3 Geography 10000 non-null object 4 Gender 10000 non-null object 5 Age 10000 non-null int64 6 Tenure 10000 non-null int64 7 Balance 10000 non-null float64 8 Num Of Products 10000 non-null int64 9 Has Credit Card 10000 non-null int64 10 Is Active Member 10000 non-null int64 11 Estimated Salary 10000 non-null float64 12 Chur 10000 non-null int64 dtypes: float64(2), int64(8), object(3) memory usage: 1015.8+ KB df.duplicated('CustomerId') 0 False
False 1 False 2 False 3 False 4 False
<pre>df = df.set_index('CustomerId') df.info() <class 'pandas.core.frame.dataframe'=""> Int64Index: 10000 entries, 15634602 to 15628319 Data columns (total 12 columns): # Column Non-Null Count Dtype</class></pre>
5 Tenure 10000 non-null int64 6 Balance 10000 non-null int64 7 Num Of Products 10000 non-null int64 8 Has Credit Card 10000 non-null int64 9 Is Active Member 10000 non-null int64 10 Estimated Salary 10000 non-null int64 11 Churn 10000 non-null int64 dtypes: float64(2), int64(7), object(3) memory usage: 1015.6+ KB Encoding df['Geography'].value_counts()
France Germany 25094 Germany 25098 Spain 2477 Name: Geography, dtype: int64 df.replace({'Geography': {'France':2, 'Germany':1, 'Spain':0}}, inplace = True) df Surname CreditScore Geography Gender Age Tenure Balance Num Of Products Has Credit Card Is Active Member Estimated Salary Churn Customerid 15634602 Hargrave 619 2 Female 42 2 0.00 1 1 1 101348.88 1 15647311 Hill 608 0 Female 41 1 83807.86 1 0 0 1 112542.58 0
15619304 Onio 502 2 Female 42 8 159660.80 3 1 0 113931.57 1 15701354 Boni 699 2 Female 39 1 0.00 2 0 93826.63 0 15737888 Mitchell 850 0 Female 43 2 125510.82 1 1 1 79084.10 0
<pre>df.replace({'Gender': {'Female':1, 'Male':0}}, inplace = True) df['Num Of Products'].value_counts() 1 5084 2 4590 3 266 4 60 Name: Num Of Products, dtype: int64 df.replace({'Num Of Products': {1:0, 2:1, 3:1, 4:1}}, inplace = True)</pre>
Surname CreditScore Geography Gender Age Tenure Balance Num Of Products Has Credit Care Section
15606229 Obijiaku 771 2 0 39 5 0.00 1 1 0 96270.64 0 15569892 Johnstone 516 2 0 35 10 57369.61 0 1 1 101699.77 0 15584532 Liu 709 2 1 36 7 0.00 0 0 1 42085.58 1 15682355 Sabbatini 772 1 0 42 3 75075.31 1 1 0 92888.52 1 15628319 Walker 792 2 1 28 4 130142.79 0 1 0 38190.78 0 df['Has Credit Card'].value_counts()
1 7055 0 2945 Name: Has Credit Card, dtype: int64 df['Is Active Member'].value_counts() 1 5151 0 4849 Name: Is Active Member, dtype: int64 df.loc[(df['Balance']==0), 'Churn'].value_counts() 0 3117 1 500 Name: Churn, dtype: int64
<pre>import numpy as np df['Zero Balance'] = np.where(df['Balance']>0, 1, 0) df['Zero Balance'].hist() <axessubplot:></axessubplot:></pre>
3000 2000 2000 2000 2000 2000 2000 2000
0 0 2064 2064 2064 2064 2064 2064 2064 2
<pre>df.columns Index(['Surname', 'CreditScore', 'Geography', 'Gender', 'Age', 'Tenure',</pre>
Handling Imbalance Data df['churn'].value_counts() 0 7963 1 2037 Name: Churn, dtype: int64 import seaborn as sns sns.countplot(x = 'Churn', data = df)
<pre><axessubplot:xlabel='churn', ylabel="count"></axessubplot:xlabel='churn',></pre> 8000 - 80
Random Under Sampling from imblearn.under_sampling import RandomUnderSampler rus = RandomUnderSampler(random_state = 2529) X_rus,Y_rus = rus.fit_resample(X,Y)
<pre>X_rus.shape, Y_rus.shape, X.shape, Y.shape ((4074, 11), (4074,), (10000, 11), (10000,)) Y.value_counts() 0 7963 1 2037 Name: Churn, dtype: int64 Y_rus.value_counts() 0 2037 1 2037 Name: Churn, dtype: int64</pre>
Random Over Sampling from imblearn.over_sampling import RandomOverSampler
ros = RandomOverSampler(random_state = 2529) X_ros,Y_ros = ros.fit_resample(X,Y) X_ros.shape, Y_ros.shape, X.shape, Y.shape ((15926, 11), (15926,), (10000, 11), (10000,)) Y.value_counts() 0 7963 1 2037 Name: Churn, dtype: int64
Y_ros.value_counts() 1 7963 0 7963 Name: Churn, dtype: int64 Y_ros.plot(kind = 'hist') <axessubplot:ylabel='frequency'> 8000 7000 6000</axessubplot:ylabel='frequency'>
Test Train Split
Split Orignal Data from sklearn.model_selection import train_test_split X_train, X_test, Y_train, Y_test = train_test_split(X,Y, test_size = 0.3, random_state = 25) Split Random Under Sample Data X_train_rus, X_test_rus, Y_train_rus, Y_test_rus = train_test_split(X_rus, Y_rus, test_size = 0.3, random_state = 25)
Split Random Over Sample Data X_train_ros, X_test_ros, Y_train_ros, Y_test_ros = train_test_split(X_ros, Y_ros, test_size = 0.3, random_state = 25) Standardize Features from sklearn.preprocessing import StandardScaler sc = StandardScaler()
Standardize Orignal Data X_train[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']] = sc.fit_transform(X_train[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']]) X_test[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']] = sc.fit_transform(X_test[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']]) Standardize Random Under Sample Data X_train_rus[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']] = sc.fit_transform(X_train_rus[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']])
<pre>X_test_rus[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']] = sc.fit_transform(X_test_rus[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']]) Standardize Random Over Sample Data X_train_ros[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']] = sc.fit_transform(X_train_ros[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']]) X_test_ros[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']] = sc.fit_transform(X_test_ros[['CreditScore', 'Age', 'Tenure', 'Balance', 'Estimated Salary']]) Modelling</pre>
<pre>from sklearn.svm import SVC model = SVC() model.fit (X_train,Y_train) v SVC svC() Y_pred = model.predict(X_test)</pre>
Model Accuracy from sklearn.metrics import confusion_matrix, classification_report confusion_matrix(Y_test,Y_pred) array([[2374, 45], [421, 160]], dtype=int64) print(classification_report(Y_test, Y_pred)) precision recall f1-score support
0 0.85 0.98 0.91 2419 1 0.78 0.28 0.41 581 accuracy 0.84 3000 macro avg 0.81 0.63 0.66 3000 weighted avg 0.84 0.84 0.81 3000 from sklearn.model_selection import GridSearchCV param_grid = {'C': [0.1,1,10], 'gamma': [1,0.1,0.01], 'kernel': ['rbf'], 'class_weight': ['balanced']}
grid = GridSearchCV(SVC(), param_grid, refit = True, verbose = 2, cv = 2) grid.fit(X_train, Y_train) Fitting 2 folds for each of 9 candidates, totalling 18 fits [CV] ENDC=0.1, class_weight=balanced, gamma=1, kernel=rbf; total time= 1.4s [CV] ENDC=0.1, class_weight=balanced, gamma=0.1, kernel=rbf; total time= 0.9s [CV] END C=0.1, class_weight=balanced, gamma=0.1, kernel=rbf; total time= 0.9s [CV] END C=0.1, class_weight=balanced, gamma=0.1, kernel=rbf; total time= 1.3s [CV] END C=0.1, class_weight=balanced, gamma=0.01, kernel=rbf; total time= 1.3s [CV] END C=0.1, class_weight=balanced, gamma=0.01, kernel=rbf; total time= 1.2s [CV] ENDC=1, class_weight=balanced, gamma=1, kernel=rbf; total time= 1.2s [CV] ENDC=1, class_weight=balanced, gamma=0.1, kernel=rbf; total time= 1.0s [CV] ENDC=1, class_weight=balanced, gamma=0.1, kernel=rbf; total time= 1.0s
[CV] END .C=1, class_weight=balanced, gamma=0.01, kernel=rbf; total time= 1.0s [CV] END .C=1, class_weight=balanced, gamma=0.01, kernel=rbf; total time= 1.0s [CV] ENDC=10, class_weight=balanced, gamma=1, kernel=rbf; total time= 1.2s [CV] ENDC=10, class_weight=balanced, gamma=1, kernel=rbf; total time= 1.1s [CV] END .C=10, class_weight=balanced, gamma=0.1, kernel=rbf; total time= 0.9s [CV] END .C=10, class_weight=balanced, gamma=0.1, kernel=rbf; total time= 0.9s [CV] END C=10, class_weight=balanced, gamma=0.01, kernel=rbf; total time= 1.0s [CV] END C=10, class_weight=balanced, gamma=0.01, kernel=rbf; total time= 0.9s [CV] END C=10, class_weight=balanced, gamma=0.01, kernel=rbf; total time= 0.8s ForidSearchCV
<pre>print(grid.best_estimator_) SVC(C=10, class_weight='balanced', gamma=1) grid_predictions = grid_predict(X_test) confusion_matrix(Y_test, grid_predictions) array([[2166, 253], [365, 216]], dtype=int64) print(classification_report(Y_test, grid_predictions))</pre>
0 0.86 0.90 0.88 2419 1 0.46 0.37 0.41 581 accuracy 0.79 3000 macro avg 0.66 0.63 0.64 3000 weighted avg 0.78 0.79 0.79 3000 Model with Random Under Sampling svc_rus = SVC()
svc_rus.fit(X_train_rus,Y_train_rus) * svc svc() Y_pred_rus = svc_rus.predict(X_test_rus) Model Accuracy confusion_matrix(Y_test_rus, Y_pred_rus)
confusion_matrix(Y_test_rus, Y_pred_rus) array([[483, 120],
print(classification_report(Y_test_rus, Y_pred_rus)) precision recall f1-score support 0 0.74 0.80 0.77 603 1 0.79 0.72 0.75 620 accuracy 0.76 1223 macro avg 0.76 0.76 0.76 1223 weighted avg 0.76 0.76 0.76 1223
print(classification_report(Y_test_rus, Y_pred_rus)) precision recall f1-score support 0 0.74 0.80 0.77 603 1 0.79 0.72 0.75 629 accuracy macro avg 0.76 0.76 1223 weighted avg 0.76 0.76 0.76 1223 Hyperparameter Tunning param_grid = {'C': [0.1,1,10],
print(classification_report(Y_test_rus, Y_pred_rus)) precision recall f1-scare support
proint(classification_report(v_test_rus, V_pred_rus)) precision recall f1-score support 0 0.74 0.80 0.77 0.72 0.75 0.90 accuraty
print(classification report(Y test rus, Y pred rus)) precision recall fi-score support 0 0.74 0.80 0.77 650 accuracy 0.78 0.78 0.78 0.78 1223 accuracy 0.78 0.78 0.78 1223 accuracy 0.78 0.78 0.78 0.78 1223 Hyperparameter Tunning paran grid = ('C': [6.1,1,1,10],
(Classification repair) Table
Procession Testal 1 - Sept 5 - Sept
Properties Pro