

1. DataPreprocessing

```
In [5]: import pandas as pd
```

1.1 Loading Data & Initial data exploration

```
In [1]: data = pd.read_excel('C:\Users\Ug\Downloads\Customer_churn_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 x 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	67	Male	New York	15	82.95	456	1
6	7	Customer_7	20	Female	Chicago	3	72.75	283	0
7	8	Customer_8	67	Female	Miami	1	97.70	150	1
8	9	Customer_9	20	Female	Miami	10	42.45	365	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
count	10000.000000	10000.000000	10000.000000	10000.000000	10000.000000
mean	5000.000000	44.027020	12.460100	69.953187	274.393902
std	5000.000000	15.852025	6.150041	21.293968	119.459363
min	1	10000.00	18.000000	1.000000	30.000000
25%	2500.000000	31.000000	6.000000	47.540000	161.000000
50%	5000.000000	44.000000	12.000000	69.000000	274.000000
75%	7500.000000	57.000000	19.000000	82.940000	367.000000
max	10000.000000	70.000000	24.000000	100.000000	500.000000

```
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:	int64

```
In [6]: data['Churn'].value_counts()
```

```
Out[6]:
```

0	55221
1	49779

Name: Churn, dtype: int64

```
In [7]: data.nunique()
```

```
Out[7]:
```

CustomerID	100000
Name	100000
Age	93
Gender	2
Location	9
Subscription_Length_Months	24
Monthly_Bill	7901
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)
```

```
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 x 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 x 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', '5'], 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 x 8 columns
```

```
In [14]: data['bill_X_GB'] = data['Monthly_Bill']/data['Total_Usage_GB']
data['bill_X_sublen'] = data['Monthly_Bill']/data['Subscription_Length_Months']
data['subs_X_bill'] = data['Subscription_Length_Months']/data['Monthly_Bill']
data['GB_X_bill'] = data['Total_Usage_GB']/data['Monthly_Bill']
data
```

```
Out[14]:
```

	CustomerID	Age	Gender	Subscription_Length_Months	Monthly_Bill	Total_Usage_GB	Churn	LocationNew	bill_X_GB	bill_X_sublen	subs_X_bill	GB_X_bill	GB_X_sublen
0	1	63	1	17	73.36	236	0	2	1732.96	0.231734	3.211022	1732.96	1732.96
1	2	62	0	1	48.76	172	0	4	8386.72	48.760000	0.020938	3.527402	8386.72
2	3	24	0	5	85.47	460	0	2	3614.20	17.044000	0.058660	5.382000	3614.20
3	4	36	0	3	97.94	297	1	3	20808.18	32.646667	0.030201	3.026469	20808.18
4	5	46	0	19	58.14	266	0	3	15465.24	3.060000	0.320797	4.751563	15465.24
...
9995	9996	33	1	23	55.13	226	1	1	12493.38	2.390567	0.417126	4.096461	12493.38
9996	9997	62	0	19	61.05	351	0	4	23819.23	2.244977	0.020016	5.495451	23819.23
9997	9998	64	1	17	96.11	251	1	0	2413.61	5.633200	0.170861	2.611361	2413.61
9998	9999	51	0	20	49.25	434	1	4	21374.60	2.483200	0.400091	8.812189	21374.60
9999	10000	27	0	19	76.57	173	1	2	12046.61	4.030000	0.248139	2.558971	12046.61

```
10000 rows x 14 columns
```

2.1 Scaling & Normalization

```
In [16]: 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[16]:
```

	CustomerID	Age	Gender	Subscription_Length_Months	Monthly_Bill	Total_Usage_GB	Churn	LocationNew	bill_X_GB	bill_X_sublen	subs_X_bill	GB_X_bill	GB_X_sublen
0	0.000000	0.983835	1.0	0.056602	0.028439	0.413333	0.0	0.50	0.320056	0.031043	0.280769	0.158873	0.520025
1	0.000000	0.983842	0.0	0.000000	0.280800	0.271133	0.0	1.00	0.143887	0.481200	0.033304	0.187387	0.143887
2	0.000000	0.113885	0.0	0.177813	0.762426	0.011111	0.0	0.00	0.060200	0.164761	0.042411	0.302607	0.760200
3	0.000000	0.346354	0.0	0.066667	0.979761	0.548889	1.0	0.75	0.586124	0.338002	0.026122	0.157123	0.586124
4	0.000004	0.538462	0.0	0.782609	0.402300	0.480000	0.0	0.76	0.287560	0.018129	0.462143	0.252844	0.287560
...
9995	0.99996	0.288462	1.0	0.856602	0.390600	0.291111	1.0	0.00	0.229602	0.011613	0.513460	0.223204	0.229602
9996	0.99997	0.983842	0.0	0.782609	0.425426	0.688889	0.0	1.00	0.433380	0.020000	0.377601	0.322100	0.433380
9997	0.99998	0.983815	1.0	0.856602	0.944429	0.446667	1.0	0.00	0.466660	0.044636	0.213131	0.131006	0.466660
9998	0.99999	0.534615	0.0	0.830007	0.275000	0.853333	1.0	1.00	0.409601	0.012277	0.501549	0.515173	0.409601
9999	1.00000	0.173077	0.0	0.782609	0.686286	0.273333	1.0	0.50	0.242170	0.020154	0.301542	0.109513	0.242170

```
10000 rows x 14 columns
```

2.2 Finding Correlation

```
In [18]: corr_matrix = normalized_df.corr()
corr_matrix['Churn'].sort_values(ascending=False)
```

```
Out[18]:
```

Churn	1.000000
LocationNew	0.808405
subs_X_bill	0.493589
Subscription_Length_Months	0.480225
Age	0.480759
bill_X_sublen	0.480759
Monthly_Bill	0.480721
bill_X_bill	0.480225
Total_Usage_GB	0.480242
GB_X_bill	0.480358
subs_X_bill	0.480358
GB_X_sublen	0.480456
CustomerID	0.480456
Name: Churn, dtype: float64	

```
In [17]: import matplotlib.pyplot as plt
plt.rcParams['font.family'] = 'serif'
```

```
In [18]: from pandas.plotting import scatter_matrix
fig, axes = plt.subplots(10, 10, figsize=(10, 10))
scatter_matrix(data[attributes], figure=fig)
```

```
Out[18]:
```

array([[<Axes: xlabel='Subscription_Length_Months', ylabel='Subscription_Length_Months',>, <Axes: xlabel='Total_Usage_GB', ylabel='Subscription_Length_Months',>, <Axes: xlabel='Churn', ylabel='Subscription_Length_Months',>, <Axes: xlabel='LocationNew', ylabel='Subscription_Length_Months',>, <Axes: xlabel='Subscription_Length_Months', ylabel='Monthly_Bill',>, <Axes: xlabel='Monthly_Bill', ylabel='Subscription_Length_Months',>, <Axes: xlabel='Churn', ylabel='Monthly_Bill',>, <Axes: xlabel='Subscription_Length_Months', ylabel='Total_Usage_GB',>, <Axes: xlabel='Total_Usage_GB', ylabel='Monthly_Bill',>, <Axes: xlabel='Churn', ylabel='Total_Usage_GB',>, <Axes: xlabel='LocationNew', ylabel='Total_Usage_GB',>, <Axes: xlabel='Subscription_Length_Months', ylabel='LocationNew',>, <Axes: xlabel='Monthly_Bill', ylabel='LocationNew',>, <Axes: xlabel='Total_Usage_GB', ylabel='LocationNew',>, <Axes: xlabel='Churn', ylabel='LocationNew',>], <Axes: xlabel='Subscription_Length_Months', ylabel='Subscription_Length_Months',>, <Axes: xlabel='Total_Usage_GB', ylabel='Subscription_Length_Months',>, <Axes: xlabel='Churn', ylabel='Subscription_Length_Months',>, <Axes: xlabel='LocationNew', ylabel='Subscription_Length_Months',>, <Axes: xlabel='Subscription_Length_Months', ylabel='Monthly_Bill',>, <Axes: xlabel='Monthly_Bill', ylabel='Subscription_Length_Months',>, <Axes: xlabel='Churn', ylabel='Monthly_Bill',>, <Axes: xlabel='Subscription_Length_Months', ylabel='Total_Usage_GB',>, <Axes: xlabel='Total_Usage_GB', ylabel='Monthly_Bill',>, <Axes: xlabel='Churn', ylabel='Total_Usage_GB',>, <Axes: xlabel='LocationNew', ylabel='Total_Usage_GB',>, <Axes: xlabel='Subscription_Length_Months', ylabel='LocationNew',>, <Axes: xlabel='Monthly_Bill', ylabel='LocationNew',>, <Axes: xlabel='Total_Usage_GB', ylabel='LocationNew',>, <Axes: xlabel='Churn', ylabel='LocationNew',>], <

