | 0.000000 | 1.000000 | 0.000000 | 11.5 |
| 25% | 584.000000 | 32.000000 | 3.000000 | 1.000000 | 0.000000 | 51002. |
| 50% | 652.000000 | 37.000000 | 5.000000 | 1.000000 | 97198.540000 | 100193.9 |
| 75% | 718.000000 | 44.000000 | 7.000000 | 2.000000 | 127644.240000 | 149388.2 |
| max | 850.000000 | 92.000000 | 10.000000 | 4.000000 | 250898.090000 | 199992.4 |

```
1 # check the feature distribution
2 # pandas.DataFrame.describe()
3 # boxplot, distplot, countplot
4 import matplotlib.pyplot as plt
5 import seaborn as sns
```

```
1 # boxplot for numerical feature
2 _,axss = plt.subplots(2,3, figsize=[20,10])
3 sns.boxplot(x='Exited', y ='CreditScore', data=churn_df, ax=axss[0][0])
4 sns.boxplot(x='Exited', y ='Age', data=churn_df, ax=axss[0][1])
5 sns.boxplot(x='Exited', y ='Tenure', data=churn_df, ax=axss[0][2])
6 sns.boxplot(x='Exited', y ='NumOfProducts', data=churn_df, ax=axss[1][0])
7 sns.boxplot(x='Exited', y ='Balance', data=churn_df, ax=axss[1][1])
8 sns.boxplot(x='Exited', y ='EstimatedSalary', data=churn_df, ax=axss[1][2])
```
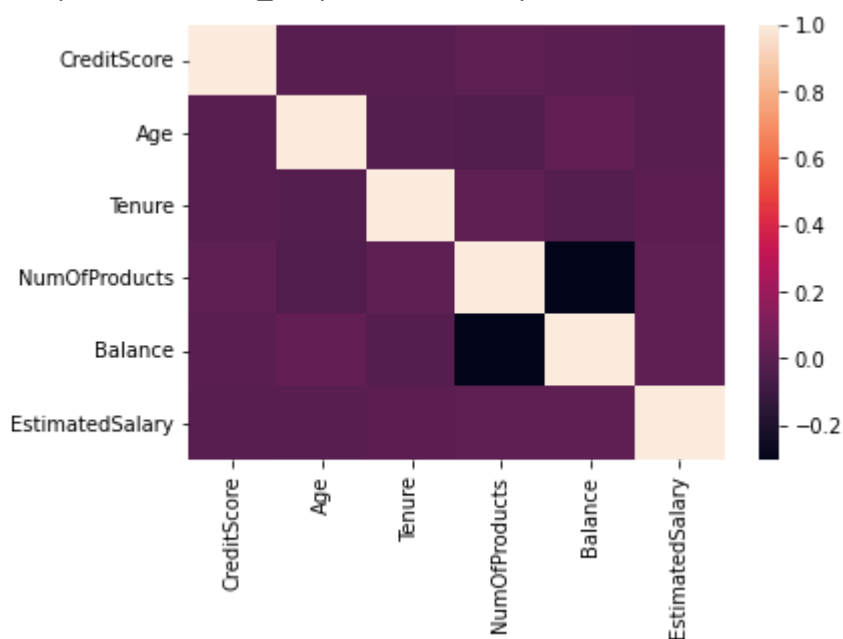
<matplotlib.axes._subplots.AxesSubplot at 0x7fb3ab028c10>



```
1 # correlations between features
2 corr_score = churn_df[['CreditScore', 'Age', 'Tenure', 'NumOfProducts','Balance', 'Estimat
3
4 # show heapmap of correlations
5 sns.heatmap(corr_score)
```
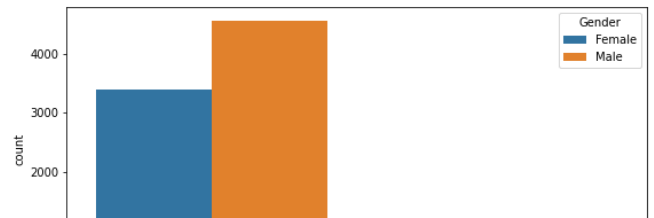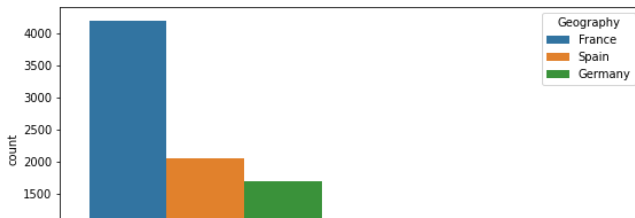
<matplotlib.axes._subplots.AxesSubplot at 0x7fb3aae68ad0>

```
1 # check the actual values of correlations
2 corr_score
```

|  | CreditScore | Age | Tenure | NumOfProducts | Balance | EstimatedSal |
|---|---|---|---|---|---|---|
| CreditScore | 1.000000 | -0.003965 | 0.000842 | 0.012238 | 0.006268 | -0.0013 |
| Age | -0.003965 | 1.000000 | -0.009997 | -0.030680 | 0.028308 | -0.0072 |
| Tenure | 0.000842 | -0.009997 | 1.000000 | 0.013444 | -0.012254 | 0.007 |
| NumOfProducts | 0.012238 | -0.030680 | 0.013444 | 1.000000 | -0.304180 | 0.014 |
| Balance | 0.006268 | 0.028308 | -0.012254 | -0.304180 | 1.000000 | 0.012 |
| EstimatedSalary | -0.001384 | -0.007201 | 0.007784 | 0.014204 | 0.012797 | 1.000 |

```
1 # understand categorical feature
2 # 'Geography', 'Gender'
3 # 'HasCrCard', 'IsActiveMember'
4 _,axss = plt.subplots(2,2, figsize=[20,10])
5 sns.countplot(x='Exited', hue='Geography', data=churn_df, ax=axss[0][0])
6 sns.countplot(x='Exited', hue='Gender', data=churn_df, ax=axss[0][1])
7 sns.countplot(x='Exited', hue='HasCrCard', data=churn_df, ax=axss[1][0])
8 sns.countplot(x='Exited', hue='IsActiveMember', data=churn_df, ax=axss[1][1])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fb3aac7fcd0>
```



# Part 2: Feature Preprocessing

```
1 # Get feature space by dropping useless feature
2 to_drop = ['RowNumber','CustomerId','Surname','Exited']
3 X = churn_df.drop(to_drop, axis=1)
```

```
1 X.head()
```

|   | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsA |
|---|-------------|-----------|--------|-----|--------|---------|---------------|-----------|-----|
| 0 | 619 | France | Female | 42 | 2 | 0.00 | 1 | 1 | |
| 1 | 608 | Spain | Female | 41 | 1 | 83807.86 | 1 | 0 | |
| 2 | 502 | France | Female | 42 | 8 | 159660.80 | 3 | 1 | |
| 3 | 699 | France | Female | 39 | 1 | 0.00 | 2 | 0 | |
| 4 | 850 | Spain | Female | 43 | 2 | 125510.82 | 1 | 1 | |

```
1 X.dtypes
```

```
CreditScore          int64
Geography            object
Gender               object
Age                  int64
Tenure               int64
Balance              float64
NumOfProducts        int64
HasCrCard            int64
IsActiveMember       int64
EstimatedSalary      float64
dtype: object
```

```
1 cat_cols = X.columns[X.dtypes == 'O']
2 num_cols = X.columns[(X.dtypes == 'float64') | (X.dtypes == 'int64')]
```

```
1 num_cols
```

```
Index(['CreditScore', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'HasCrCard',
       'IsActiveMember', 'EstimatedSalary'],
```

```
        dtype='object')
```

```
1 cat_cols
```

```
    Index(['Geography', 'Gender'], dtype='object')
```

## Split dataset

```
 1 # Splite data into training and testing
 2 # 100 -> 75:y=1, 25:y=0
 3 # training(80): 60 y=1; 20 y=0
 4 # testing(20):  15 y=1; 5 y=0
 5
 6 from sklearn import model_selection
 7
 8 # Reserve 25% for testing
 9 # stratify example:
10 # 100 -> y: 80 '0', 20 '1' -> 4:1
11 # 80% training 64: '0', 16:'1' -> 4:1
12 # 20% testing  16:'0', 4: '1' -> 4:1
13 X_train, X_test, y_train, y_test = model_selection.train_test_split(X, y, test_size=0.25,
14
15 print('training data has ' + str(X_train.shape[0]) + ' observation with ' + str(X_train.sh
16 print('test data has ' + str(X_test.shape[0]) + ' observation with ' + str(X_test.shape[1]
```

```
    training data has 7500 observation with 10 features
    test data has 2500 observation with 10 features
```

- 10000 -> 8000 '0' + 2000 '1'

- 25% test 75% training

---

without stratified sampling:

## ▾ • extreme case:

1. testing: 2000 '1' + 500 '0'
2. training: 7500 '0'

---

with stratified sampling:

1. testing: 2000 '0' + 500 '1'
2. training: 6000 '0' + 1500 '1'

Read more for handling [categorical feature](#), and there is an awesome package for [encoding](#).

```
1 X_train.head()
```

|      | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | ] |
|------|-------------|-----------|--------|-----|--------|---------|---------------|-----------|---|
| 7971 | 633 | Spain | Male | 42 | 10 | 0.00 | 1 | 0 | |
| 9152 | 708 | Germany | Female | 23 | 4 | 71433.08 | 1 | 1 | |
| 6732 | 548 | France | Female | 37 | 9 | 0.00 | 2 | 0 | |
| 902 | 645 | France | Female | 48 | 7 | 90612.34 | 1 | 1 | |
| 2996 | 729 | Spain | Female | 45 | 7 | 91091.06 | 2 | 1 | |

```
1 # One hot encoding
2 # another way: get_dummies
3 from sklearn.preprocessing import OneHotEncoder
4
5 def OneHotEncoding(df, enc, categories):
6   transformed = pd.DataFrame(enc.transform(df[categories]).toarray(), columns=enc.get_feat
7   return pd.concat([df.reset_index(drop=True), transformed], axis=1).drop(categories, axis
8
9 categories = ['Geography']
10 enc_ohe = OneHotEncoder()
11 enc_ohe.fit(X_train[categories])
12
13 X_train = OneHotEncoding(X_train, enc_ohe, categories)
14 X_test = OneHotEncoding(X_test, enc_ohe, categories)
15
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: F
  warnings.warn(msg, category=FutureWarning)
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: F
  warnings.warn(msg, category=FutureWarning)
```

```
1 X_train.head()
```

|   | CreditScore | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember |
|---|-------------|--------|-----|--------|---------|---------------|-----------|----------------|
| 0 | 633 | Male | 42 | 10 | 0.00 | 1 | 0 | 1 |
| 1 | 708 | Female | 23 | 4 | 71433.08 | 1 | 1 | 0 |
| 2 | 548 | Female | 37 | 9 | 0.00 | 2 | 0 | 0 |
| 3 | 645 | Female | 48 | 7 | 90612.34 | 1 | 1 | 1 |
| 4 | 729 | Female | 45 | 7 | 91091.06 | 2 | 1 | 0 |

```
1 # Ordinal encoding
2 from sklearn.preprocessing import OrdinalEncoder
3
4 categories = ['Gender']
5 enc_oe = OrdinalEncoder()
6 enc_oe.fit(X_train[categories])
7
8 X_train[categories] = enc_oe.transform(X_train[categories])
9 X_test[categories] = enc_oe.transform(X_test[categories])
```

```
1 X_train.head()
```

|   | CreditScore | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember |
|---|---|---|---|---|---|---|---|---|
| **0** | 633 | 1.0 | 42 | 10 | 0.00 | 1 | 0 | 1 |
| **1** | 708 | 0.0 | 23 | 4 | 71433.08 | 1 | 1 | 0 |
| **2** | 548 | 0.0 | 37 | 9 | 0.00 | 2 | 0 | 0 |
| **3** | 645 | 0.0 | 48 | 7 | 90612.34 | 1 | 1 | 1 |
| **4** | 729 | 0.0 | 45 | 7 | 91091.06 | 2 | 1 | 0 |

## Standardize/Normalize Data

```
1 # Scale the data, using standardization
2 # standardization (x-mean)/std
3 # normalization (x-x_min)/(x_max-x_min) ->[0,1]
4
5 # 1. speed up gradient descent
6 # 2. same scale
7 # 3. algorithm requirments
8
9 # for example, use training data to train the standardscaler to get mean and std
10 # apply mean and std to both training and testing data.
11 # fit_transform does the training and applying, transform only does applying.
12 # Because we can't use any info from test, and we need to do the same modification
13 # to testing data as well as training data
14
15 # https://scikit-learn.org/stable/auto_examples/preprocessing/plot_all_scaling.html#sphx-g
16 # https://scikit-learn.org/stable/modules/preprocessing.html
17
18
19 # min-max example: (x-x_min)/(x_max-x_min)
20 # [1,2,3,4,5,6,100] -> fit(min:1, max:6) (scalar.min = 1, scalar.max = 6) -> transform [(1
21 # scalar.fit(train) -> min:1, max:100
22 # scalar.transform(apply to x) -> apply min:1, max:100 to X_train
23 # scalar.transform -> apply min:1, max:100 to X_test
24
```

```
25 # scalar.fit -> mean:1, std:100
26 # scalar.transform -> apply mean:1, std:100 to X_train
27 # scalar.transform -> apply mean:1, std:100 to X_test
28
29 from sklearn.preprocessing import StandardScaler
30 scaler = StandardScaler()
31 scaler.fit(X_train[num_cols])
32 X_train[num_cols] = scaler.transform(X_train[num_cols])
33 X_test[num_cols] = scaler.transform(X_test[num_cols])
```

```
1 X_train.head()
```

|   | CreditScore | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember |
|---|---|---|---|---|---|---|---|---|
| **0** | 633 | 1.0 | 42 | 10 | 0.00 | 1 | 0 | 1 |
| **1** | 708 | 0.0 | 23 | 4 | 71433.08 | 1 | 1 | 0 |
| **2** | 548 | 0.0 | 37 | 9 | 0.00 | 2 | 0 | 0 |
| **3** | 645 | 0.0 | 48 | 7 | 90612.34 | 1 | 1 | 1 |
| **4** | 729 | 0.0 | 45 | 7 | 91091.06 | 2 | 1 | 0 |

## ▾ Part 3: Model Training and Result Evaluation

## ▾ Part 3.1: Model Training

```
1 #@title build models                          build models
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.neighbors import KNeighborsClassifier
4 from sklearn.linear_model import LogisticRegression
5
6 # Logistic Regression
7 classifier_logistic = LogisticRegression()
8
9 # K Nearest Neighbors
10 classifier_KNN = KNeighborsClassifier()
11
12 # Random Forest
13 classifier_RF = RandomForestClassifier()
```

```
1 # Train the model
2 classifier_logistic.fit(X_train, y_train)
```

```
    LogisticRegression()
```

```
1 # Prediction of test data
2 classifier_logistic.predict(X_test)
```

```
    array([0, 0, 0, ..., 0, 0, 0])
```

```
1 # Accuracy of test data
2 classifier_logistic.score(X_test, y_test)
```

```
    0.7936
```

```
1 # Use 5-fold Cross Validation to get the accuracy for different models
2 model_names = ['Logistic Regression','KNN','Random Forest']
3 model_list = [classifier_logistic, classifier_KNN, classifier_RF]
4 count = 0
5
6 for classifier in model_list:
7     cv_score = model_selection.cross_val_score(classifier, X_train, y_train, cv=5)
8     print(cv_score)
9     print('Model accuracy of ' + model_names[count] + ' is ' + str(cv_score.mean()))
10    count += 1
```

```
    [0.78533333 0.78933333 0.78866667 0.792      0.78933333]
    Model accuracy of Logistic Regression is 0.7889333333333333
    [0.76333333 0.772      0.75733333 0.76466667 0.758     ]
    Model accuracy of KNN is 0.7630666666666668
    [0.88133333 0.864      0.85266667 0.85533333 0.86533333]
    Model accuracy of Random Forest is 0.8637333333333332
```

## ▾ (Optional) Part 3.2: Use Grid Search to Find Optimal Hyperparameters

alternative: random search

```
1 #Loss/cost function --> (wx + b - y) ^2 + λ * |w| --> λ is a hyperparameter
```

```
1 from sklearn.model_selection import GridSearchCV
2
3 # helper function for printing out grid search results
4 def print_grid_search_metrics(gs):
5     print ("Best score: " + str(gs.best_score_))
6     print ("Best parameters set:")
7     best_parameters = gs.best_params_
8     for param_name in sorted(best_parameters.keys()):
9         print(param_name + ':' + str(best_parameters[param_name]))
```

## ▾ Part 3.2.1: Find Optimal Hyperparameters - LogisticRegression

```python
1 # Possible hyperparamter options for Logistic Regression Regularization
2 # Penalty is choosed from L1 or L2
3 # C is the 1/lambda value(weight) for L1 and L2
4 # solver: algorithm to find the weights that minimize the cost function
5
6 # ('l1', 0.01)('l1', 0.05) ('l1', 0.1) ('l1', 0.2)('l1', 1)
7 # ('12', 0.01)('l2', 0.05) ('l2', 0.1) ('l2', 0.2)('l2', 1)
8 parameters = {
9     'penalty':('l1', 'l2'),
10    'C':(0.01, 0.05, 0.1, 0.2, 1)
11 }
12 Grid_LR = GridSearchCV(LogisticRegression(solver='liblinear'),parameters, cv=5)
13 Grid_LR.fit(X_train, y_train)
```

```
GridSearchCV(cv=5, estimator=LogisticRegression(solver='liblinear'),
             param_grid={'C': (0.01, 0.05, 0.1, 0.2, 1),
                         'penalty': ('l1', 'l2')})
```

```python
1 # the best hyperparameter combination
2 # C = 1/lambda
3 print_grid_search_metrics(Grid_LR)
```

```
Best score: 0.8130666666666666
Best parameters set:
C:0.05
penalty:l1
```

```python
1 # best model
2 best_LR_model = Grid_LR.best_estimator_
```

```python
1 best_LR_model.predict(X_test)
```

```
array([0, 0, 0, ..., 0, 0, 0])
```

```python
1 best_LR_model.score(X_test, y_test)
```

```
0.8084
```

```python
1 LR_models = pd.DataFrame(Grid_LR.cv_results_)
2 res = (LR_models.pivot(index='param_penalty', columns='param_C', values='mean_test_score')
3           )
4 _ = sns.heatmap(res, cmap='viridis')
```

## Part 3.2.2: Find Optimal Hyperparameters: KNN

```
1 # Possible hyperparamter options for KNN
2 # Choose k
3 parameters = {
4     'n_neighbors':[1,3,5,7,9]
5 }
6 Grid_KNN = GridSearchCV(KNeighborsClassifier(),parameters, cv=5)
7 Grid_KNN.fit(X_train, y_train)
```

```
    GridSearchCV(cv=5, estimator=KNeighborsClassifier(),
                 param_grid={'n_neighbors': [1, 3, 5, 7, 9]})
```

```
1 # best k
2 print_grid_search_metrics(Grid_KNN)
```

```
    Best score: 0.7856
    Best parameters set:
    n_neighbors:9
```

```
1 best_KNN_model = Grid_KNN.best_estimator_
```

## Part 3.2.3: Find Optimal Hyperparameters: Random Forest

```
1 # Possible hyperparamter options for Random Forest
2 # Choose the number of trees
3 parameters = {
4     'n_estimators' : [60,80,100],
5     'max_depth': [1,5,10]
6 }
7 Grid_RF = GridSearchCV(RandomForestClassifier(),parameters, cv=5)
8 Grid_RF.fit(X_train, y_train)
```

```
    GridSearchCV(cv=5, estimator=RandomForestClassifier(),
                 param_grid={'max_depth': [1, 5, 10],
                             'n_estimators': [60, 80, 100]})
```

```
1 # best number of tress
2 print_grid_search_metrics(Grid_RF)
```

```
    Best score: 0.8667999999999999
    Best parameters set:
    max_depth:10
    n_estimators:100
```

```
1 # best random forest
2 best_RF_model = Grid_RF.best_estimator_
```

```
1 best_RF_model
```

```
    RandomForestClassifier(max_depth=10)
```

## ▼ Part 3.3: Model Evaluation - Confusion Matrix (Precision, Recall, Accuracy)

class of interest as positive

TP: correctly labeled real churn

Precision(PPV, positive predictive value): tp / (tp + fp); Total number of true predictive churn divided by the total number of predictive churn; High Precision means low fp, not many return users were predicted as churn users.

Recall(sensitivity, hit rate, true positive rate): tp / (tp + fn) Predict most postive or churn user correctly. High recall means low fn, not many churn users were predicted as return users.

```
 1 from sklearn.metrics import confusion_matrix
 2 from sklearn.metrics import classification_report
 3 from sklearn.metrics import precision_score
 4 from sklearn.metrics import recall_score
 5
 6 # calculate accuracy, precision and recall, [[tn, fp],[]]
 7 def cal_evaluation(classifier, cm):
 8     tn = cm[0][0]
 9     fp = cm[0][1]
10     fn = cm[1][0]
11     tp = cm[1][1]
12     accuracy  = (tp + tn) / (tp + fp + fn + tn + 0.0)
13     precision = tp / (tp + fp + 0.0)
14     recall = tp / (tp + fn + 0.0)
15     print (classifier)
16     print ("Accuracy is: " + str(accuracy))
17     print ("precision is: " + str(precision))
18     print ("recall is: " + str(recall))
19     print ()
20
```

```
21 # print out confusion matrices
22 def draw_confusion_matrices(confusion_matricies):
23     class_names = ['Not','Churn']
24     for cm in confusion_matrices:
25         classifier, cm = cm[0], cm[1]
26         cal_evaluation(classifier, cm)
```

```
1 # Confusion matrix, accuracy, precison and recall for random forest and logistic regressio
2 confusion_matrices = [
3     ("Random Forest", confusion_matrix(y_test,best_RF_model.predict(X_test))),
4     ("Logistic Regression", confusion_matrix(y_test,best_LR_model.predict(X_test))),
5     ("K nearest neighbor", confusion_matrix(y_test, best_KNN_model.predict(X_test)))
6 ]
7
8 draw_confusion_matrices(confusion_matrices)
```

```
Random Forest
Accuracy is: 0.8604
precision is: 0.8007518796992481
recall is: 0.41846758349705304

Logistic Regression
Accuracy is: 0.8084
precision is: 0.6056338028169014
recall is: 0.16895874263261296

K nearest neighbor
Accuracy is: 0.776
precision is: 0.2
recall is: 0.03339882121807466
```

## ▾ Part 3.4: Model Evaluation - ROC & AUC

RandomForestClassifier, KNeighborsClassifier and LogisticRegression have predict_prob() function

## ▾ Part 3.4.1: ROC of RF Model

```
1 from sklearn.metrics import roc_curve
2 from sklearn import metrics
3
4 # Use predict_proba to get the probability results of Random Forest
5 y_pred_rf = best_RF_model.predict_proba(X_test)[:, 1]
6 fpr_rf, tpr_rf, _ = roc_curve(y_test, y_pred_rf)
```
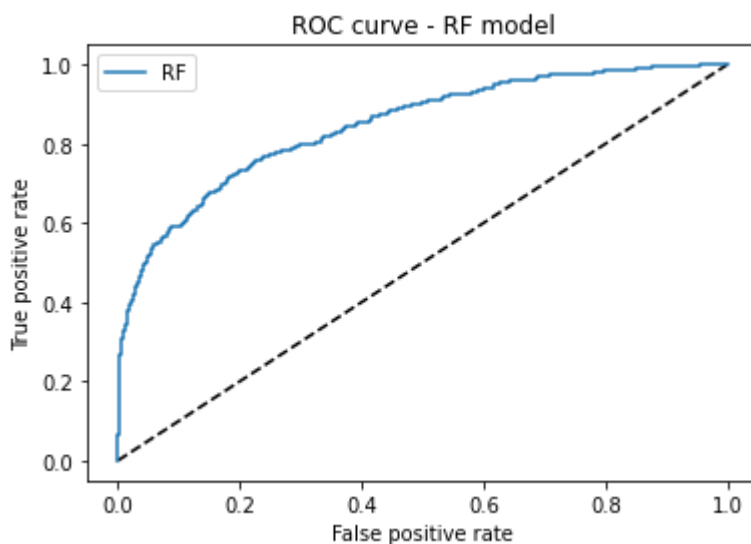
```
1 best_RF_model.predict_proba(X_test)
```

```
array([[0.69937921, 0.30062079],
       [0.95239486, 0.04760514],
       [0.74105858, 0.25894142],
       ...,
       [0.85967291, 0.14032709],
       [0.95083241, 0.04916759],
       [0.90647722, 0.09352278]])
```

```
1 # ROC curve of Random Forest result
2 import matplotlib.pyplot as plt
3 plt.figure(1)
4 plt.plot([0, 1], [0, 1], 'k--')
5 plt.plot(fpr_rf, tpr_rf, label='RF')
6 plt.xlabel('False positive rate')
7 plt.ylabel('True positive rate')
8 plt.title('ROC curve - RF model')
9 plt.legend(loc='best')
10 plt.show()
```



```
1 from sklearn import metrics
2
3 # AUC score
4 metrics.auc(fpr_rf,tpr_rf)
```
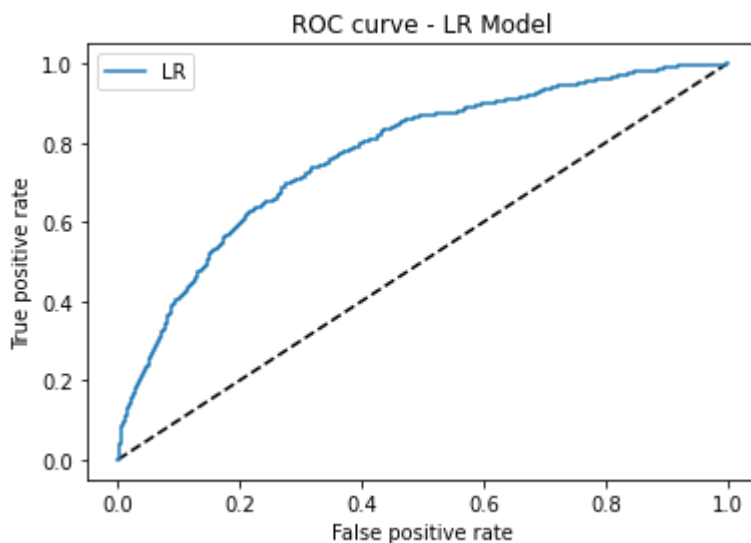
```
0.8457044914295074
```

## Part 3.4.1: ROC of LR Model

```
1 # Use predict_proba to get the probability results of Logistic Regression
2 y_pred_lr = best_LR_model.predict_proba(X_test)[:, 1]
3 fpr_lr, tpr_lr, thresh = roc_curve(y_test, y_pred_lr)
```

```
1 best_LR_model.predict_proba(X_test)
```

```
array([[0.83732252, 0.16267748],
       [0.91667203, 0.08332797],
       [0.85713621, 0.14286379],
       ...,
       [0.72754602, 0.27245398],
       [0.88726134, 0.11273866],
       [0.86330938, 0.13669062]])
```

```
1 # ROC Curve
2 plt.figure(1)
3 plt.plot([0, 1], [0, 1], 'k--')
4 plt.plot(fpr_lr, tpr_lr, label='LR')
5 plt.xlabel('False positive rate')
6 plt.ylabel('True positive rate')
7 plt.title('ROC curve - LR Model')
8 plt.legend(loc='best')
9 plt.show()
```



```
1 # AUC score
2 metrics.auc(fpr_lr,tpr_lr)
```

```
0.7703625055381831
```

# Part 4: Model Extra Functionality

## Part 4.1: Logistic Regression Model

The corelated features that we are interested in

```
1 X_with_corr = X.copy()
2
3 X_with_corr = OneHotEncoding(X_with_corr, enc_ohe, ['Geography'])
4 X_with_corr['Gender'] = enc_oe.transform(X_with_corr[['Gender']])
5 X_with_corr['SalaryInRMB'] = X_with_corr['EstimatedSalary'] * 6.4
6 X_with_corr.head()
```

/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: F
  warnings.warn(msg, category=FutureWarning)

| | CreditScore | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember |
|---|---|---|---|---|---|---|---|---|
| 0 | 619 | 0.0 | 42 | 2 | 0.00 | 1 | 1 | 1 |
| 1 | 608 | 0.0 | 41 | 1 | 83807.86 | 1 | 0 | 1 |
| 2 | 502 | 0.0 | 42 | 8 | 159660.80 | 3 | 1 | 0 |
| 3 | 699 | 0.0 | 39 | 1 | 0.00 | 2 | 0 | 0 |
| 4 | 850 | 0.0 | 43 | 2 | 125510.82 | 1 | 1 | 1 |

```
1 # add L1 regularization to logistic regression
2 # check the coef for feature selection
3 scaler = StandardScaler()
4 X_l1 = scaler.fit_transform(X_with_corr)
5 LRmodel_l1 = LogisticRegression(penalty="l1", C = 0.04, solver='liblinear')
6 LRmodel_l1.fit(X_l1, y)
7
8 indices = np.argsort(abs(LRmodel_l1.coef_[0]))[::-1]
9
10 print ("Logistic Regression (L1) Coefficients")
11 for ind in range(X_with_corr.shape[1]):
12   print ("{0} : {1}".format(X_with_corr.columns[indices[ind]],round(LRmodel_l1.coef_[0][in
```

Logistic Regression (L1) Coefficients
Age : 0.7307
IsActiveMember : -0.5046
Geography_Germany : 0.3121
Gender : -0.2409
Balance : 0.151
CreditScore : -0.0457
NumOfProducts : -0.0438
Tenure : -0.0271
EstimatedSalary : 0.0055
Geography_France : -0.0043
SalaryInRMB : 0.0037
HasCrCard : -0.0022
Geography_Spain : 0.0

```
 1 # add L2 regularization to logistic regression
 2 # check the coef for feature selection
 3 np.random.seed()
 4 scaler = StandardScaler()
 5 X_l2 = scaler.fit_transform(X_with_corr)
 6 LRmodel_l2 = LogisticRegression(penalty="l2", C = 0.1, solver='liblinear', random_state=42
 7 LRmodel_l2.fit(X_l2, y)
 8 LRmodel_l2.coef_[0]
 9
10 indices = np.argsort(abs(LRmodel_l2.coef_[0]))[::-1]
11
12 print ("Logistic Regression (L2) Coefficients")
13 for ind in range(X_with_corr.shape[1]):
14   print ("{0} : {1}".format(X_with_corr.columns[indices[ind]],round(LRmodel_l2.coef_[0][in
```

```
Logistic Regression (L2) Coefficients
Age : 0.751
IsActiveMember : -0.5272
Gender : -0.2591
Geography_Germany : 0.2279
Balance : 0.162
Geography_France : -0.1207
Geography_Spain : -0.089
CreditScore : -0.0637
NumOfProducts : -0.0586
Tenure : -0.0452
HasCrCard : -0.0199
SalaryInRMB : 0.0137
EstimatedSalary : 0.0137
```

## Part 4.2: Random Forest Model - Feature Importance Discussion

```
1 X_RF = X.copy()
2
3 X_RF = OneHotEncoding(X_RF, enc_ohe, ['Geography'])
4 X_RF['Gender'] = enc_oe.transform(X_RF[['Gender']])
5
6 X_RF.head()
```

```
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: F
  warnings.warn(msg, category=FutureWarning)
```

| CreditScore | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsActiveMember |
|---|---|---|---|---|---|---|---|

```
 1 # check feature importance of random forest for feature selection
 2 forest = RandomForestClassifier()
 3 forest.fit(X_RF, y)
 4
 5 importances = forest.feature_importances_
 6
 7 indices = np.argsort(importances)[::-1]
 8
 9 # Print the feature ranking
10 print("Feature importance ranking by Random Forest Model:")
11 for ind in range(X.shape[1]):
12   print ("{0} : {1}".format(X_RF.columns[indices[ind]],round(importances[indices[ind]], 4)
```

```
Feature importance ranking by Random Forest Model:
Age : 0.2372
EstimatedSalary : 0.1483
Balance : 0.1429
CreditScore : 0.1426
NumOfProducts : 0.1302
Tenure : 0.0822
IsActiveMember : 0.0396
Geography_Germany : 0.0217
HasCrCard : 0.0186
Gender : 0.0182
```

✓  0s    completed at 10:30 AM                                                    ● ✕