

1. Data Preprocessing

```
In [1]: Import pandas as pd
```

1.1 Loading Data & Initial data exploration

```
In [2]: data=pd.read_csv["C:\\Users\\Nig\\Downloads\\churn_large_dataset.xlsx"]
data

Out[2]:
```

	CustomerID	Name	Age	Gender	Location	Subscription_Length_Months	Monthly_Bill	Total_Usage_GB	Churn
0	1	Customer_1	63	Male	Los Angeles	17	73.36	236	0
1	2	Customer_2	62	Female	New York	1	48.76	172	0
2	3	Customer_3	24	Female	Los Angeles	5	85.47	460	0
3	4	Customer_4	36	Female	Miami	3	97.94	297	1
4	5	Customer_5	46	Female	Miami	19	58.14	266	0
...
9995	9996	Customer_9996	33	Male	Houston	23	55.13	226	1
9996	9997	Customer_9997	62	Female	New York	19	61.05	351	0
9997	9998	Customer_9998	64	Male	Chicago	17	96.11	251	1
9998	9999	Customer_9999	51	Female	New York	20	49.25	434	1
9999	10000	Customer_10000	27	Female	Los Angeles	19	76.57	173	1
10000	rows > 9 columns								

```
In [3]: data.head(5)

Out[3]:
```

	CustomerID	Name	Age	Gender	Location	Subscription_Length_Months	Monthly_Bill	Total_Usage_GB	Churn
0	1	Customer_1	63	Male	Los Angeles	17	73.36	236	0
1	2	Customer_2	62	Female	New York	1	48.76	172	0
2	3	Customer_3	24	Female	Los Angeles	5	85.47	460	0
3	4	Customer_4	36	Female	Miami	3	97.94	297	1
4	5	Customer_5	46	Female	Miami	19	58.14	266	0
5	6	Customer_6	37	Male	New York	15	85.85	496	1
6	7	Customer_7	30	Female	Chicago	3	72.79	299	0
7	8	Customer_8	67	Female	Miami	1	97.70	395	1
8	9	Customer_9	30	Female	Miami	10	42.45	150	1
9	10	Customer_10	53	Female	Los Angeles	12	64.49	383	1

```
In [4]: data.describe()

Out[4]:
```

	CustomerID	Age	Subscription_Length_Months	Monthly_Bill	Total_Usage_GB	Churn
count	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000
mean	5000.000000	44.070200	12.400100	65.051907	274.263650	0.497700
std	2867.607700	15.963000	6.430601	20.208866	139.450361	0.498966
min	1.000000	16.000000	1.000000	30.000000	50.000000	0.000000
25%	2500.000000	31.000000	6.000000	47.540000	161.000000	0.000000
50%	5000.000000	44.000000	12.000000	65.000000	274.000000	0.000000
75%	7500.000000	57.000000	19.000000	81.500000	387.000000	1.000000
max	10000.000000	70.000000	24.000000	100.000000	500.000000	1.000000

```
In [5]: data.dtypes

Out[5]:
CustomerID      int64
Name            object
Age             int64
Gender          object
Location        object
Subscription_Length_Months  int64
Monthly_Bill     float64
Total_Usage_GB   int64
Churn           object
dtype: object

In [6]: data["Churn"].value_counts()

Out[6]:
0    55221
1    48779
Name: Churn, dtype: int64

In [7]: data.nunique()

Out[7]:
CustomerID      10000
Name             10000
Age               53
Gender            24
Location         10
Subscription_Length_Months  2
Monthly_Bill     7961
Total_Usage_GB   451
Churn            2
dtype: int64

1.2 Handling missing data

In [8]: data.isna().sum()

Out[8]:
CustomerID      0
Name             0
Age             0
Gender          0
Location        0
Subscription_Length_Months  0
Monthly_Bill    0
Total_Usage_GB  0
Churn           0
dtype: int64

1.3 Encoding Categorical variable

In [9]: data["Gender"] = data["Gender"].apply(lambda x: 1 if x == "Male" else 0)
data

Out[9]:
```

	CustomerID	Name	Age	Gender	Location	Subscription_Length_Months	Monthly_Bill	Total_Usage_GB	Churn
0	1	Customer_1	63	1	Los Angeles	17	73.36	236	0
1	2	Customer_2	62	0	New York	1	48.76	172	0
2	3	Customer_3	24	0	Los Angeles	5	85.47	460	0
3	4	Customer_4	36	0	Miami	3	97.94	297	1
4	5	Customer_5	46	0	Miami	19	58.14	266	0
...
9995	9996	Customer_9996	33	1	Houston	23	55.13	226	1
9996	9997	Customer_9997	62	0	New York	19	61.05	351	0
9997	9998	Customer_9998	64	1	Chicago	17	96.11	251	1
9998	9999	Customer_9999	51	0	New York	20	49.25	434	1
9999	10000	Customer_10000	27	0	Los Angeles	19	76.57	173	1
10000	rows > 9 columns								

```
In [10]: from sklearn.preprocessing import LabelEncoder

label_encoder=LabelEncoder()
data["LocationNew"] = label_encoder.fit_transform(data["Location"])
data

Out[10]:
```

	CustomerID	Name	Age	Gender	Location	Subscription_Length_Months	Monthly_Bill	Total_Usage_GB	Churn	LocationNew
0	1	Customer_1	63	1	Los Angeles	17	73.36	236	0	2
1	2	Customer_2	62	0	New York	1	48.76	172	0	4
2	3	Customer_3	24	0	Los Angeles	5	85.47	460	0	2
3	4	Customer_4	36	0	Miami	3	97.94	297	1	3
4	5	Customer_5	46	0	Miami	19	58.14	266	0	3
...
9995	9996	Customer_9996	33	1	Houston	23	55.13	226	1	1
9996	9997	Customer_9997	62	0	New York	19	61.05	351	0	4
9997	9998	Customer_9998	64	1	Chicago	17	96.11	251	1	0
9998	9999	Customer_9999	51	0	New York	20	49.25	434	1	4
9999	10000	Customer_10000	27	0	Los Angeles	19	76.57	173	1	2
10000	rows > 10 columns									

```
In [11]: city_names = data["LocationNew"].apply(lambda x: str(x).split(',')[0].strip()).unique()
city_names

Out[11]: array(['2', '4', '3', '0', '1'], dtype=object)

In [12]: city_names = data["Location"].apply(lambda x: str(x).split(',')[0].strip()).unique()
city_names

Out[12]: array(['Los Angeles', 'New York', 'Miami', 'Chicago', 'Houston'], dtype=object)

2. Feature Engineering

In [13]: data.drop(['Name', 'Location'], axis=1, inplace=True)
data

Out[13]:
```

	CustomerID	Age	Gender	Subscription_Length_Months	Monthly_Bill	Total_Usage_GB	Churn	LocationNew
0	1	63	1	17	73.36	236	0	2
1	2	62	0	1	48.76	172	0	4
2	3	24	0	5	85.47	460	0	2
3	4	36	0	3	97.94	297	1	3
4	5	46	0	19	58.14	266	0	3
...
9995	9996	33	1	23	55.13	226	1	1
9996	9997	62	0	19	61.05	351	0	4
9997	9998	64	1	17	96.11	251	1	0
9998	9999	51	0	20	49.25	434	1	4
9999	10000	27	0	19	76.57	173	1	2
10000	rows > 8 columns							

```
In [14]: data["bill_X_GB"] = data["Monthly_Bill"] / data["Total_Usage_GB"]
data["bill_Y_sulen"] = data["Monthly_Bill"] / data["Subscription_Length_Months"]
data["subs_X_bill"] = data["Subscription_Length_Months"] / data["Monthly_Bill"]
data["GB_Y_bill"] = data["Total_Usage_GB"] / data["Monthly_Bill"]
data["GB_X_sulen"] = data["Subscription_Length_Months"] / data["Monthly_Bill"]
data

Out[14]:
```

	CustomerID	Age	Gender	Subscription_Length_Months	Monthly_Bill	Total_Usage_GB	Churn	LocationNew	bill_X_GB	bill_Y_sulen	subs_X_bill	GB_Y_bill	GB_X_sulen
0	1	63	1	17	73.36	236	0	2	1.712136	4.332504	0.231734	3.217002	1.712136
1	2	62	0	1	48.76	172	0	4	0.368172	46.760000	0.020550	5.017460	0.368172
2	3	24	0	5	85.47	460	0	2	38.6160	17.084000	0.056560	6.382005	38.6160
3	4	36	0	3	97.94	297	1	3	29.0818	32.646667	0.030361	3.021469	29.0818
4	5	46	0	19	58.14	266	0	3	15.46524	0.320779	4.575103	15.46524	15.46524
...
9995	9996	33	1	23	55.13	226	1	1	1.349318	2.386567	0.417199	4.096461	1.349318
9996	9997	62	0	19	61.05	351	0	4	2.638935	3.244737	0.306105	5.893451	2.638935
9997	9998	64	1	17	96.11	251	1	0	24.12341	5.613529	0.176851	2.411551	24.12341
9998	9999	51	0	20	49.25	434	1	4	2.127450	2.482200	0.460391	8.832283	2.127450
9999	10000	27	0	19	76.57	173	1	2	1.026165	4.090000	0.248139	3.280171	1.026165
10000	rows > 14 columns												

2.1 Scaling & Normalization

```
In [15]: from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()
normalized_data = scaler.fit_transform(data)
normalized_df = pd.DataFrame(normalized_data, columns=data.columns)
normalized_df

Out[15]:
```

	CustomerID	Age	Gender	Subscription_Length_Months	Monthly_Bill	Total_Usage_GB	Churn	LocationNew	bill_X_GB	bill_Y_sulen	subs_X_bill	GB_Y_bill	GB_X_sulen
0	0.00000	0.980385	1.0	0.85602	0.648429	0.413321	0.0	0.50	0.320696	0.010403	0.260709	0.518873	0.320696
1	0.00000	0.980454	0.0	0.00000	0.761429	0.070111	0.0	1.00	0.148857	0.461296	0.023504	0.181857	0.148857
2	0.00000	0.311385	0.0	0.179613	0.794429	0.911111	0.0	0.50	0.780220	0.196475	0.045411	0.303907	0.780220
3	0.00000	0.314554	0.0	0.068607	0.970571	0.548889	1.0	0.75	0.589124	0.318802	0.028122	0.157223	0.589124
4	0.00004	0.534862	0.0	0.762009	0.403000	0.480000	0.0	0.75	0.287560	0.018129	0.401143	0.218664	0.287560
...
9995	0.99996	0.288462	1.0	0.85602	0.393900	0.391111	1.0	0.25	0.25982	0.018163	0.516469	0.222324	0.25982
9996	0.99997	0.984154	0.0	0.762009	0.401241	0.608889	0.0	1.00	0.413383	0.022220	0.377383	0.322220	0.413383
9997	0.99998	0.984615	1.0	0.85602	0.944429	0.449527	1.0	0.50	0.468660	0.044996	0.213131	0.131000	0.468660
9998	0.99999	0.984615	0.0	0.00000	0.286000	0.070111	1.0	1.00	0.469002	0.010777	0.043566	0.010742	0.469002
9999	1.00000	0.179377	0.0	0.762009	0.406126	0.273333	1.0	0.50	0.242176	0.000140	0.301540	0.010153	0.242176
10000	rows > 14 columns												

2.2 Finding Correlation

```
In [16]: corr_matrix = normalized_df.corr()
corr_matrix["Churn"].sort_values(ascending=False)

Out[16]:
Churn      1.000000
LocationNew 0.404605
subs_X_bill 0.403366
Subscription_Length_Months 0.402339
Gender      0.002121
Age         0.001559
bill_X_sulen 0.000761
Monthly_Bill 0.000711
bill_Y_sulen 0.000710
Total_Usage_GB 0.000710
bill_X_GB   0.000369
GB_Y_bill   0.000368
GB_X_sulen  0.000368
CustomerID  0.000368
Name: Churn, dtype: float64

In [17]: import seaborn as plt
sns.heatmap(corr_matrix)

Out[17]:
```

```
In [18]: from pandas.plotting import scatter_matrix
attributes = ["Subscription_Length_Months", "Monthly_Bill", "Total_Usage_GB", "Churn", "LocationNew"]
scatter_matrix(attributes, figsize=(10,10), alpha=0.3)

Out[18]:
```

```
In [19]: data.plot(kind='scatter', x='LocationNew', y='Churn', alpha=0)

Out[19]:
```

```
In [20]: city = data.loc[data["Churn"]==1, "LocationNew"]
city.groupby(city).size()

Out[20]:
city
0    2
1    2
2    4
3    2
4    2
Name: LocationNew, dtype: int64

In [21]: import numpy as np
print("Has infinite values", np.any(np.isinf(normalized_df[["bill_X_GB", "bill_Y_sulen", "subs_X_bill"]]))
Has infinite values: False

In [22]: normalized_df.isna().sum()

Out[22]:
CustomerID      0
Age             0
Gender          0
Subscription_Length_Months  0
Monthly_Bill    0
Total_Usage_GB  0
Churn           0
LocationNew     0
bill_X_GB       0
bill_Y_sulen    0
subs_X_bill     0
GB_Y_bill       0
GB_X_sulen      0
dtype: int64

In [23]: normalized_df.plot(kind='scatter', x='subs_X_bill', y='bill_X_GB', alpha=0)

Out[23]:
```

```
In [24]: corr_matrix_2 = normalized_df.corr()
corr_matrix_2["Churn"].sort_values(ascending=False)

Out[24]:
Churn      1.000000
LocationNew 0.404605
subs_X_bill 0.403366
Subscription_Length_Months 0.402339
Gender      0.002121
Age         0.001559
bill_X_sulen 0.000761
Monthly_Bill 0.000711
bill_Y_sulen 0.000710
Total_Usage_GB 0.000710
bill_X_GB   0.000369
GB_Y_bill   0.000368
GB_X_sulen  0.000368
CustomerID  0.000368
Name: Churn, dtype: float64

2.3 Splitting Train & Test data

In [25]: from sklearn.model_selection import train_test_split

x = normalized_df[["Subscription_Length_Months", "Monthly_Bill", "LocationNew", "Total_Usage_GB", "bill_X_GB"]]
y = normalized_df["Churn"]

X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=42)
print(X_train.shape)
print(X_test.shape)
print(y_train.shape)
print(y_test.shape)

Out[25]:
70000
30000
70000
30000

In [26]: import numpy as np
print("Has infinite values", np.any(np.isinf(X)))
Has infinite values: False

In [27]: normalized_df.isna().sum()

Out[27]:
CustomerID      0
Age             0
Gender          0
Subscription_Length_Months  0
Monthly_Bill    0
Total_Usage_GB  0
Churn           0
LocationNew     0
bill_X_GB       0
bill_Y_sulen    0
subs_X_bill     0
GB_Y_bill       0
GB_X_sulen      0
dtype: int64

In [28]: normalized_df.plot(kind='scatter', x='subs_X_bill', y='bill_X_GB', alpha=0)

Out[28]:
```

```
In [29]: corr_matrix_2 = normalized_df.corr()
corr_matrix_2["Churn"].sort_values(ascending=False)

Out[29]:
Churn      1.000000
LocationNew 0.404605
subs_X_bill 0.403366
Subscription_Length_Months 0.402339
Gender      0.002121
Age         0.001559
bill_X_sulen 0.000761
Monthly_Bill 0.000711
bill_Y_sulen 0.000710
Total_Usage_GB 0.000710
bill_X_GB   0.000369
GB_Y_bill   0.000368
GB_X_sulen  0.000368
CustomerID  0.000368
Name: Churn, dtype: float64

2.3 Splitting Train & Test data

In [30]: from sklearn.model_selection import train_test_split

x = normalized_df[["Subscription_Length_Months", "Monthly_Bill", "LocationNew", "Total_Usage_GB", "bill_X_GB"]]
y = normalized_df["Churn"]

X
```


