

## Python Code

main.py

```
import pandas as pd
from missing_value_handling import handle_missing
from outlier_handling import handle_outliers
from data_normalisation_handling import handle_normalisation
from feature_handling import handle_feature

#Reading csv
df = pd.read_csv('customer_insights_raw.csv')

#calling functions
handle_missing(df.copy())
handle_outliers(df.copy())
handle_normalisation(df.copy())
handle_feature(df.copy())
```

## Section 1: Missing Value Handling

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

pd.set_option(pat='display.max_columns', val=None)

def handle_missing(df): 2 usages

#generating summary of data
print('Summary of csv data \n', df.info())

#Finding missing columns
missing_cols = df.columns[df.isnull().any()]

#dividing missing columns into numeric ones and category ones
numeric_col = df[missing_cols].select_dtypes(include=['int64', 'float64']).columns.tolist()
category_col = df[missing_cols].select_dtypes(include=['object']).columns.tolist()

#Median and mean dataset for cols
df_mean = df.copy()
df_median = df.copy()

#filling mean values to numeric cols
for col in numeric_col:
    df_mean[col].fillna(df_mean[col].mean(), inplace=True)

#filling median values to numeric cols
for col in numeric_col:
    df_median[col].fillna(df[col].median(), inplace=True)
```

```
#filling mode values to category cols
for col in category_col:
    df_mean[col].fillna(df[col].mode()[0], inplace=True)
    df_median[col].fillna(df[col].mode()[0], inplace=True)

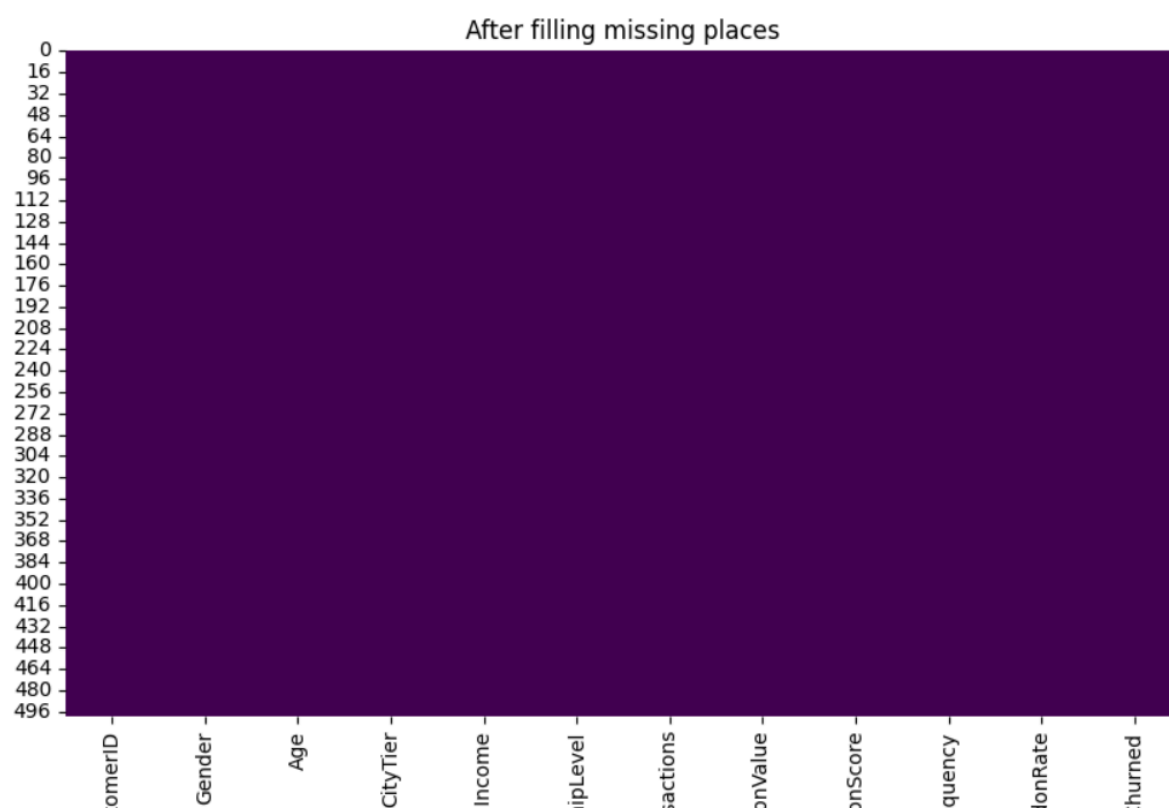
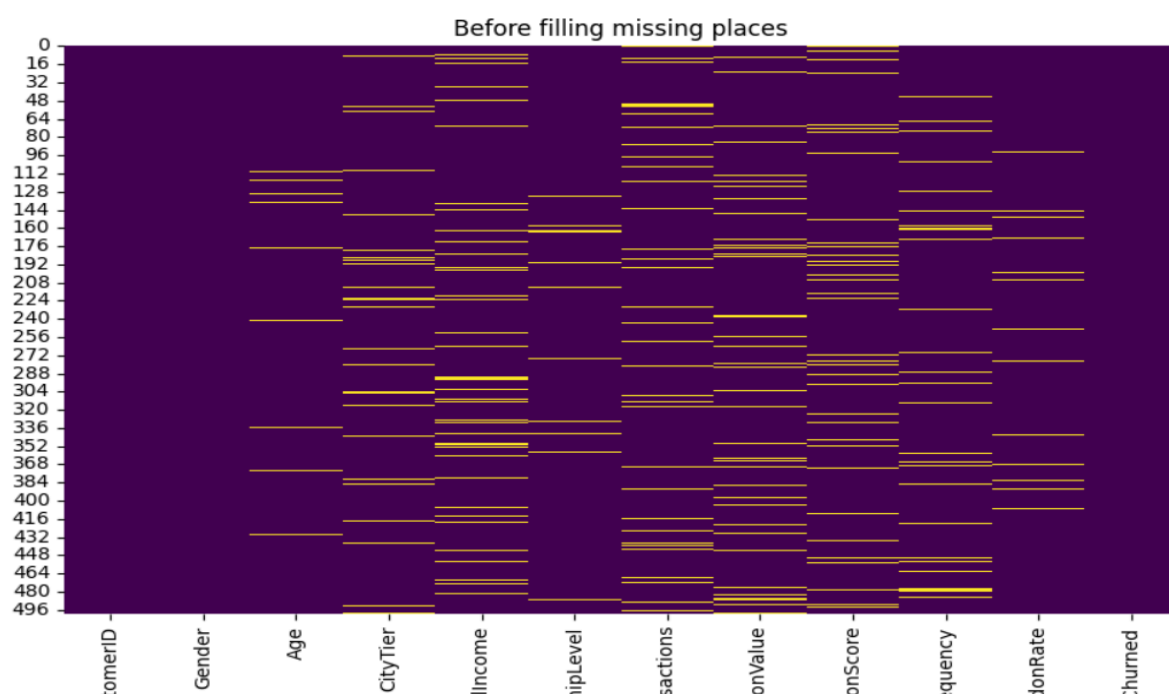
#Comparing mean and median imputations in missing fields
print('\nMean imputed dataset \n', df_mean.describe())
print('\nMedian imputed dataset \n', df_median.describe())

#dropping rows where more than 2 values are missing in the row
df_drop_rows = df_median[df.isnull().sum(axis = 1) <= 2]

#making new csv with cleaned data
df_drop_rows.to_csv('customer_cleaned_imputed.csv', index=False)

#plotting heatmap
plot(title='Before filling missing places', df)
plot(title='After filling missing places', df_drop_rows)

#using functions so that we can reuse it
def plot(title, data): 2 usages
    plt.figure(figsize = (10,6))
    sns.heatmap(
        data.isnull(),
        cbar = False,
        cmap = 'viridis'
    )
    plt.title(title)
    plt.show()
```



## Section 2: Outlier Detection and Handling

```
import pandas as pd
from scipy.stats import zscore
import seaborn as sns
import matplotlib.pyplot as plt

def handle_outliers(df): 2 usages
|
|
|     #taking the numerical cols
|     numeric_cols = ['AnnualIncome', 'AvgTransactionValue', 'LoginFrequency']
|
|     #Applying zscore
|     z_score = df[numeric_cols].apply(zscore)
|
|     threshold = 3
|     outliers_zscore = (z_score.abs() > threshold)
|     # print(outliers_zscore.sum())
|
|     data_outliers_zscore = df[outliers_zscore.any(axis=1)]
|
|     print('Outliers in data frame : \n', data_outliers_zscore)
|
|     #Applying IQR
|
|     iqr_bounds = {}
|     for col in numeric_cols:
|         q1 = df[col].quantile(0.25)
|         q3 = df[col].quantile(0.75)
|
|         iqr = q3 - q1
|
|         lower_bound = q1 - 1.5 * iqr
|         upper_bound = q3 + 1.5 * iqr
```

```
iqr_bounds[col] = (lower_bound, upper_bound)

print(f'{col}: {lower_bound:.2f} - {upper_bound:.2f}')

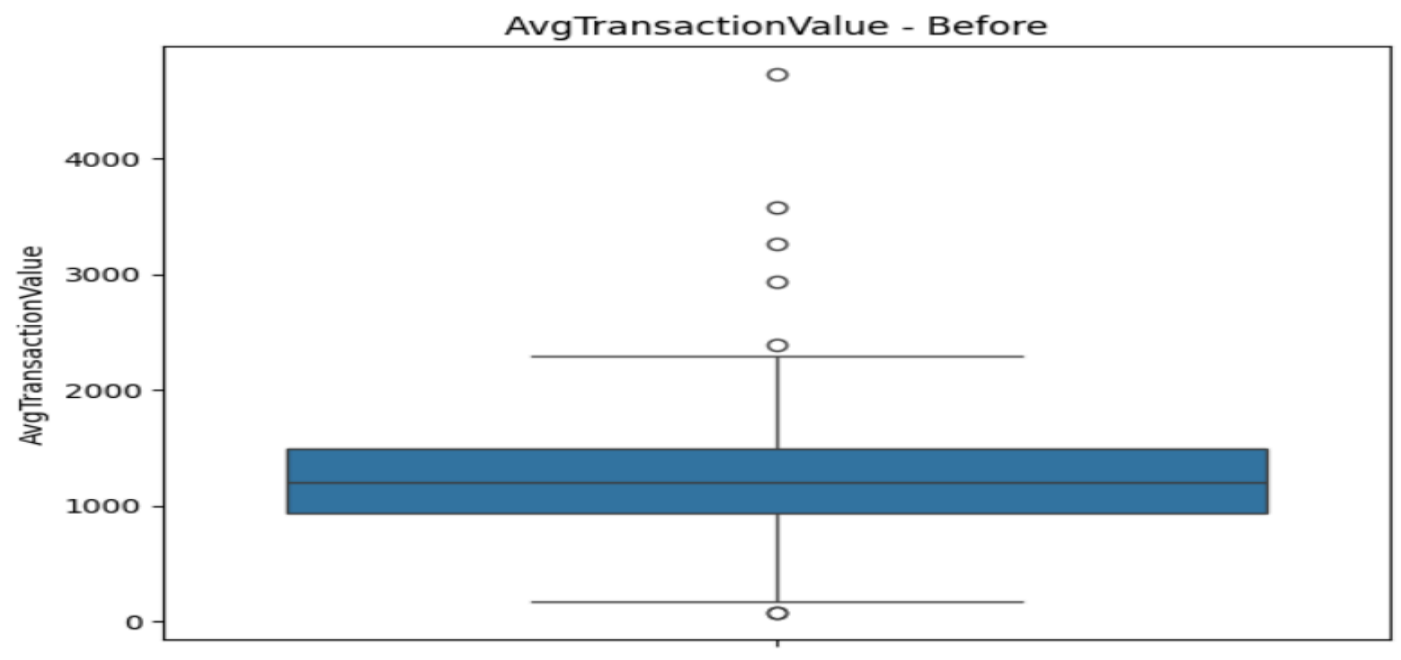
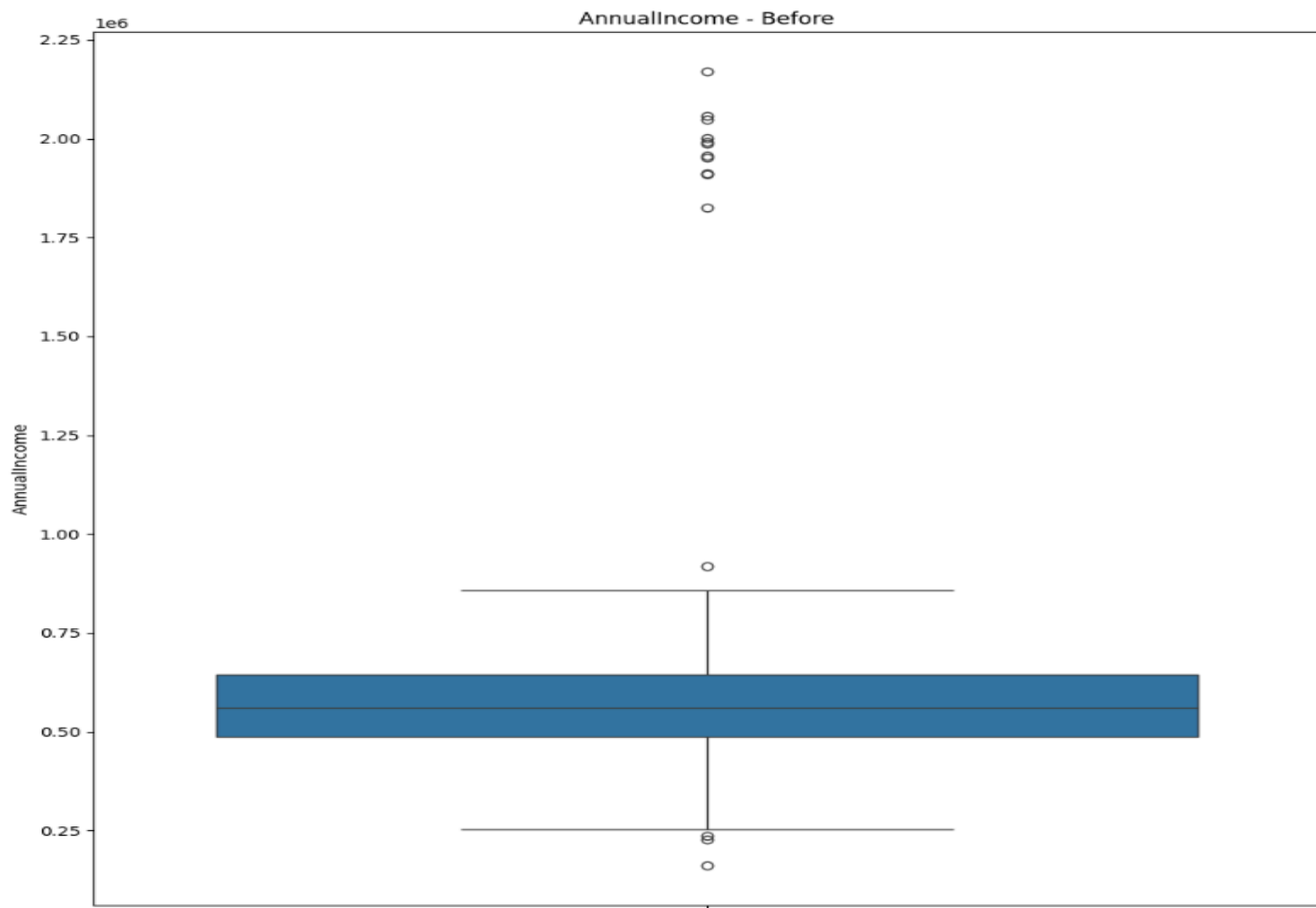
capping_data = df.copy()

for col in numeric_cols:
    lower_bound, upper_bound = iqr_bounds[col]
    capping_data[col] = capping_data[col].apply(
        lambda x : lower_bound if x < lower_bound else upper_bound if x > upper_bound else x
    )

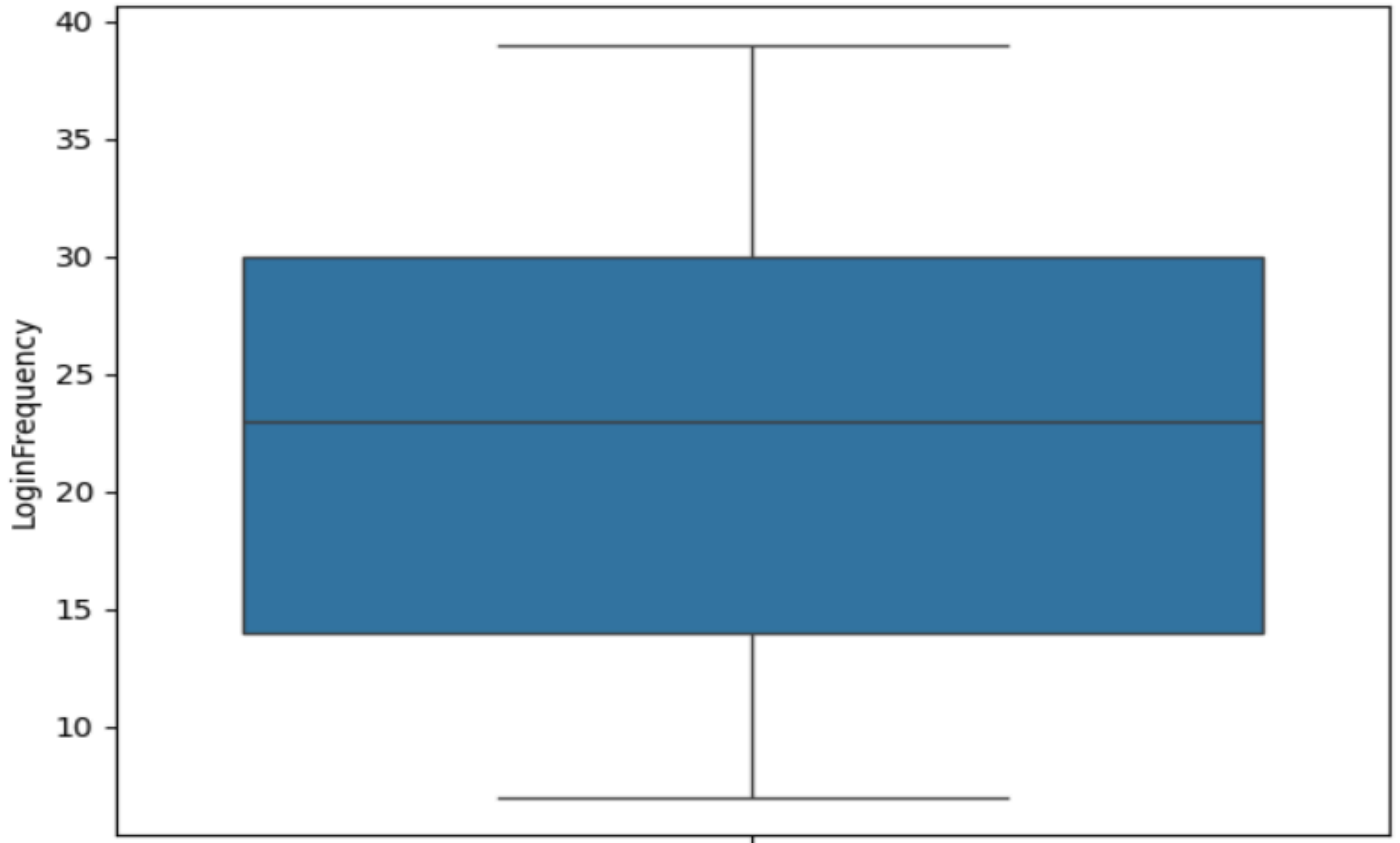
#Before Outlier
plt.figure(figsize = (10,10))
for col in numeric_cols:
    plot(df, title: f'{col} - Before', col)

#After Outlier
for col in numeric_cols:
    plot(capping_data, title: f'{col} - After', col)

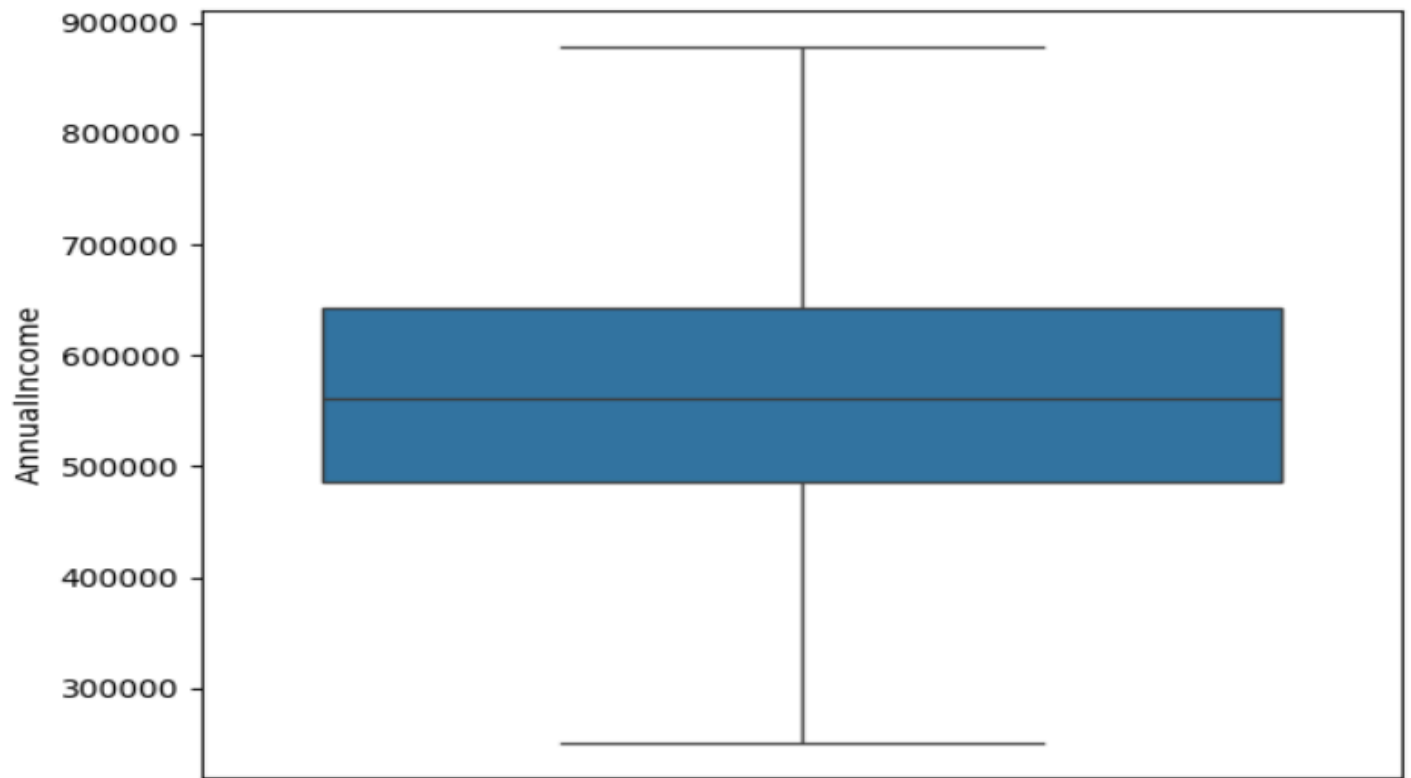
def plot(data, title, col): 2 usages
    sns.boxplot(
        data,
        y = col
    )
    plt.title(title)
    plt.tight_layout()
    plt.show()
```

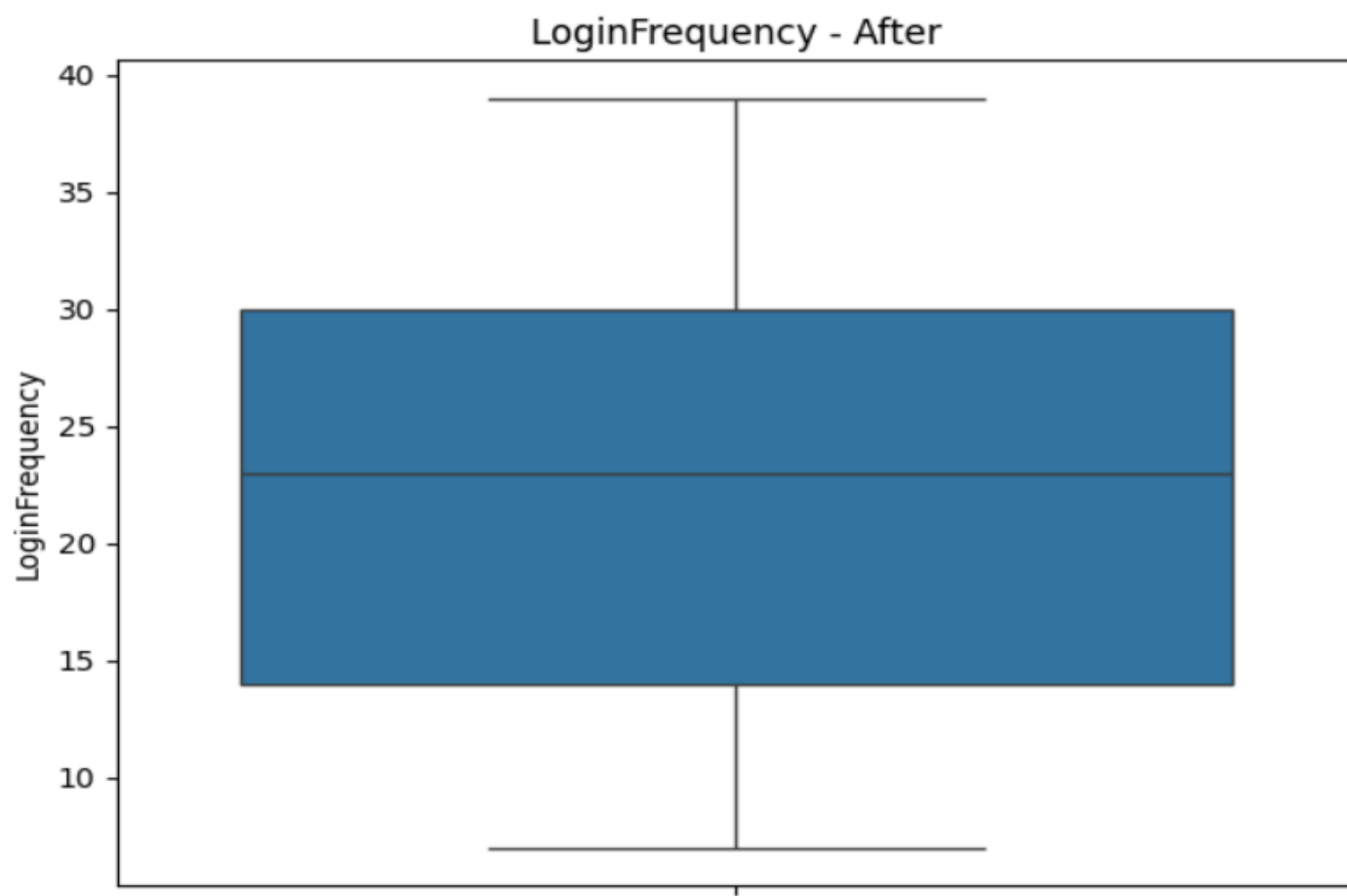
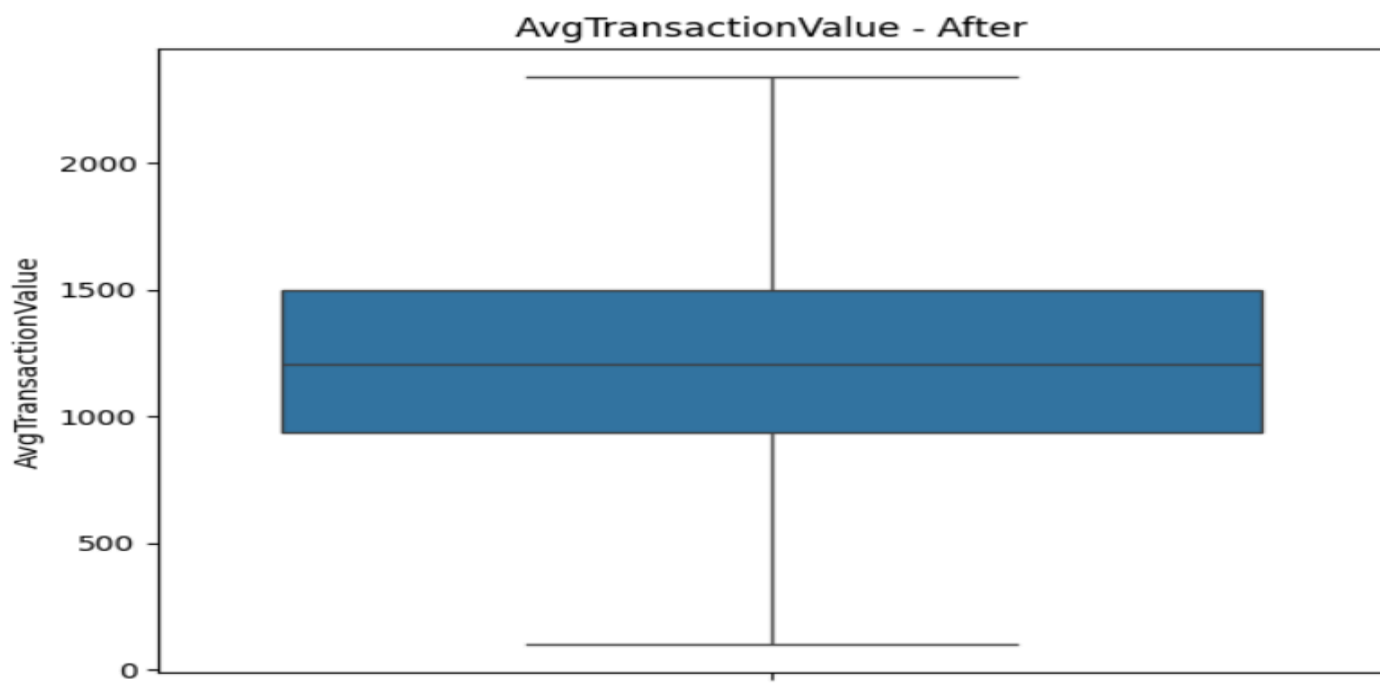


LoginFrequency - Before



AnnualIncome - After





## Section 3: Data Normalization

```
import pandas as pd
from matplotlib import pyplot as plt
from sklearn.preprocessing import MinMaxScaler, StandardScaler
```

```
pd.set_option(pat: 'display.max_columns', val: None)
```

```
def handle_normalisation(df): 2 usages
```

```
💡 #Selecting the columns for normalisation|
cols = ['AnnualIncome', 'AvgTransactionValue', 'LoginFrequency']
```

```
#Min-Max scaling
minmax_scaler = MinMaxScaler()
df_minmax = pd.DataFrame(
    minmax_scaler.fit_transform(df[cols]),
    columns = [f'{col}_MinMax' for col in cols]
)
```

```
#Z-Score standardisation
zscore_standardisation = StandardScaler()
df_zscore = pd.DataFrame(
    zscore_standardisation.fit_transform(df[cols]),
    columns=[f'{col}_zscore' for col in cols]
)
```

```
#making them into data set to push further to csv files
df_minmax = pd.concat(objs: [df, df_minmax], axis=1)
df_zscore = pd.concat(objs: [df, df_zscore], axis=1)
```

```
#plotting histogram
plt.figure(figsize = (10, 10))
```

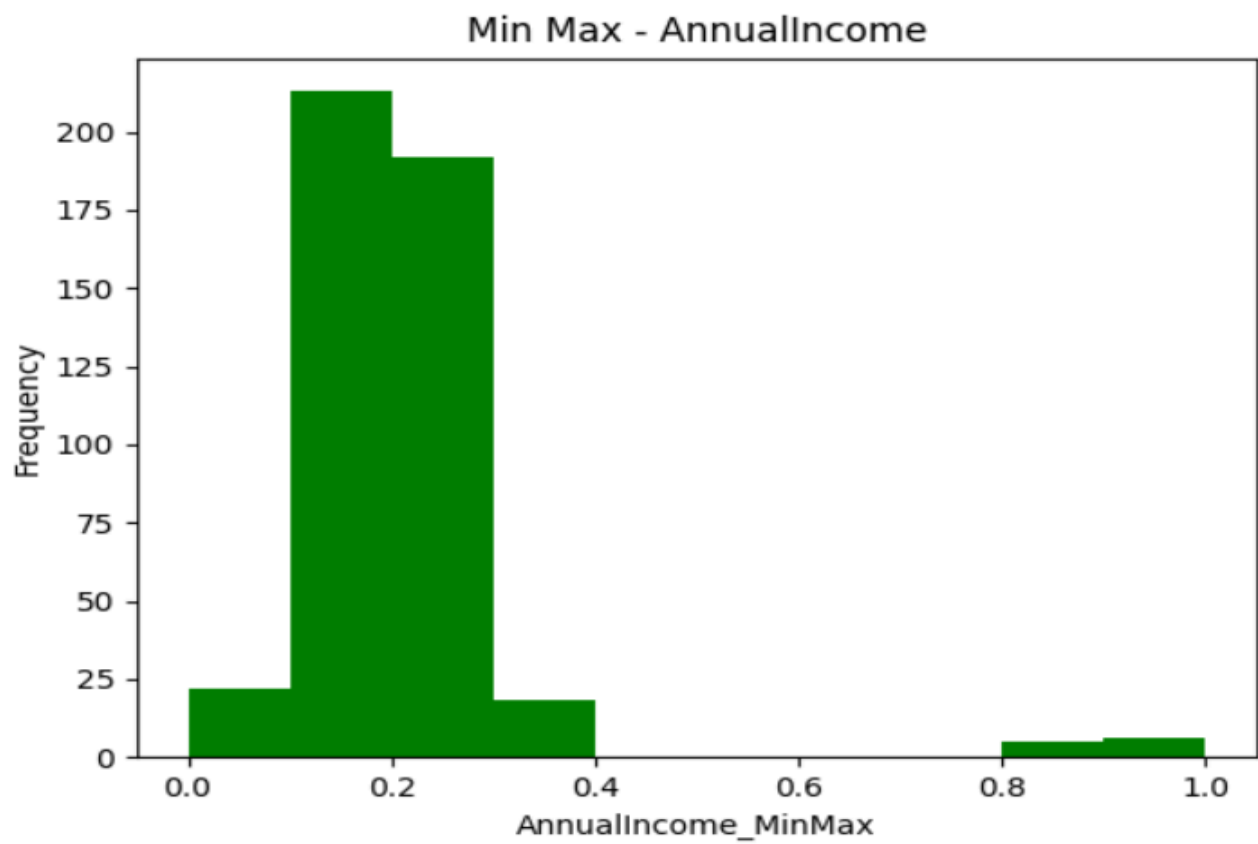
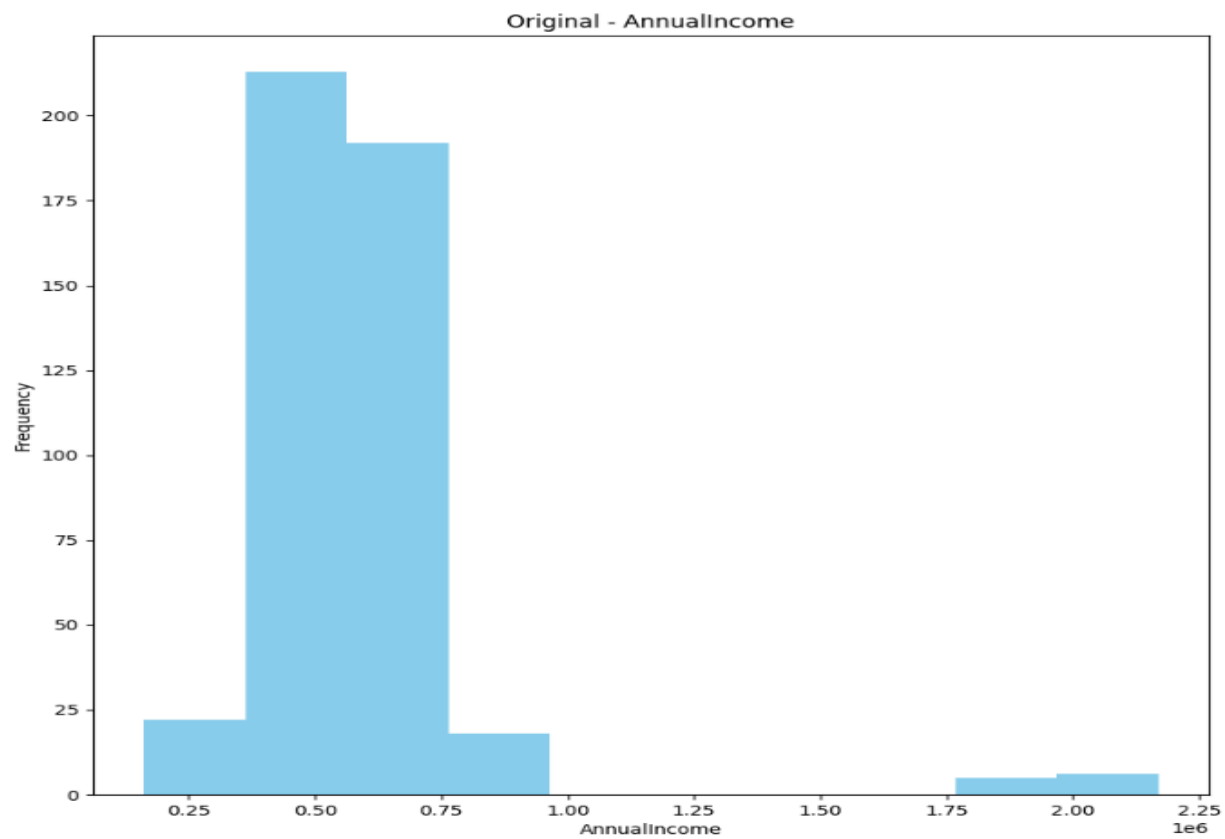
```
for col in cols:
    plot(df, col, title: f'Original - {col}', clr: 'skyblue')
    plot(df_minmax, col: f'{col}_MinMax', title: f'Min Max - {col}', clr: 'green')
    plot(df_zscore, col: f'{col}_zscore', title: f'Z-Score - {col}', clr: 'orange')
```

```
#Making a new csv with the data frames
df_minmax.to_csv(path_or_buf: 'customer_minmax_scaled.csv', index = False)
df_zscore.to_csv(path_or_buf: 'customer_standard_scaled.csv', index = False)
```

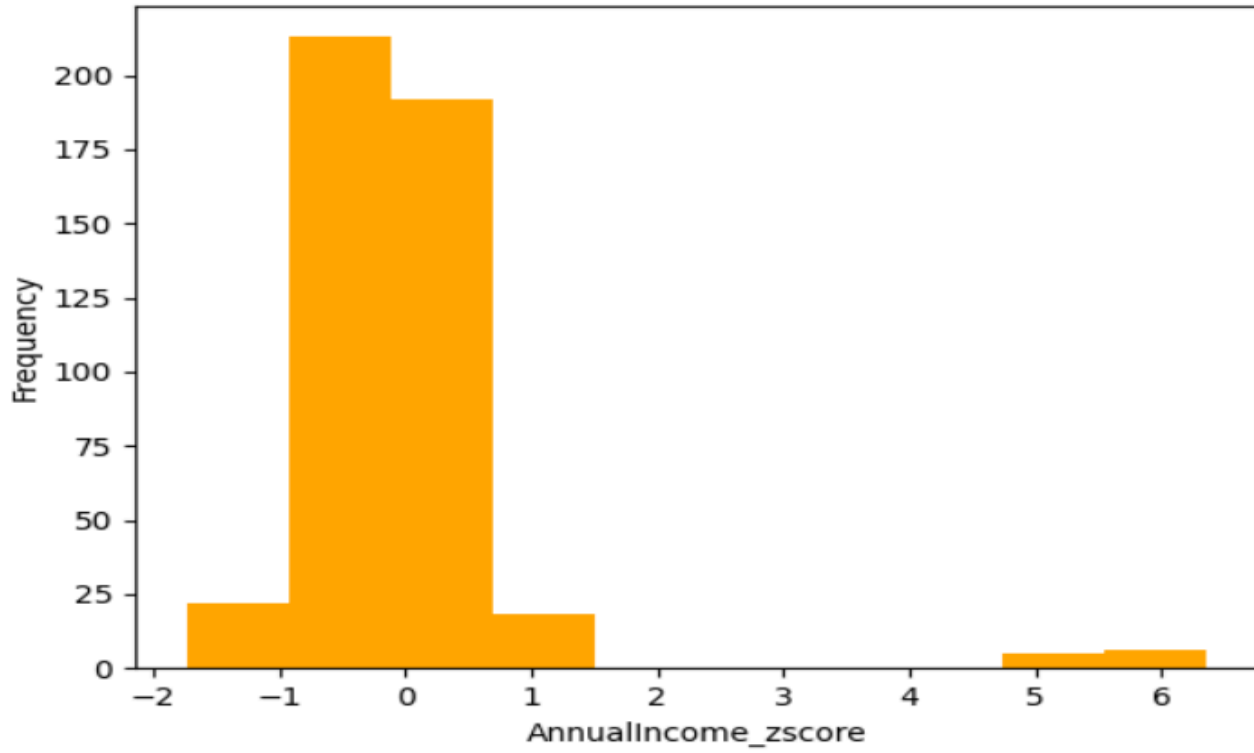
```
def plot(df, col, title, clr): 3 usages
```

```
plt.hist(
    df[col],
    bins = 10,
    color = clr
)
plt.title(title)
plt.xlabel(col)
plt.ylabel('Frequency')
plt.show()
```

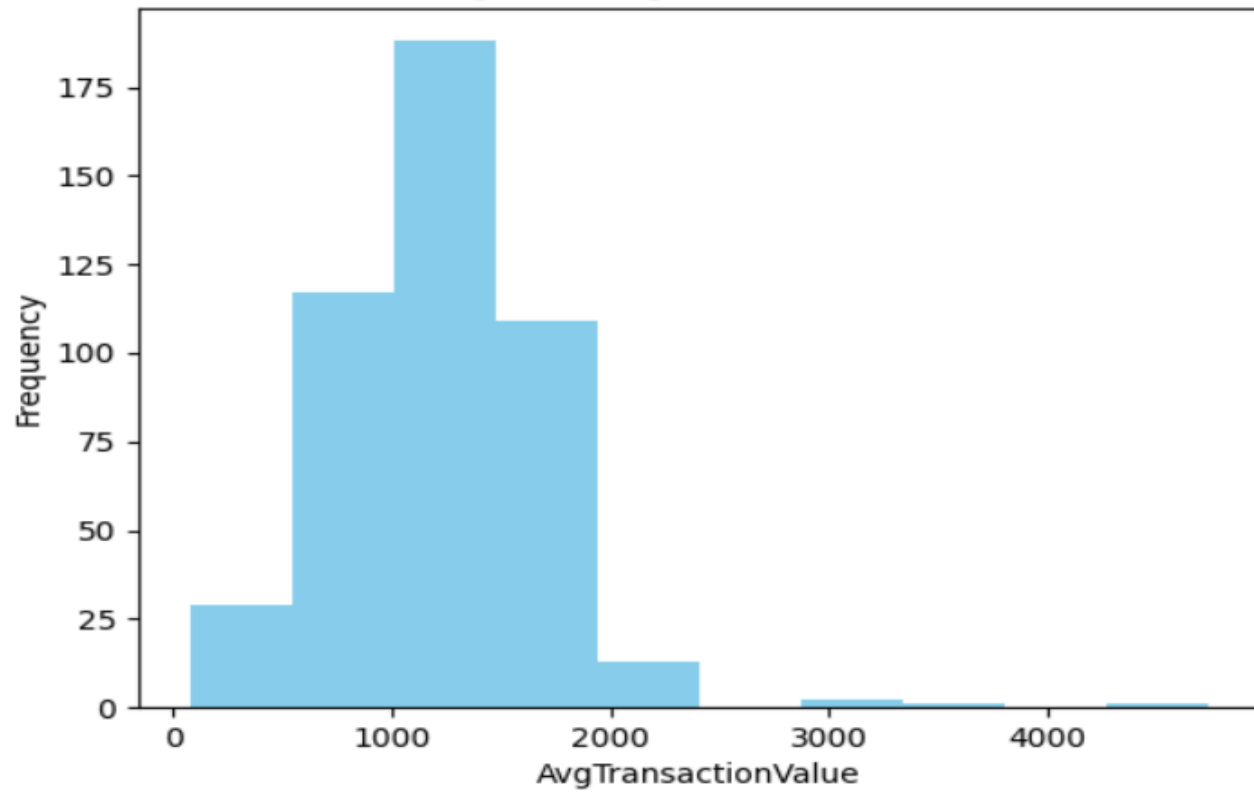




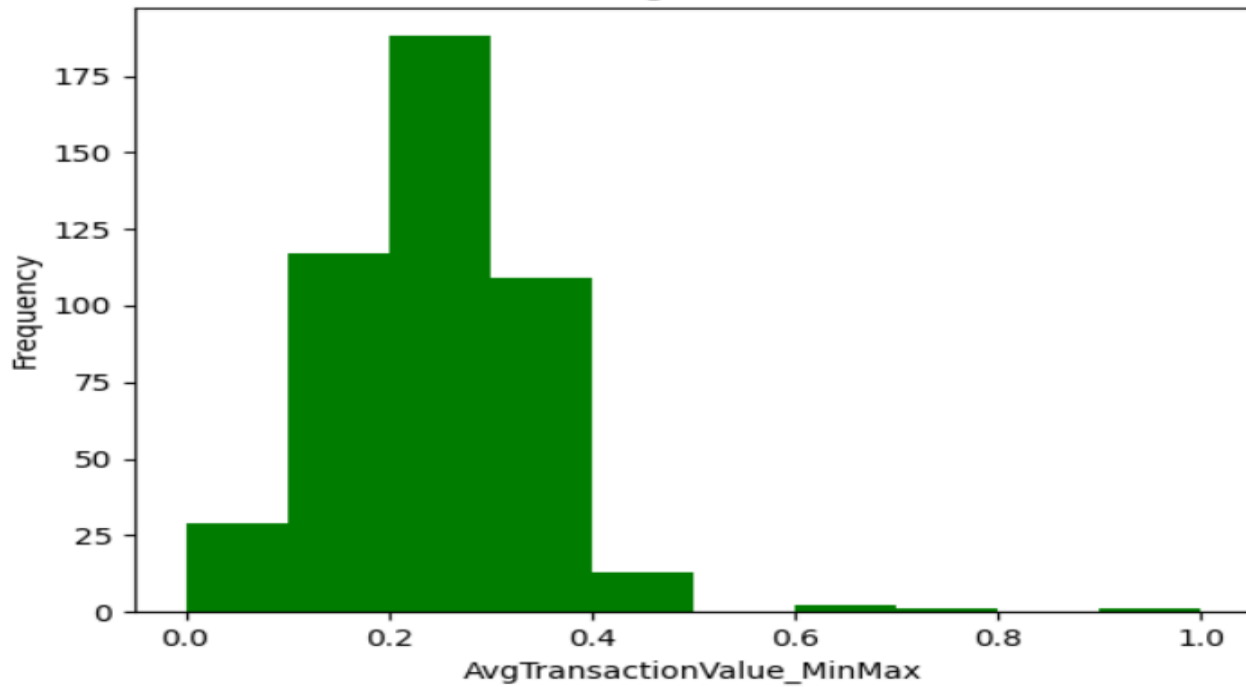
Z-Score - AnnualIncome



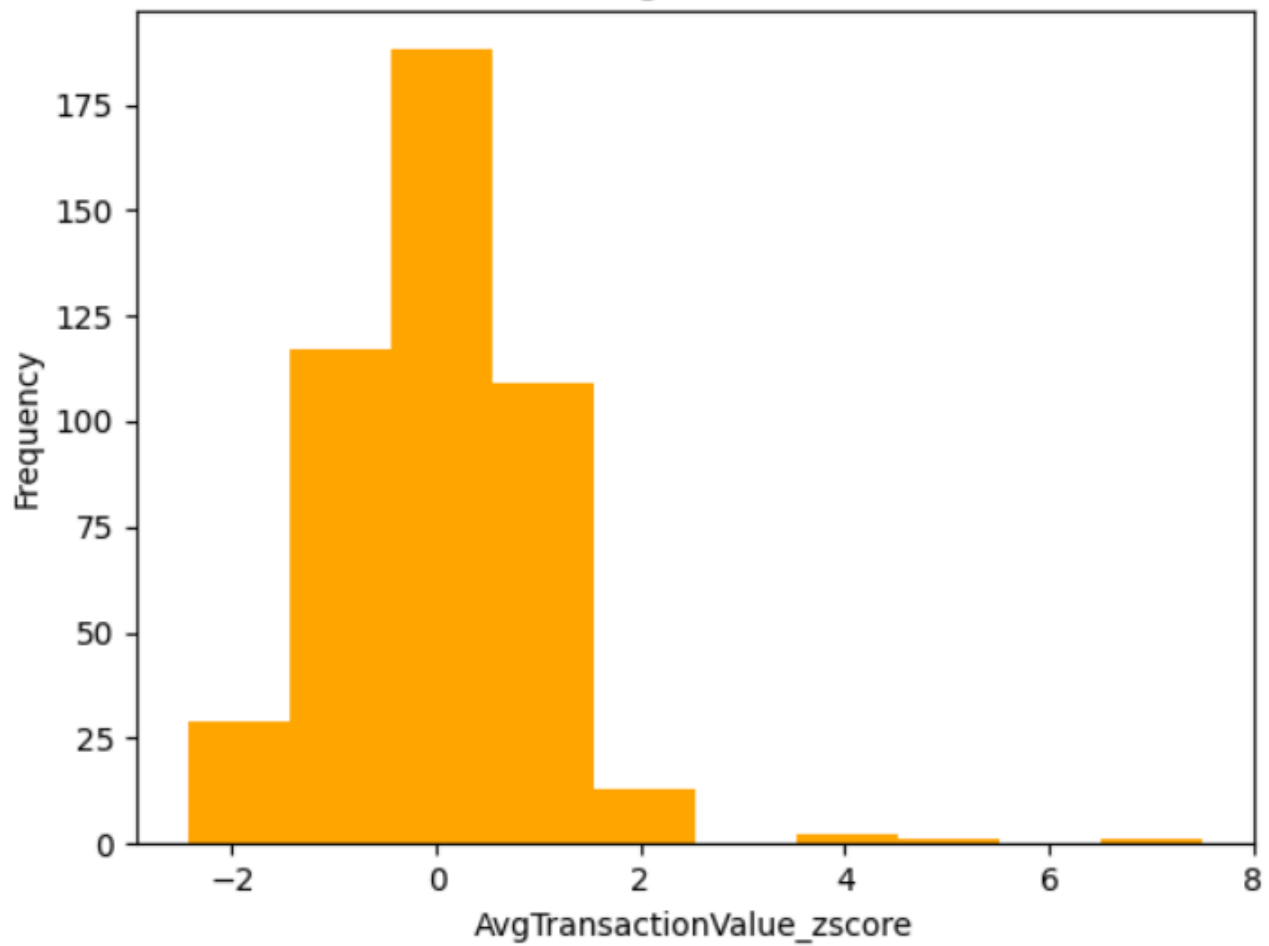
Original - AvgTransactionValue



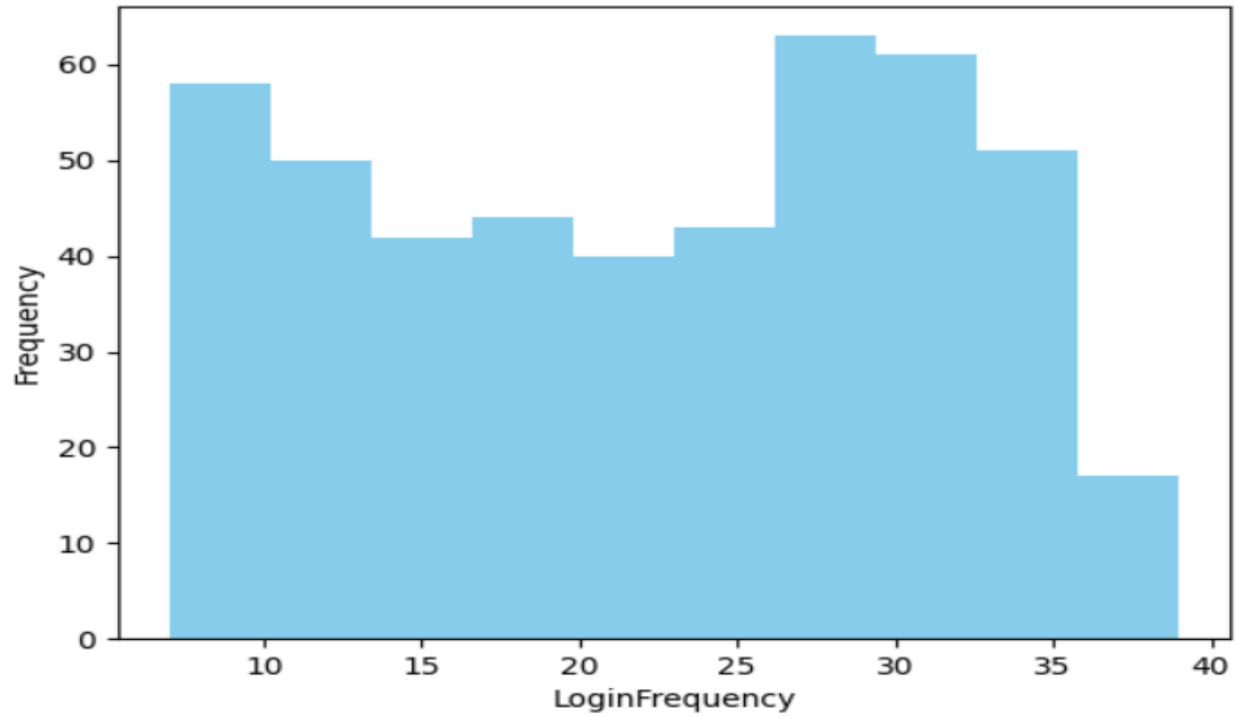
Min Max - AvgTransactionValue



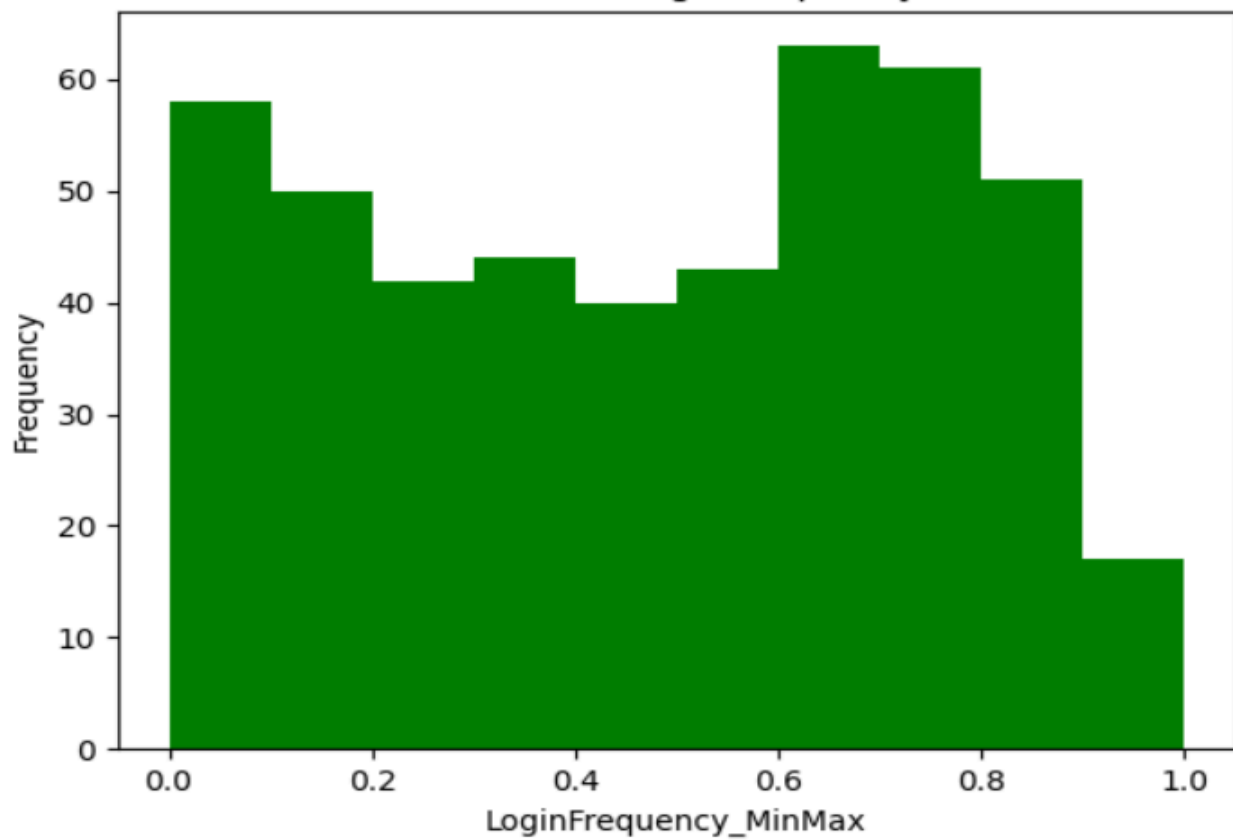
Z-Score - AvgTransactionValue



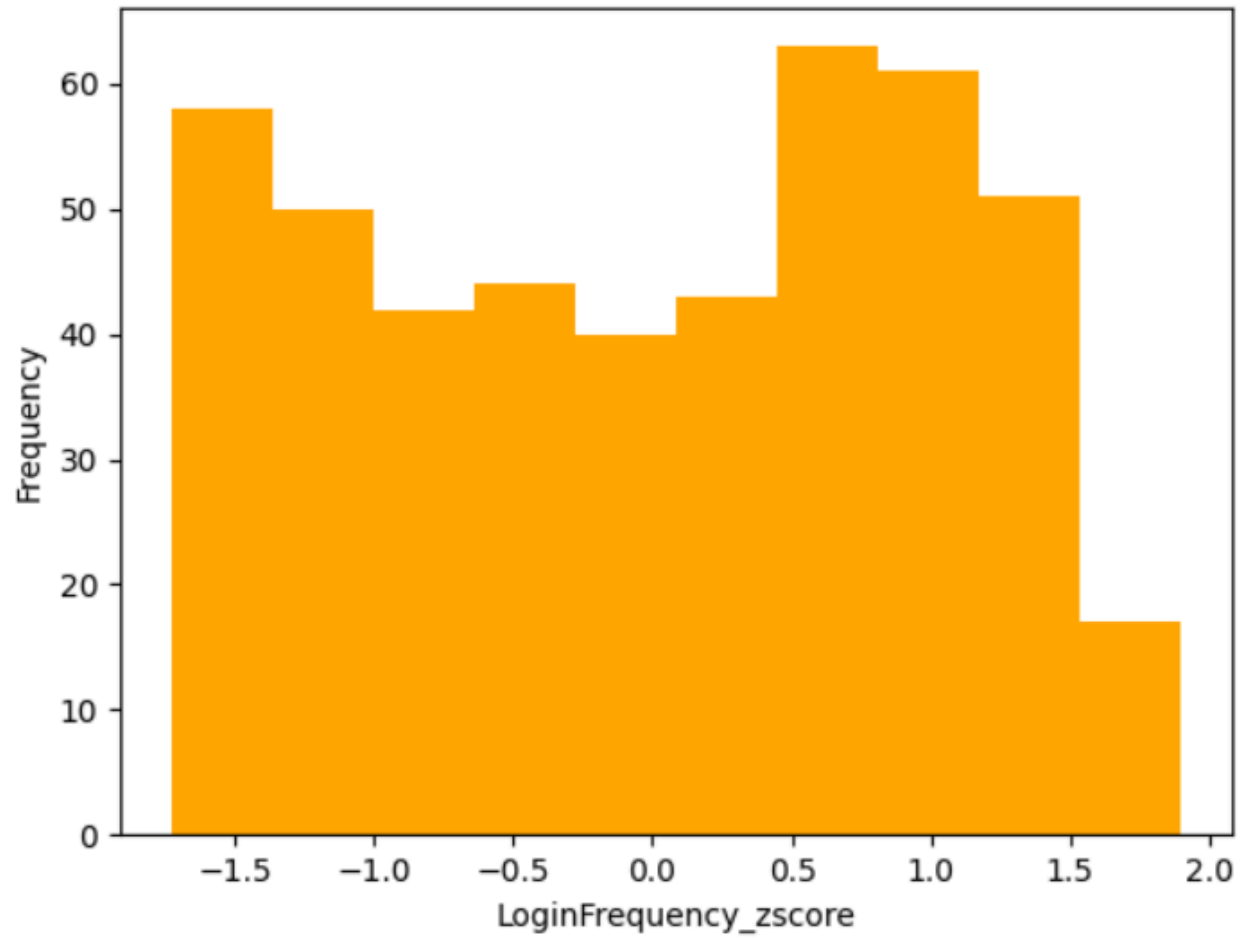
Original - LoginFrequency



Min Max - LoginFrequency



Z-Score - LoginFrequency



## Section 4: Feature Engineering

```
import pandas as pd
pd.set_option(pat: 'display.max_columns', val: None)

def handle_feature(df): 2 usages

    #creating new feature
    df['TotalSpend'] = df['TotalTransactions'] * df['AvgTransactionValue']

    #creating age groups
    df['AgeGroup'] = df['Age'].apply(
        lambda x : 'Child' if x < 18 else
        'Young' if x < 25 else 'Young Adults'
        if x < 36 else 'Mid-Age Adults' if x < 46 else 'Old'
    )

    #converting isChurned
    df['IsChurned'] = df['IsChurned'].apply(
        lambda x : 1 if x.lower() == 'yes' else 0
    )

    #average transaction frequency per month
    df['MonthlyTransactionRate'] = df['TotalTransactions'] / (df['LoginFrequency'] / 4)

    #Derive a new score
    df['CustomerScore'] = (df['TotalSpend'] * df['SatisfactionScore']) / df['CartAbandonRate']

    df.to_csv('customer_features_enriched.csv', index=False)
```