

Retention of Credit card customers

```
In [1]: #Importing libraries
import pandas as pd
import matplotlib.pyplot as plt
sns.set()

#Importing Lib
import sys
from sklearn.model_selection import cross_val_score
from sklearn.preprocessing import OneHotEncoder
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
import warnings
warnings.filterwarnings('ignore')
from sklearn.metrics import f1_score, precision_score, roc_auc_score
```

```
In [2]: import numpy as np

In [3]: #reading data
df = pd.read_csv('customer_retention.csv')
```

```
In [4]: df.shape

Out[4]: (10127, 23)
```

```
In [5]: df.head(10)

Out[5]:
```

	CLIENTNUM	Attrition_Flag	Customer_Age	Gender	Dependent_count	Education_Level	Marital_Status	Income_Category	Card_Ca
0	768805383	Existing Customer	45	M	3	High School	Married	60K--80K	
1	818770008	Existing Customer	49	F	5	Graduate	Single	Less than \$40K	
2	713982108	Existing Customer	51	M	3	Graduate	Married	80K--120K	
3	769911858	Existing Customer	40	F	4	High School	Unknown	Less than \$40K	
4	709106358	Existing Customer	40	M	3	Uneducated	Married	60K--80K	
5	713061558	Existing Customer	44	M	2	Graduate	Married	40K--60K	
6	810347208	Existing Customer	51	M	4	Unknown	Married	\$120K +	
7	819063208	Existing Customer	32	M	0	High School	Unknown	60K--80K	
8	710930508	Existing Customer	37	M	3	Uneducated	Single	60K--80K	
9	719661558	Existing Customer	48	M	2	Graduate	Single	80K--120K	

```
10 rows x 23 columns

In [6]: #dropping unnecessary columns
df = df.drop(['Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Educati
```

```
In [7]: df.columns

Out[7]:
```

```
Index(['CLIENTNUM', 'Attrition_Flag', 'Customer_Age', 'Gender', 'Dependent_count', 'Education_Level', 'Marital_Status', 'Income_Category', 'Card_Category', 'Months_on_book', 'Total_Relationship_Count', 'Months_Inactive_12_mon', 'Contacts_Count_12_mon', 'Credit_Limit', 'Total_Revolving_Bal', 'Avg_Open_To_Buy', 'Total_Amt_Chng_Q4_Q1', 'Total_Trans_Amt', 'Total_Trans_Ct', 'Total_Ct_Chng_Q4_Q1', 'Avg_Utilization_Ratio', 'Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Education_Level_Months_Inactive_12_mon_2'],
      dtype='object')
```

```
In [8]: df = df.drop(['Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Educati
```

```
In [9]: df.columns

Out[9]:
```

```
Index(['CLIENTNUM', 'Attrition_Flag', 'Customer_Age', 'Gender', 'Dependent_count', 'Education_Level', 'Marital_Status', 'Income_Category', 'Card_Category', 'Months_on_book', 'Total_Relationship_Count', 'Months_Inactive_12_mon', 'Contacts_Count_12_mon', 'Credit_Limit', 'Total_Revolving_Bal', 'Avg_Open_To_Buy', 'Total_Amt_Chng_Q4_Q1', 'Total_Trans_Amt', 'Total_Trans_Ct', 'Total_Ct_Chng_Q4_Q1', 'Avg_Utilization_Ratio'],
      dtype='object')
```

```
In [10]: df.shape

Out[10]: (10127, 21)
```

```
In [11]: for col in df.select_dtypes('category').columns.to_list():
          print(col + " : " + strX[col].cat.categories.to_list())
```

```
In [12]: #checking for missing values
df.isna().sum()
```

```
Out[12]: CLIENTNUM      0
Attrition_Flag      0
Customer_Age        0
Gender              0
Dependent_count     0
Education_Level     0
Marital_Status      0
Income_Category     0
Card_Category       0
Months_on_book      0
Total_Relationship_Count  0
Months_Inactive_12_mon  0
Contacts_Count_12_mon  0
Credit_Limit        0
Total_Revolving_Bal  0
Avg_Open_To_Buy     0
Total_Amt_Chng_Q4_Q1  0
Total_Trans_Amt     0
Total_Trans_Ct      0
Total_Ct_Chng_Q4_Q1  0
Avg_Utilization_Ratio  0
dtype: int64
```

```
In [13]: df2_corr = df.corr()
df2_corr
```

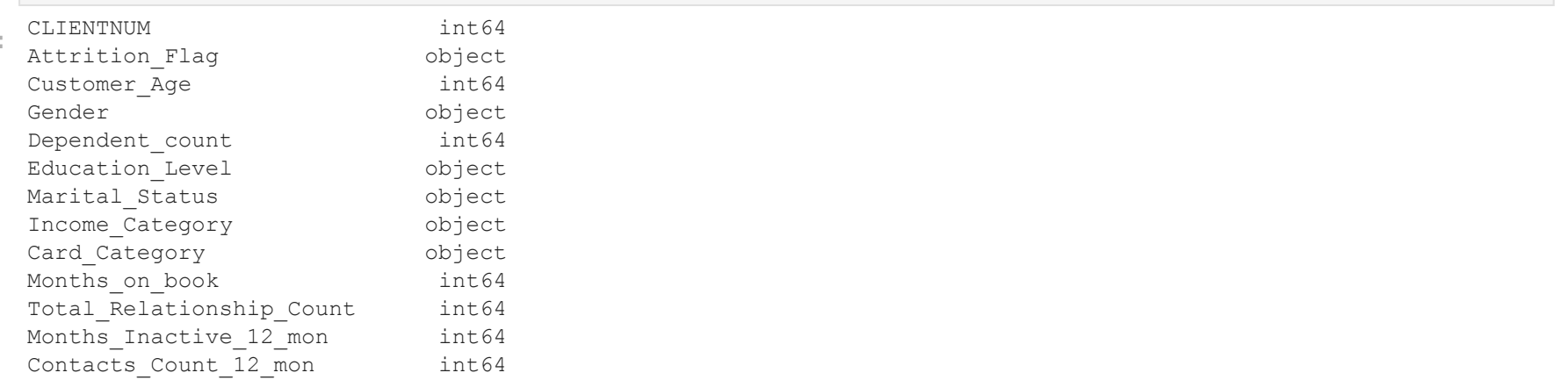
```
Out[13]:
```

	CLIENTNUM	Customer_Age	Dependent_count	Months_on_book	Total_Relationship_Count	Months_inactive_12_mon
CLIENTNUM	1.000000	0.007613	0.006772	0.134588	0.006907	0.005633
Customer_Age	0.007613	1.000000	-0.122254	0.788912	-0.010931	-0.062042
Dependent_count	0.006772	-0.122254	1.000000	-0.103062	-0.039076	-0.009203
Months_on_book	0.134588	0.788912	-0.103062	1.000000	-0.009203	-0.009203
Total_Relationship_Count	0.006907	-0.010931	-0.039076	-0.009203	1.000000	-0.009203
Months_inactive_12_mon	0.005729	0.054361	-0.010768	0.074164	-0.003675	1.000000
Contacts_Count_12_mon	0.005694	-0.018452	-0.040505	-0.010774	0.055203	-0.009203
Credit_Limit	0.005708	0.002476	0.068065	0.007507	-0.071386	-0.009203
Total_Revolving_Bal	0.000825	0.014780	-0.002688	0.008623	0.013726	-0.009203
Avg_Open_To_Buy	0.005633	0.001151	0.068291	0.006732	-0.072601	-0.009203
Total_Amt_Chng_Q4_Q1	0.017369	-0.062042	-0.035439	-0.048959	0.050119	-0.009203
Total_Trans_Amt	-0.019692	-0.046446	0.025046	-0.038591	-0.347229	-0.009203
Total_Trans_Ct	-0.002961	-0.067097	0.049912	-0.049819	-0.241891	-0.009203
Total_Ct_Chng_Q4_Q1	0.007896	-0.012143	0.011087	-0.014072	0.040831	-0.009203
Avg_Utilization_Ratio	0.000266	0.007114	-0.037135	-0.007541	0.067663	-0.009203

Visualizing data

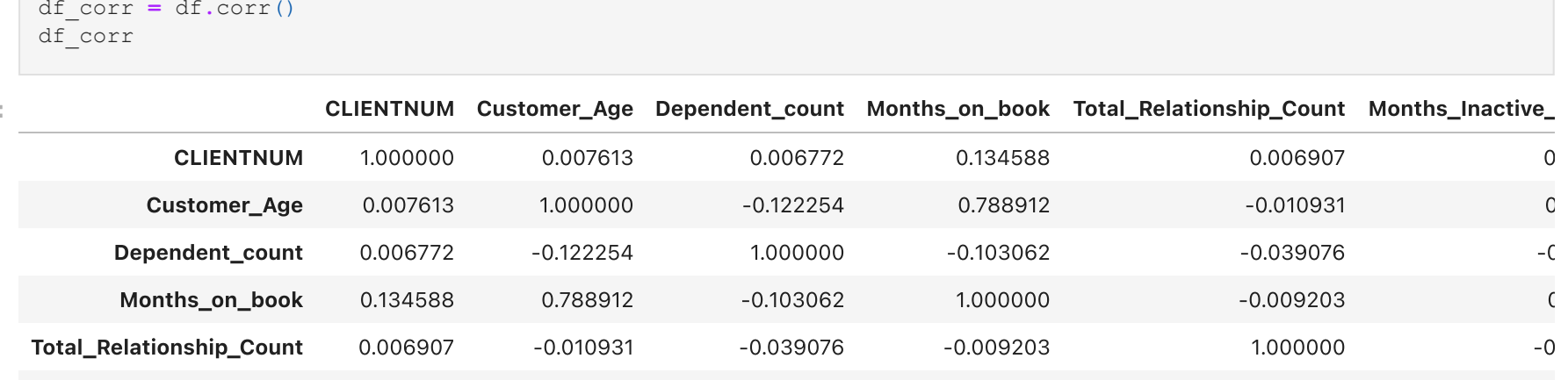
```
In [14]: df['Attrition_Flag'].value_counts().plot.bar()
```

```
Out[14]: <AxesSubplot>
```



```
In [15]: df['Education_Level'].value_counts().plot.bar(color = ['green','blue','magenta','cyan','black','brown','red'])
```

```
Out[15]: <AxesSubplot>
```



```
In [16]: df['Gender'].value_counts().plot.pie()
```

```
Out[16]: <AxesSubplot:label='Gender'>
```



```
In [17]: ins, att = df.shape
ins
```

```
Out[17]: 10127
```

```
In [18]: df.dtypes
```

```
Out[18]: CLIENTNUM      int64
Attrition_Flag      object
Customer_Age        int64
Gender              object
Dependent_count     int64
Education_Level     object
Marital_Status      object
Income_Category     object
Card_Category       object
Months_on_book      int64
Total_Relationship_Count  int64
Months_Inactive_12_mon  int64
Contacts_Count_12_mon  int64
Credit_Limit        float64
Total_Revolving_Bal  float64
Avg_Open_To_Buy     float64
Total_Amt_Chng_Q4_Q1  float64
Total_Trans_Amt     int64
Total_Trans_Ct      int64
Total_Ct_Chng_Q4_Q1  float64
Avg_Utilization_Ratio  float64
dtype: object
```

```
In [19]: #calculating correlation matrix
df_corr = df.corr()
df_corr
```

```
Out[19]:
```

	CLIENTNUM	Customer_Age	Dependent_count	Months_on_book	Total_Relationship_Count	Months_inactive_12_mon
CLIENTNUM	1.000000	0.007613	0.006772	0.134588	0.006907	0.005633
Customer_Age	0.007613	1.000000	-0.122254	0.788912	-0.010931	-0.062042
Dependent_count	0.006772	-0.122254	1.000000	-0.103062	-0.039076	-0.009203
Months_on_book	0.134588	0.788912	-0.103062	1.000000	-0.009203	-0.009203
Total_Relationship_Count	0.006907	-0.010931	-0.039076	-0.009203	1.000000	-0.009203
Months_inactive_12_mon	0.005729	0.054361	-0.010768	0.074164	-0.003675	1.000000
Contacts_Count_12_mon	0.005694	-0.018452	-0.040505	-0.010774	0.055203	-0.009203
Credit_Limit	0.005708	0.002476	0.068065	0.007507	-0.071386	-0.009203
Total_Revolving_Bal	0.000825	0.014780	-0.002688	0.008623	0.013726	-0.009203
Avg_Open_To_Buy	0.005633	0.001151	0.068291	0.006732	-0.072601	-0.009203
Total_Amt_Chng_Q4_Q1	0.017369	-0.062042	-0.035439	-0.048959	0.050119	-0.009203
Total_Trans_Amt	-0.019692	-0.046446	0.025046	-0.038591	-0.347229	-0.009203
Total_Trans_Ct	-0.002961	-0.067097	0.049912	-0.049819	-0.241891	-0.009203
Total_Ct_Chng_Q4_Q1	0.007896	-0.012143	0.011087	-0.014072	0.040831	-0.009203
Avg_Utilization_Ratio	0.000266	0.007114	-0.037135	-0.007541	0.067663	-0.009203

```
In [20]: df1= df.rename(columns={'CLIENTNUM':'1', 'Attrition_Flag':'2', 'Customer_Age':'3', 'Gender':'4',
                                'Dependent_count':'5', 'Education_Level':'6', 'Marital_Status':'7',
                                'Income_Category':'8', 'Card_Category':'9', 'Months_on_book':'10',
                                'Total_Relationship_Count':'11', 'Months_Inactive_12_mon':'12',
                                'Contacts_Count_12_mon':'13', 'Credit_Limit':'14', 'Total_Revolving_Bal':'15',
                                'Avg_Open_To_Buy':'16', 'Total_Amt_Chng_Q4_Q1':'17', 'Total_Trans_Amt':'18',
                                'Total_Trans_Ct':'19', 'Total_Ct_Chng_Q4_Q1':'20', 'Avg_Utilization_Ratio':'21'})
df1.columns
```

```
Out[20]: Index(['1', '2', '3', '4', '5', '6', '7', '8', '9', '10', '11', '12', '13', '14', '15', '16', '17', '18', '19', '20', '21'],
              dtype='object')
```

```
In [21]: df1_corr = df1.corr()
df1_corr
```

```
Out[21]:
```

	1	3	5	10	11	12	13	14	15	16	17
1	1.000000	0.007613	0.006772	0.134588	0.006907	0.005694	0.005708	0.000825	0.005633	0.017369	-0.01
3	0.007613	1.000000	-0.122254	0.788912	-0.010931	0.054361	-0.018452	0.002476	0.014780	0.001151	-0.062042
5	0.006772	-0.122254	1.000000	-0.103062	-0.039076	-0.010768	-0.040505	0.068065	-0.002688	0.068291	-0.035439
10	0.134588	0.788912	-0.103062	1.000000	-0.009203	0.074164	-0.010774	0.008623	0.006732	-0.048959	-0.03
11	0.006907	-0.010931	-0.039076	-0.009203	1.000000	-0.003675	0.055203	-0.071386	0.013726	-0.072601	-0.009203
12	0.005694	-0.018452	-0.040505	-0.010774	0.055203	1.000000	0.029493	-0.020394	-0.042210	-0.016605	-0.032247
13	0.005694	-0.018452	-0.040505	-0.010774	0.055203	0.029493	1.000000	0.020817	-0.053913	0.025646	-0.024445
14	0.005708	0.002476	0.068065	0.007507	-0.071386	-0.020394	0.20817	1.000000	0.042493	0.995981	0.012813
15	0.000825	0.014780	-0.002688	0.008623	0.013726	-0.042210	-0.053913	0.042493	1.000000	-0.047167	0.058174
16	0.005633	0.001151	0.068291	0.006732	-0.072601	-0.016605	0.025646	0.042493	-0.047167	1.000000	0.007995
17	0.017369	-0.062042	-0.035439	-0.048959	0.050119	-0.032247	-0.024445	0.012813	0.058174	0.007995	1.000000
18	-0.019692	-0.046446	0.025046	-0.038591	-0.347229	-0.036982	-0.172774	0.171730	0.064370	0.165923	0.039678
19	-0.002961	-0.067097	0.049912	-0.049819	-0.241891	-0.042787	-0.152213	0.075927	0.056060	0.070885	0.005469
20	0.007896	-0.012143	0.011087	-0.014072	0.040831	-0.039899	-0.094997	-0.002020	0.089861	-0.010076	0.384189
21	0.000266	0.007114	-0.037135	-0.007541	0.067663	-0.007503	-0.005541	-0.482965	0.024022	-0.538808	0.035235

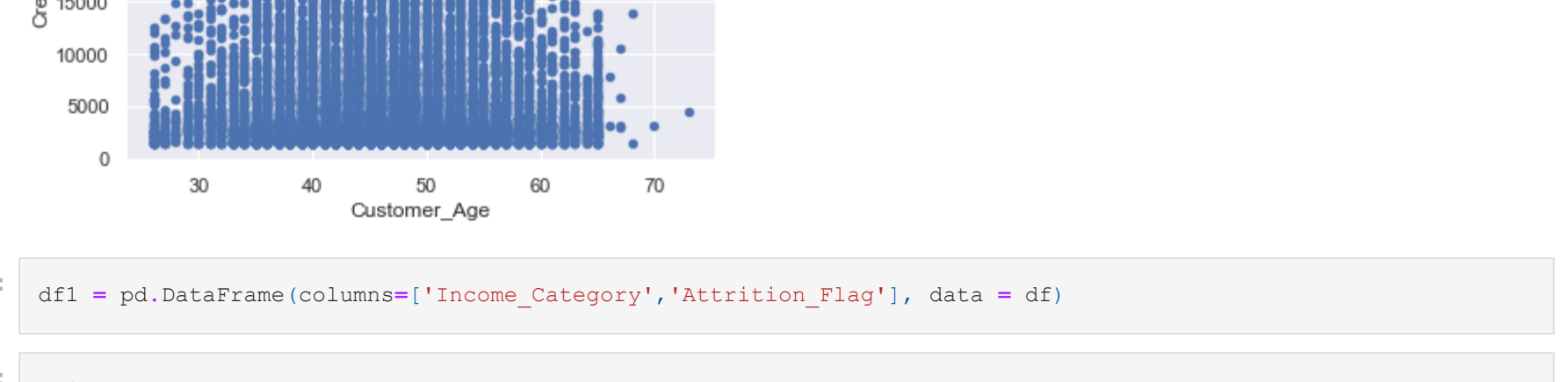
```
In [22]: df.columns

Out[22]:
```

```
Index(['CLIENTNUM', 'Attrition_Flag', 'Customer_Age', 'Gender', 'Dependent_count', 'Education_Level', 'Marital_Status', 'Income_Category', 'Card_Category', 'Months_on_book', 'Total_Relationship_Count', 'Months_Inactive_12_mon', 'Contacts_Count_12_mon', 'Credit_Limit', 'Total_Revolving_Bal', 'Avg_Open_To_Buy', 'Total_Amt_Chng_Q4_Q1', 'Total_Trans_Amt', 'Total_Trans_Ct', 'Total_Ct_Chng_Q4_Q1', 'Avg_Utilization_Ratio'],
      dtype='object')
```

```
In [23]: #plotting heatmap for correlation matrix
df_corr = df.corr()
mask = np.triu(np.ones_like(df_corr, dtype=bool))
f, ax = plt.subplots(figsize=(11, 9))
cmap = sns.diverging_palette(230, 20, as_cmap=True)
sns.heatmap(df_corr, annot=True, mask=mask, cmap=cmap, vmax=3, center=0,
            square=True, linewidths=.5, cbar_kws={"shrink": .5})
```

```
Out[23]: <AxesSubplot>
```

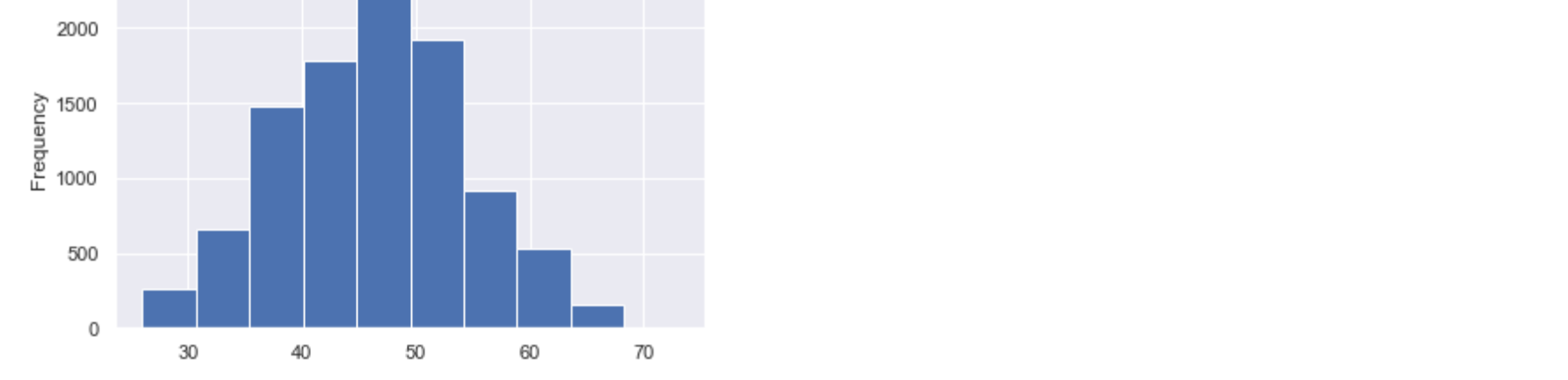


```
In [24]: df['Education_Level'].value_counts()
```

```
Out[24]: Graduate      2013
High School    1519
Unknown        1487
Uneducated     1013
Post-graduate    516
Doctorate       451
Name: Education_Level, dtype: int64
```

```
In [25]: df.plot.scatter('Customer_Age', 'Credit_Limit')
```

```
Out[25]:
```



```
In [26]: df1 = pd.DataFrame(columns=['Income_Category', 'Attrition_Flag'], data = df)
```

```
Out[26]:
```

```
In [27]: df1.head()

Out[27]:
```

	Income_Category	Attrition_Flag
0	60K--80K	Existing Customer
1	Less than \$40K	Existing Customer
2	80K--120K	Existing Customer
3	Less than \$40K	Existing Customer
4	60K--80K	Existing Customer

```
In [28]: df['Income_Category'].value_counts()
```

```
Out[28]: Less than $40K    3561
$40K -- $60K    1790
$60K -- $120K    1539
$60K -- $80K    1402
Unknown    1122
$120K +    727
Name: Income_Category, dtype: int64
```

```
In [29]: df['Customer_Age'].plot(kind='hist')
```

```
Out[29]: <AxesSubplot:label='Frequency'>
```

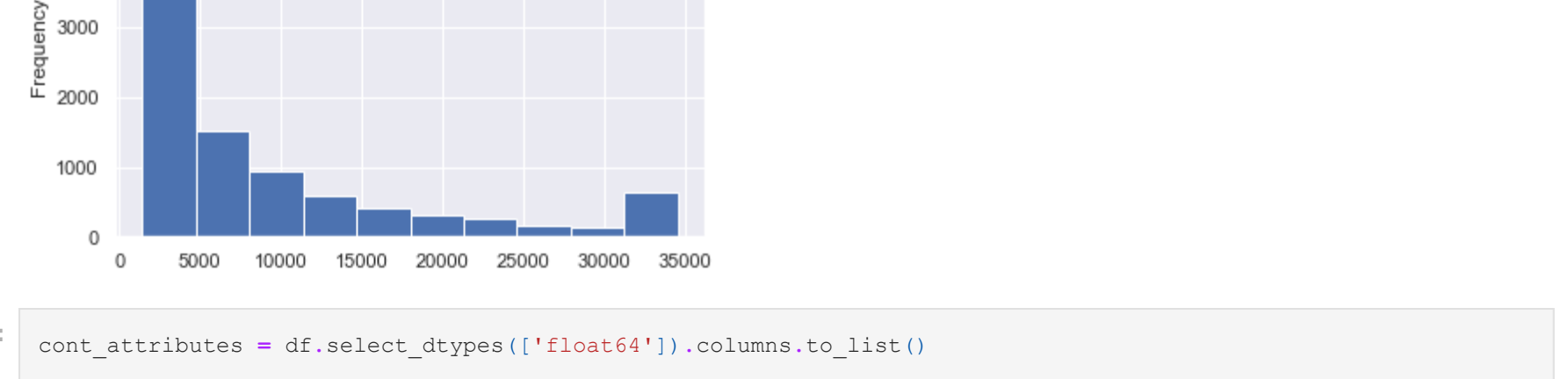


```
In [30]: pd.cut(df['Customer_Age'],bins = 4).value_counts()
```

```
Out[30]: (37.5, 45.5]    5137
(45.5, 51.25]    3218
(51.25, 57.75]    1400
(57.75, 75.0]    312
Name: Customer_Age, dtype: int64
```

```
In [31]: df.loc[df['Customer_Age'] > 65, 'Customer_Age'] = np.mean(df['Customer_Age'])
df['Customer_Age'].plot.box()
```

```
Out[31]: <AxesSubplot>
```



```
In [32]: bins = [0, 35, 45, 55, 70]
cut_labels = ['0-35', '35-45', '45-55', '55-70']
df['Customer_Age'] = pd.cut(df['Customer_Age'], bins=bins, labels=cut_labels)
```

```
In [33]: df.head()
```

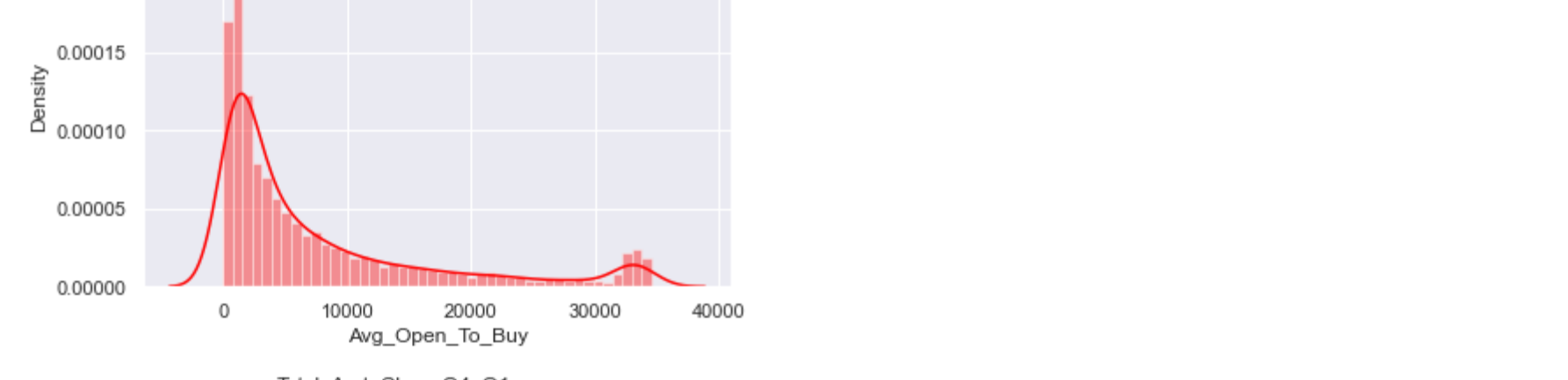
```
Out[33]:
```

	CLIENTNUM	Attrition_Flag	Customer_Age	Gender	Dependent_count	Education_Level	Marital_Status	Income_Category	Card_Ca
0	768805383	Existing Customer	35_45	M	3	High School	Married	60K--80K	
1	818770008	Existing Customer	45_55	F	5	Graduate	Single	Less than \$40K	
2	713982108	Existing Customer	45_55	M	3	Graduate	Married	80K--120K	
3	769911858	Existing Customer	35_45	F	4	High School	Unknown	Less than \$40K	
4	709106358	Existing Customer	35_45	M	3	Uneducated	Married	60K--80K	

```
5 rows x 21 columns

In [34]: df['Customer_Age'].value_counts().plot(kind = 'bar')
```

```
Out[34]: <AxesSubplot>
```



```
In [35]: df['Credit_Limit'].plot(kind = 'hist')
```

```
Out[35]: <AxesSubplot:label='Frequency'>
```



```
In [36]: cont_attributes = df.select_dtypes(['float64']).columns.to_list()
for cont_attribute in cont_attributes:
    plt.figure()
    plt.title(cont_attribute)
    ax = sns.distplot(df[cont_attribute], color = 'red')
```




```
[99]: print("Classification report:")
rep = classification_report(Y_test,Y_test_pred)
print(rep)
print(f"dtree_score={Y_test}")

Classification report:
              precision    recall  f1-score   support

      0               0.96         0.96         0.96         850
      1               0.87         0.97         0.92         163

 accuracy               0.93
 macro avg              0.92
 weighted avg           0.93

Accuracy:
0.93089321816387

Final score of Decision Tree is 93.83% on the test data

KNN

In [100]: from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import f1_score

In [101]: model2 = KNeighborsClassifier(n)

In [102]: model2.fit(X_train,Y_train)
k_pred = model2.predict(X_test)
x_pred

Out[102]: array([1, 0, 0, 0, ..., 0, 0, ..., 0], dtype=object)

In [103]: knn_score = model2.score(X_test,Y_test)

In [104]: def test(k):
for i in k:
test_model = KNeighborsClassifier(n=neighbors)
test_model.fit(X_train,Y_train)
test_pred = test_model.predict(X_test)
test_error.append(np.mean(pred_i != Y_test))
return test_error

In [105]: k = range(1,30)

In [106]: test_model = elbow(k)

In [107]: plt.plot(k,test_model)
plt.xlabel('n_neighbors')
plt.ylabel('error')
plt.show()

Out[107]: Text(0.5, 1.0, 'Elbow curve')

In [108]: knn_score

Out[108]: 0.81639693978283

In [109]: knn_score=knn_score

In [110]: knn_y_train_pred = model2.predict(X_train)
knn_y_test_pred = model2.predict(X_test)

In [111]: knn_cm_train = confusion_matrix(Y_train,knn_y_train_pred,labels=model2.classes_)
knn_cm_test = confusion_matrix(Y_test,knn_y_test_pred,labels=model2.classes_)

In [112]: disp = ConfusionMatrixDisplay(confusion_matrix=knn_cm_test,display_labels=model2.classes_)

In [113]: disp.plot()
plt.grid(False)

In [114]: print("Classification report:")
rep = classification_report(Y_test,knn_y_test_pred)
print(rep)

Classification report:
              precision    recall  f1-score   support

      0               0.84         0.96         0.90         850
      1               0.23         0.06         0.10         163

 accuracy               0.82
 macro avg              0.51
 weighted avg           0.50

In [115]: knn_f1_score = f1_score(model2,X_test,Y_test,scoring='accuracy')
knn_cross_val = knn_cv1.mean()

Out[115]: 0.80055194874355

In [116]: knn_f1_score = (10**0.97,11**0.10)
knn_prec = (10**0.94,11**0.23)

In [117]: knn_estimator = KNeighborsClassifier()

In [118]: parameters_RNN = {'n_neighbors': [1,100, 1], 'p': [1,2], 'weights': ['uniform', 'distance'], 'metric': ['minkowski']}

In [119]: grid_search_RNN = GridSearchCV(
estimator=knn_estimator,
param_grid=parameters_RNN,
scoring = 'accuracy',
n_jobs = -1,
cv = 5
)

In [120]: grid_search_RNN.fit(X_train,Y_train)

Out[120]: GridSearchCV(cv=5, estimator=KNeighborsClassifier(), n_jobs=-1,
param_grid={'metric': ['minkowski'], 'chebychev'},
n_neighbors': [1, 100, 1], 'p': [1, 2],
weights': ['uniform', 'distance']},
scoring='accuracy')

In [121]: grid_search_RNN.best_params_

Out[121]: {'metric': 'minkowski', 'n_neighbors': 100, 'p': 1, 'weights': 'uniform'}

In [122]: grid_search_RNN.best_estimator_

In [123]: KNeighborsClassifier(n_neighbors=100, p=1)

Out[123]: grid_search_RNN.best_score_
0.8393680456997519

Random-Forest classifier

In [124]: from sklearn.ensemble import RandomForestClassifier

In [125]: RFC = RandomForestClassifier()

In [126]: RFC.fit(X_train,Y_train)

Out[126]: RandomForestClassifier()

In [127]: RFC.score(X_test,Y_test)

Out[127]: 0.95859393089321

In [128]: RFC_score = RFC.score(X_test,Y_test)

In [129]: RFC.score(X_train,X_train)

Out[129]: 1.0

In [130]: RFC.feature_importances_

Out[130]: array([0.02798511, 0.01468376, 0.02600532, 0.06236768, 0.02679781,
0.02683043, 0.03668299, 0.10325041, 0.03489332, 0.06135704,
0.17403662, 0.13166084, 0.10031076, 0.05166333, 0.00654652,
0.00407256, 0.00513198, 0.00531398, 0.00380278, 0.00792597, 0.00540371,
0.00231961, 0.00224902, 0.00386414, 0.0035403 , 0.00229237,
0.00323397, 0.00318467, 0.00071038, 0.00587001, 0.00508476,
0.00174717 , 0.00309023, 0.00264331, 0.00294235, 0.00378429,
0.00146572, 0.00186663, 0.00100223, 0.00023083, 0.00160891])

In [131]: feat_importances = pd.Series(RFC.feature_importances_.index*X.columns)

In [132]: import time
import numpy as np
start_time = time.time()
importances = RFC.feature_importances_
std = np.std([tree.feature_importances_ for tree in RFC.estimators_], axis=0)
elapsed_time = time.time() - start_time
print(f"Elapsed time to compute the importances: (elapsed_time:.3f) seconds")

Elapsed time to compute the importances: 0.013 seconds

In [133]: fig, ax = plt.subplots(figsize=(6,6))
feat_importances.nlargest(10).sort_values().plot(kind='barh')
plt.title("Top 10 Important Features")
plt.show()

In [134]: Top 10 Important Features

In [135]: Final score of Random-Forest Classifier on test data is 95.11%

In [136]: rf_y_test_pred = RFC.predict(X_test)

In [137]: rf_y_train_pred = RFC.predict(X_train)

In [138]: rf_cm_train = confusion_matrix(Y_train, rf_y_train_pred,labels=RFC.classes_)
rf_cm_test = confusion_matrix(Y_test, rf_y_test_pred,labels=RFC.classes_)

In [139]: print("Training confusion matrix:")
print(rf_cm_train)

Training confusion matrix:
[[7650 0]
 [ 0 1464]]

In [140]: print("Test confusion matrix:")
print(rf_cm_test)

Test confusion matrix:
[[839 11]
 [ 31 132]]

In [141]: disp = ConfusionMatrixDisplay(confusion_matrix=rf_cm_test,display_labels=RFC.classes_)

In [142]: disp.plot()
plt.grid(False)

In [143]: rf_score = classification_report(Y_test,rf_y_test_pred)

In [144]: print("RFC performance report:")
print(rf_score)

RFC performance report:
              precision    recall  f1-score   support

      0               0.96         0.99         0.98         850
      1               0.92         0.81         0.86         163

 accuracy               0.94
 macro avg              0.90
 weighted avg           0.96

In [145]: rf_f1_score = (10**0.97,11**0.85)
rfc_prec = (10**0.96,11**0.91)

In [146]: rf_cv1 = cross_val_score(RFC,X_test,Y_test,scoring='accuracy')
rf_cross_val = rf_cv1.mean()
rfc_cross_val

Out[146]: 0.9012437204311963

In [147]: rf = RandomForestClassifier(n_jobs=-1,max_features='sqrt',n_estimators=50, oob_score = True)

In [148]: rf_grid = GridSearchCV(estimator=rf, param_grid=param_grid, cv=5)

In [149]: rf_grid.fit(X_train,Y_train)

Out[149]: GridSearchCV(cv=5, estimator=RandomForestClassifier(max_features='sqrt',
n_estimators=50, oob_score=True,
param_grid={'max_features': ('auto', 'sqrt', 'log2'),
n_estimators': [200, 700]}),
scoring='accuracy')

In [150]: rf_grid.best_params_

Out[150]: {'max_features': 'auto', 'n_estimators': 700}

In [151]: rf_grid.best_estimator_

Out[151]: RandomForestClassifier(n_estimators=700, n_jobs=-1, oob_score=True)

In [152]: rf_grid.best_score_

Out[152]: 0.9220896785976

SVM

In [153]: from sklearn import svm

In [154]: model3 = svm.SVC(C=1.0, kernel='rbf', gamma='scale', random_state=None)

In [155]: model3.fit(X_train,Y_train)
svc_y_train_pred = model3.predict(X_train)
svc_y_test_pred = model3.predict(X_test)

In [156]: SVM_score = model3.score(X_test,Y_test)

In [157]: svc_cm_train = confusion_matrix(Y_train, svc_y_train_pred,labels=model3.classes_)
svc_cm_test = confusion_matrix(Y_test, svc_y_test_pred,labels=model3.classes_)

In [158]: disp = ConfusionMatrixDisplay(confusion_matrix=svc_cm_test,display_labels=model3.classes_)

In [159]: disp.plot()
plt.grid(False)

In [160]: svc_score = classification_report(Y_test,svc_y_test_pred)
print("Classification report for SVC:")
print(svc_rep)

Classification report for SVC:
              precision    recall  f1-score   support

      0               0.84         1.00         0.91         850
      1               0.00         0.00         0.00         163

 accuracy               0.42
 macro avg              0.50
 weighted avg           0.46

In [161]: svm_f1_score = (10**0.91,11**0.00)
svm_prec = (10**0.84,11**0.00)

In [162]: svm_cv1 = cross_val_score(svm,X_test,Y_test,scoring='accuracy')
svm_cross_val = svm_cv1.mean()

Out[162]: 0.8390967175353288

Performing GridsearchCV

In [163]: from sklearn.model_selection import GridSearchCV

In [164]: param_grid = {'C': [0.1, 1, 10, 100, 1000], 'gamma': [1, 0.1, 0.01, 0.001, 0.0001], 'kernel': ['rbf']}

In [165]: grid = GridSearchCV(svm.SVC(), param_grid, refit = True, verbose = 3)

In [166]: grid.fit(X_train,Y_train)

Fitting 5 folds for each of 25 candidates, totalling 125 fits
(CV 2/5) END .....C=0.1, gamma=1, kernel=rbf, score=0.839 total time= 4.3s
(CV 3/5) END .....C=0.1, gamma=1, kernel=rbf, score=0.839 total time= 4.2s
(CV 4/5) END .....C=0.1, gamma=1, kernel=rbf, score=0.839 total time= 4.2s
(CV 5/5) END .....C=0.1, gamma=1, kernel=rbf, score=0.840 total time= 4.3s
(CV 2/5) END .....C=0.1, gamma=0.1, kernel=rbf, score=0.839 total time= 4.8s
(CV 3/5) END .....C=0.1, gamma=0.1, kernel=rbf, score=0.839 total time= 4.8s
(CV 4/5) END .....C=0.1, gamma=0.1, kernel=rbf, score=0.839 total time= 4.8s
(CV 5/5) END .....C=0.1, gamma=0.1, kernel=rbf, score=0.839 total time= 4.8s
(CV 2/5) END .....C=0.1, gamma=0.01, kernel=rbf, score=0.839 total time= 5.3s
(CV 3/5) END .....C=0.1, gamma=0.01, kernel=rbf, score=0.839 total time= 4.6s
(CV 4/5) END .....C=0.1, gamma=0.01, kernel=rbf, score=0.839 total time= 4.5s
(CV 5/5) END .....C=0.1, gamma=0.01, kernel=rbf, score=0.839 total time= 4.5s
(CV 2/5) END .....C=0.1, gamma=0.001, kernel=rbf, score=0.839 total time= 4.4s
(CV 3/5) END .....C=0.1, gamma=0.001, kernel=rbf, score=0.839 total time= 4.3s
(CV 4/5) END .....C=0.1, gamma=0.001, kernel=rbf, score=0.839 total time= 4.3s
(CV 5/5) END .....C=0.1, gamma=0.001, kernel=rbf, score=0.839 total time= 4.3s
(CV 2/5) END .....C=1, gamma=1, kernel=rbf, score=0.839 total time= 4.7s
(CV 3/5) END .....C=1, gamma=1, kernel=rbf, score=0.839 total time= 4.7s
(CV 4/5) END .....C=1, gamma=1, kernel=rbf, score=0.839 total time= 4.7s
(CV 5/5) END .....C=1, gamma=1, kernel=rbf, score=0.839 total time= 4.7s
(CV 2/5) END .....C=1, gamma=0.1, kernel=rbf, score=0.840 total time= 4.9s
(CV 3/5) END .....C=1, gamma=0.1, kernel=rbf, score=0.839 total time= 4.6s
(CV 4/5) END .....C=1, gamma=0.1, kernel=rbf, score=0.839 total time= 4.6s
(CV 5/5) END .....C=1, gamma=0.1, kernel=rbf, score=0.839 total time= 4.6s
(CV 2/5) END .....C=1, gamma=0.01, kernel=rbf, score=0.839 total time= 4.4s
(CV 3/5) END .....C=1, gamma=0.01, kernel=rbf, score=0.839 total time= 4.3s
(CV 4/5) END .....C=1, gamma=0.01, kernel=rbf, score=0.839 total time= 4.3s
(CV 5/5) END .....C=1, gamma=0.01, kernel=rbf, score=0.839 total time= 4.3s
(CV 2/5) END .....C=1, gamma=0.001, kernel=rbf, score=0.839 total time= 4.2s
(CV 3/5) END .....C=1, gamma=0.001, kernel=rbf, score=0.839 total time= 4.2s
(CV 4/5) END .....C=1, gamma=0.001, kernel=rbf, score=0.839 total time= 4.2s
(CV 5/5) END .....C=1, gamma=0.001, kernel=rbf, score=0.839 total time= 4.2s
(CV 2/5) END .....C=10, gamma=1, kernel=rbf, score=0.839 total time= 8.2s
(CV 3/5) END .....C=10, gamma=1, kernel=rbf, score=0.839 total time= 8.2s
(CV 4/5) END .....C=10, gamma=1, kernel=rbf, score=0.839 total time= 8.2s
(CV 5/5) END .....C=10, gamma=1, kernel=rbf, score=0.839 total time= 8.2s
(CV 2/5) END .....C=10, gamma=0.1, kernel=rbf, score=0.839 total time= 8.2s
(CV 3/5) END .....C=10, gamma=0.1, kernel=rbf, score=0.839 total time= 8.2s
(CV 4/5) END .....C=10, gamma=0.1, kernel=rbf, score=0.839 total time= 8.2s
(CV 5/5) END .....C=10, gamma=0.1, kernel=rbf, score=0.839 total time= 8.2s
(CV 2/5) END .....C=10, gamma=0.01, kernel=rbf, score=0.839 total time= 8.2s
(CV 3/5) END .....C=10, gamma=0.01, kernel=rbf, score=0.839 total time= 8.2s
(CV 4/5) END .....C=10, gamma=0.01, kernel=rbf, score=0.839 total time= 8.2s
(CV 5/5) END .....C=10, gamma=0.01, kernel=rbf, score=0.839 total time= 8.2s
(CV 2/5) END .....C=10, gamma=0.001, kernel=rbf, score=0.839 total time= 8.2s
(CV 3/5) END .....C=10, gamma=0.001, kernel=rbf, score=0.839 total time= 8.2s
(CV 4/5) END .....C=10, gamma=0.001, kernel=rbf, score=0.839 total time= 8.2s
(CV 5/5) END .....C=10, gamma=0.001, kernel=rbf, score=0.839 total time= 8.2s
(CV 2/5) END .....C=100, gamma=1, kernel=rbf, score=0.839 total time= 7.8s
(CV 3/5) END .....C=100, gamma=1, kernel=rbf, score=0.839 total time= 7.8s
(CV 4/5) END .....C=100, gamma=1, kernel=rbf, score=0.839 total time= 7.8s
(CV 5/5) END .....C=100, gamma=1, kernel=rbf, score=0.839 total time= 7.8s
(CV 2/5) END .....C=100, gamma=0.1, kernel=rbf, score=0.840 total time= 8.2s
(CV 3/5) END .....C=100, gamma=0.1, kernel=rbf, score=0.839 total time= 8.2s
(CV 4/5) END .....C=100, gamma=0.1, kernel=rbf, score=0.839 total time= 8.2s
(CV 5/5) END .....C=100, gamma=0.1, kernel=rbf, score=0.839 total time= 8.2s
(CV 2/5) END .....C=100, gamma=0.01, kernel=rbf, score=0.839 total time= 8.2s
(CV 3/5) END .....C=100, gamma=0.01, kernel=rbf, score=0.839 total time= 8.2s
(CV 4/5) END .....C=100, gamma=0.01, kernel=rbf, score=0.839 total time= 8.2s
(CV 5/5) END .....C=100, gamma=0.01, kernel=rbf, score=0.839 total time= 8.2s
(CV 2/5) END .....C=100, gamma=0.001, kernel=rbf, score=0.839 total time= 8.2s
(CV 3/5) END .....C=100, gamma=0.001, kernel=rbf, score=0.839 total time= 8.2s
(CV 4/5) END .....C=100, gamma=0.001, kernel=rbf, score=0.839 total time= 8.2s
(CV 5/5) END .....C=100, gamma=0.001, kernel=rbf, score=0.839 total time= 8.2s
(CV 2/5) END .....C=1000, gamma=1, kernel=rbf, score=0.839 total time= 8.5s
(CV 3/5) END .....C=1000, gamma=1, kernel=rbf, score=0.839 total time= 8.5s
(CV 4/5) END .....C=1000, gamma=1, kernel=rbf, score=0.839 total time= 8.5s
(CV 5/5) END .....C=1000, gamma=1, kernel=rbf, score=0.839 total time= 8.5s
(CV 2/5) END .....C=1000, gamma=0.1, kernel=rbf, score=0.839 total time= 8.7s
(CV 3/5) END .....C=1000, gamma=0.1, kernel=rbf, score=0.839 total time= 8.7s
(CV 4/5) END .....C=1000, gamma=0.1, kernel=rbf, score=0.839 total time= 8.7s
(CV 5/5) END .....C=1000, gamma=0.1, kernel=rbf, score=0.839 total time= 8.7s
(CV 2/5) END .....C=1000, gamma=0.01, kernel=rbf, score=0.839 total time= 8.7s
(CV 3/5) END .....C=1000, gamma=0.01, kernel=rbf, score=0.839 total time= 8.7s
(CV 4/5) END .....C=1000, gamma=0.01, kernel=rbf, score=0.839 total time= 8.7s
(CV 5/5) END .....C=1000, gamma=0.01, kernel=rbf, score=0.839 total time= 8.7s
(CV 2/5) END .....C=1000, gamma=0.001, kernel=rbf, score=0.839 total time= 8.7s
(CV 3/5) END .....C=1000, gamma=0.001, kernel=rbf, score=0.839 total time= 8.7s
(CV 4/5) END .....C=1000, gamma=0.001, kernel=rbf, score=0.839 total time= 8.7s
(CV 5/5) END .....C=1000, gamma=0.001, kernel=rbf, score=0.839 total time= 8.7s
(CV 2/5) END .....C=1000, gamma=0.0001, kernel=rbf, score=0.839 total time= 8.8s
(CV 3/5) END .....C=1000, gamma=0.0001, kernel=rbf, score=0.839 total time= 8.8s
(CV 4/5) END .....C=1000, gamma=0.0001, kernel=rbf, score=0.839 total time= 8.8s
(CV 5/5) END .....C=1000, gamma=0.0001, kernel=rbf, score=0.839 total time= 8.8s
GridSearchCV(estimator=SVC(),
param_grid={'C': [0.1, 1, 10, 100, 1000],
'gamma': [1, 0.1, 0.01, 0.001, 0.0001],
'kernel': ['rbf']},
verbose=3)

In [167]: print(grid.best_params_)

{'C': 0.1, 'gamma': 1, 'kernel': 'rbf'}

In [168]: print(grid.best_estimator_)

SVC(C=0.1, gamma=1)

Comparing all model scores

In [169]: model_scores = {"Support_Vector_Machine":SVM_score,
"Random_Forest_Classifier":RFC_score,
"k_Nearest_Neighbours":knn_score,
"Decision_Tree":DT_score}

In [170]: keys = model_scores.keys()
values = model_scores.values()

In [171]: plt.figure(figsize=(10,4))
plt.bar(keys,values,width=0.4)
plt.xticks(rotation=90)
plt.show()

In [172]: print('Model Scores')
for m,s in model_scores.items():
print(f'{m} ({f"{s:.2f}"})')

Model Scores
Support_Vector_Machine 0.8390965153011
Random_Forest_Classifier 0.958593893089321
k_Nearest_Neighbours 0.816386963978283
Decision_Tree 0.93089321816387

In [173]: model_scores.keys() = ("Support_Vector_Machine",svm_cross_val,
"Random_Forest_Classifier",rfc_cross_val,
"k_Nearest_Neighbours",knn_cross_val,
"Decision_Tree",dtree_cross_val)

In [174]: keys = model_scores.keys()
values = model_scores.values()

In [175]: plt.figure(figsize=(10,4))
plt.bar(keys,values,width=0.4,color = 'brown')
plt.xticks(rotation=0)
plt.show()

In [176]: plt.figure(figsize=(10,4))
plt.bar(X,axis, label_0, color = 'r',
width = width, edgecolor = 'black',
label='label_0')
figure2 = plt.bar(X,axis + width, label_1, color = 'g',
width = width, edgecolor = 'black',
label='label_1')

plt.grid(linestyle='--')
plt.xticks(X,axis + width/2,X)
plt.legend()
plt.xticks(X,axis,X)
plt.xlabel('Models')
plt.ylabel('F1 measures')
plt.title('F1 measures of models')
plt.legend()
plt.show()

In [177]: F1 scores of models

In [178]: X = ['SVM','KNN','DecisionTree','RFC']
label_0 = [svm_prec['0'],knn_prec['0'],dtree_prec['0'],rfc_prec['0']]
label_1 = [svm_prec['1'],knn_prec['1'],dtree_prec['1'],rfc_prec['1']]
n = 4
X_axis = np.arange(n)
from pylab import rcParams
rcParams['figure.figsize'] = 10, 8

figure1 = plt.bar(X,axis, label_0, color = 'b',
width = width, edgecolor = 'black',
label='label_0')
figure2 = plt.bar(X,axis + width, label_1, color = 'g',
width = width, edgecolor = 'black',
label='label_1')

plt.grid(linestyle='--')
plt.xticks(X,axis + width/2,X)
plt.legend()
plt.xticks(X,axis,X)
plt.xlabel('Precision')
plt.ylabel('F1 measures')
plt.title('Precision of models')
plt.legend()
plt.show()

In [179]: Precision of models

In [180]: RandomForest Classifier performs best comparitively with all other models

In [181]:
```