

Bank Churn Customer Analysis: Customer Segmentation

[Segmenting Bank Customers and Recommend Potential New Products or Services for each Segment]

Objective 1: Preparing the Data for Modeling

Our First Objective is to Prepare the Data for Modeling by Selecting a Subset of Fields, making sure they are Numeric, looking at their Distributions, and Engineering a new feature.

```
In [1]: # Importing the Required Libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

In [2]: # First, make sure your Google Drive is mounted in Colab:
# from google.colab import drive
# drive.mount('/content/drive')

Mounted at /content/drive

In [2]: # churn_curt_info = pd.read_excel("Bank_Churn_Messy.xlsx")
# data = pd.read_csv("/content/drive/MyDrive/Data Analytics & BI Career Path/Batch 2/Python/Bank Churn Customer Project/Bank_Churn.csv")
data = pd.read_csv("Bank_Churn.csv")

Out [2]:
```

	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	15534602	Huigang	619	France	Female	42	2	0.00	1	1	1	101348.88	1
1	15534701	Wu	608	Spain	Female	41	1	83807.86	1	0	1	112542.58	0
2	15619394	Ono	502	France	Female	42	8	159660.80	3	1	0	113931.57	1
3	15719354	Rane	699	France	Female	39	1	0.00	2	0	0	93826.63	0
4	15737889	Mitchell	850	Spain	Female	43	2	125510.82	1	1	1	79084.10	0

```
In [3]: data.info()

Out [3]:
Out [4]:
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	619	France	Female	42	2	0.00	1	1	1	101348.88	1
1	608	Spain	Female	41	1	83807.86	1	0	1	112542.58	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

```
In [5]: data_clean = data_subset.copy()
data_clean.head()
```

```
Out [5]:
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	Exited
0	619	France	Female	42	2	0.00	1	1	1	101348.88	1
1	608	Spain	Female	41	1	83807.86	1	0	1	112542.58	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

```
In [6]: data_clean.Geography.value_counts()

Out [6]:
```

Geography	count
France	5014
Germany	2509
Spain	2473

```
dtype: int64

In [7]: data_clean.Gender.value_counts()

Out [7]:
```

Gender	count
Male	5457
Female	4543

```
dtype: int64

In [8]: data_clean.Gender = np.where(data_clean.Gender == 'Female', 1, 0)
data_clean.head()
```

```
Out [8]:
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	France	Germany	Spain
0	619	France	1	42	2	0.00	1	1	1	101348.88	1	0	0
1	608	Spain	1	41	1	83807.86	1	0	1	112542.58	0	0	1
2	502	France	1	42	8	159660.80	3	1	0	113931.57	1	0	0
3	699	France	1	39	1	0.00	2	0	0	93826.63	1	0	0
4	850	Spain	1	43	2	125510.82	1	1	1	79084.10	0	0	1

```
In [12]: data_clean = pd.get_dummies(data_clean, columns = ['Geography'], dtype = 'int', prefix = '', prefix_sep = '')
data_clean.head()
```

```
Out [12]:
```

	CreditScore	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	France	Germany	Spain	
0	619	1	42	2	0.00	1	1	1	1	101348.88	1	0	0
1	608	1	41	1	83807.86	1	0	1	1	112542.58	0	0	1
2	502	1	42	8	159660.80	3	1	0	1	113931.57	1	0	0
3	699	1	39	1	0.00	2	0	0	0	93826.63	1	0	0
4	850	1	43	2	125510.82	1	1	1	1	79084.10	0	0	1

Let's keep a subset of the Data

Create a DataFrame containing all fields except "CustomerId", "Surname" and "Exited".

```
In [4]: data_subset = data[['CreditScore', 'Geography', 'Gender',
                           'Age', 'Tenure', 'Balance', 'NumOfProducts',
                           'HasCrCard', 'IsActiveMember', 'EstimatedSalary']]
data_subset.head()
```

```
Out [4]:
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
0	619	France	Female	42	2	0.00	1	1	1	101348.88
1	608	Spain	Female	41	1	83807.86	1	0	1	112542.58
2	502	France	Female	42	8	159660.80	3	1	0	113931.57
3	699	France	Female	39	1	0.00	2	0	0	93826.63
4	850	Spain	Female	43	2	125510.82	1	1	1	79084.10

Making Text Fields Numeric

```
In [5]: data_clean = data_subset.copy()
data_clean.head()
```

```
Out [5]:
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary
0	619	France	Female	42	2	0.00	1	1	1	101348.88
1	608	Spain	Female	41	1	83807.86	1	0	1	112542.58
2	502	France	Female	42	8	159660.80	3	1	0	113931.57
3	699	France	Female	39	1	0.00	2	0	0	93826.63
4	850	Spain	Female	43	2	125510.82	1	1	1	79084.10

```
In [6]: data_clean.Geography.value_counts()

Out [6]:
```

Geography	count
France	5014
Germany	2509
Spain	2473

```
dtype: int64

In [7]: data_clean.Gender.value_counts()

Out [7]:
```

Gender	count
Male	5457
Female	4543

```
dtype: int64

In [8]: data_clean.Gender = np.where(data_clean.Gender == 'Female', 1, 0)
data_clean.head()
```

```
Out [8]:
```

	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	France	Germany	Spain
0	619	France	1	42	2	0.00	1	1	1	101348.88	1	0	0
1	608	Spain	1	41	1	83807.86	1	0	1	112542.58	0	0	1
2	502	France	1	42	8	159660.80	3	1	0	113931.57	1	0	0
3	699	France	1	39	1	0.00	2	0	0	93826.63	1	0	0
4	850	Spain	1	43	2	125510.82	1	1	1	79084.10	0	0	1

```
In [12]: data_clean = pd.get_dummies(data_clean, columns = ['Geography'], dtype = 'int', prefix = '', prefix_sep = '')
data_clean.head()
```

```
Out [12]:
```

	CreditScore	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	France	Germany	Spain	
0	619	1	42	2	0.00	1	1	1	1	101348.88	1	0	0
1	608	1	41	1	83807.86	1	0	1	1	112542.58	0	0	1
2	502	1	42	8	159660.80	3	1	0	1	113931.57	1	0	0
3	699	1	39	1	0.00	2	0	0	0	93826.63	1	0	0
4	850	1	43	2	125510.82	1	1	1	1	79084.10	0	0	1

Exploring the Data

Exploring the Data by looking at the Min/Max values and the Distribution of each Column.

```
In [13]: data_clean.describe().round(2)

Out [13]:
```

	CreditScore	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	France	Germany	Spain
count	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0
mean	651.0	0.0	39.0	5.0	76486.0	2.0	1.0	1.0	100090.0	1.0	0.0	0.0
std	87.0	0.0	10.0	3.0	62397.0	1.0	0.0	0.0	57510.0	1.0	0.0	0.0
min	350.0	0.0	18.0	0.0	0.0	1.0	0.0	0.0	12.0	0.0	0.0	0.0
25%	584.0	0.0	32.0	3.0	0.0	1.0	0.0	0.0	51002.0	0.0	0.0	0.0
50%	652.0	0.0	37.0	5.0	97199.0	1.0	1.0	1.0	100194.0	1.0	0.0	0.0
75%	718.0	1.0	44.0	7.0	127644.0	2.0	1.0	1.0	149388.0	1.0	1.0	0.0
max	850.0	1.0	92.0	10.0	250898.0	4.0	1.0	1.0	199992.0	1.0	1.0	1.0

```
In [14]: sns.pairplot(data_clean, corner = True);
```

Let's Engineer a New Feature

Engineer a new feature called "ProductsPerYear".

```
In [15]: data_clean['ProductsPerYear'] = np.where(data_clean.Tenure == 0, data_clean.NumOfProducts, data_clean.NumOfProducts / data_clean.Tenure)
data_clean.head()
```

```
Out [15]:
```

	CreditScore	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	France	Germany	Spain	ProductsPerYear
0	619	1	42	2	0.00	1	1	1	101348.88	1	0	0	0.500
1	608	1	41	1	83807.86	1	0	1	112542.58	0	0	1	1.000
2	502	1	42	8	159660.80	3	1	0	113931.57	1	0	0	0.375
3	699	1	39	1	0.00	2	0	0	93826.63	1	0	0	2.000
4	850	1	43	2	125510.82	1	1	1	79084.10	0	0	1	0.500

```
In [16]: data_clean.describe().round(2)

Out [16]:
```

	CreditScore	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	France	Germany	Spain	ProductsPerYear
count	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0
mean	651.0	0.0	39.0	5.0	76486.0	2.0	1.0	1.0	100090.0	1.0	0.0	0.0	1.000
std	87.0	0.0	10.0	3.0	62397.0	1.0	0.0	0.0	57510.0	1.0	0.0	0.0	1.0
min	350.0	0.0	18.0	0.0	0.0	1.0	0.0	0.0	12.0	0.0	0.0	0.0	0.0
25%	584.0	0.0	32.0	3.0	0.0	1.0	0.0	0.0	51002.0	0.0	0.0	0.0	0.0
50%	652.0	0.0	37.0	5.0	97199.0	1.0	1.0	1.0	100194.0	1.0	0.0	0.0	0.0
75%	718.0	1.0	44.0	7.0	127644.0	2.0	1.0	1.0	149388.0	1.0	1.0	0.0	1.0
max	850.0	1.0	92.0	10.0	250898.0	4.0	1.0	1.0	199992.0	1.0	1.0	1.0	4.0

Objective 2: Cluster the Customers (Round 1)

Our Second Objective is to Segment the Customers using K-Means Clustering, including Standardizing the Data, Creating an Inertia Plot, and interpreting the Clusters.

Scale the Data using Standardization

Standardize the data so that each column has a Mean of 0 and Standard Deviation of 1.

```
In [19]: from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
data_scaled = pd.DataFrame(scaler.fit_transform(data_clean), columns = data_clean.columns)
data_scaled.head()
```

```
Out [19]:
```

	CreditScore	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	France	Germany	Spain	ProductsPerYear
0	-0.326221	1.099988	0.293517	-1.041760	-1.225848	-0.911583	0.646092	0.970243	0.021886	0.997204	-0.578736	-0.578909	-0.016781
1	-0.404003	1.099988	0.198164	-1.387538	0.117350	-0.911583	-1.547768	0.970243	0.216534	-1.002804	-0.578736	1.742740	0.967674
2	-1.536794	1.099988	0.293517	1.032908	1.333053	2.527057	0.646092	-1.030670	0.240687	0.997204	-0.578736	-0.578909	-0.262994
3	0.501521	1.099988	0.007457	-1.387538	-1.225848	0.807737	-1.547768	-1.030670	-0.108918	0.997204	-0.578736	-0.578909	2.936584
4	2.063884	1.099988	0.388971	-1.041760	0.785728	-0.911583	0.646092	0.970243	-0.365276	-1.002804	-0.578736	1.742740	-0.016781

```
In [20]: df_scaled.describe().round(1)

Out [20]:
```

	CreditScore	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	EstimatedSalary	France	Germany	Spain	ProductsPerYear
count	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0	10000.0
mean	-0.0	0.0	0.0	0.0	0.0	0.0	-0.0	0.0	-0.0	-0.0	0.0	-0.0	0.0
std	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
min	-3.0	-1.0	-2.0	-2.0	-1.0	-1.0	-2.0	-1.0	-2.0	-1.0	-1.0	-1.0	-1.0
25%	-1.0	-1.0	-1.0	-1.0	-1.0	-1.0	-2.0	-1.0	-1.0	-1.0	-1.0	-1.0	-1.0
50%	0.0	-1.0	-0.0	-0.0	0.0	-1.0	1.0	1.0	0.0	1.0	-1.0	-1.0	0.0
75%	1.0	1.0	0.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	2.0	1.0	0.0
max	2.0	1.0	5.0	2.0	3.0	4.0	1.0	1.0	2.0	1.0	2.0	2.0	7.0

Fit K-Means Models with 2-15 Clusters

Fit K-Means Clustering Models on the Standardized Data with 2-15 Clusters to create an Inertia Plot.

```
In [20]: # Import kmeans and write a loop to fit models with 2 to 15 clusters
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score

# Create an empty list to hold many Inertia and Silhouette Values
inertia
```