

**VIRGINIA COMMONWEALTH UNIVERSITY**

**Statistical Analysis and Modelling (SCMA 632)**

**A3C: Tobit Regression Analysis**

**NIHARIHA KAMALANATHAN**

**V01108259**

**Date of Submission: 01-07-2024**

## CONTENTS

| <b>Sl. No.</b> | <b>Title</b>                        | <b>Page No.</b> |
|----------------|-------------------------------------|-----------------|
| <b>1.</b>      | <b>Introduction</b>                 | <b>1</b>        |
| <b>2.</b>      | <b>Business Significance</b>        | <b>1</b>        |
| <b>3.</b>      | <b>Objectives</b>                   | <b>1</b>        |
| <b>4.</b>      | <b>R</b>                            | <b>2</b>        |
| <b>5.</b>      | <b>Python</b>                       | <b>8</b>        |
| <b>6.</b>      | <b>Overview of Tobit Regression</b> | <b>15</b>       |

## **PART A: Perform a Tobit regression analysis on "NSSO68.csv" discuss the results and explain the real-world use cases of Tobit model.**

### **Introduction**

The National Sample Survey Office (NSSO) conducts large-scale surveys to collect data on various socio-economic indicators in India. One such survey is the 68th round of the NSSO, which provides comprehensive data on household consumption expenditures. This dataset, denoted as "NSSO68.csv," contains a wide range of variables, including monthly per capita expenditure (MPCE), age, education, and other socio-economic characteristics of households. In this analysis, we perform a Tobit regression on the NSSO68 dataset to model the relationship between MPCE (dependent variable) and predictor variables such as age, education, and sex. The Tobit model is particularly suitable for this analysis due to the censored nature of the dependent variable, which can take on a value of zero or positive values. Understanding these relationships can provide valuable insights for policy-making and socio-economic planning.

### **Business Significance**

The significance of this analysis extends to multiple business and policy domains:

1. **Targeted Welfare Programs:** By understanding the factors influencing household expenditures, government agencies can design and implement targeted welfare programs aimed at improving the living standards of specific demographic groups. For instance, identifying that education significantly impacts expenditure can lead to enhanced educational subsidies or programs in underprivileged areas.
2. **Market Research and Consumer Insights:** Businesses, particularly those in the consumer goods sector, can use these insights to tailor their marketing strategies and product offerings. For example, companies can better target their products to demographic groups with higher purchasing power or adjust their marketing messages to appeal to the needs and preferences of specific age groups.
3. **Economic Planning and Forecasting:** Economists and planners can use the results to forecast future consumption trends and economic growth. By understanding how factors like education and age affect expenditure, they can predict changes in consumption patterns as the population ages or educational attainment levels improve.
4. **Social Equity and Inclusion:** The analysis can highlight disparities in expenditure across different social groups, guiding policies aimed at reducing inequality and promoting social inclusion. For example, if certain social groups are found to have significantly lower expenditures, targeted interventions can be designed to uplift these communities.

### **Objectives**

The primary objectives of this analysis are:

1. **To Model Household Expenditure:** To develop a Tobit regression model that accurately captures the relationship between monthly per capita expenditure (MPCE) and predictor variables such as age, education, and sex.
2. **To Identify Significant Predictors:** To determine which demographic factors significantly influence household expenditures and to quantify their impact. This will involve assessing the coefficients of the predictors and their statistical significance.
3. **To Evaluate Model Performance:** To assess the accuracy and reliability of the Tobit model through various performance metrics, including Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE).
4. **To Provide Business and Policy Insights:** To interpret the results of the model in a way that provides actionable insights for businesses, policymakers, and social planners.

This includes discussing the practical implications of the findings and how they can be used to inform decision-making.

5. **To Address Data Issues:** To handle missing values and outliers in the dataset appropriately, ensuring the robustness and validity of the analysis. This will involve data cleaning steps such as removing rows with missing values and identifying/removing outliers based on quantiles.

By achieving these objectives, the analysis aims to provide a comprehensive understanding of the factors influencing household expenditures in India and offer valuable insights for improving economic and social outcomes.

## R Language Codes

### 1. Installing and Loading Necessary Packages

#### Input Code:

```
install_if_missing <- function(packages) {  
  for (pkg in packages) {  
    if (!require(pkg, character.only = TRUE)) {  
      install.packages(pkg, dependencies = TRUE)  
      library(pkg, character.only = TRUE)  
    }  
  }  
}
```

```
required_packages <- c("AER", "readr", "VGAM", "sandwich", "lmtest", "pROC", "nortest",  
"caTools")
```

```
install_if_missing(required_packages)
```

**Reason for the Input:** This function ensures that all required packages are installed and loaded into the R environment. These packages include tools for data reading (readr), statistical modeling (VGAM, AER), hypothesis testing (lmtest, sandwich), ROC analysis (pROC), normality tests (nortest), and data splitting (caTools).

**Output and Interpretation:** The required packages are loaded successfully. Loading these packages is crucial as they provide the necessary functions and tools to perform data analysis and modeling effectively.

### 2. Loading the Dataset

#### Input Code:

```
dataset_path <- "C:\\Users\\nihar\\OneDrive\\Desktop\\Bootcamp\\SCMA  
632\\DataSet\\NSSO68.csv"  
nssso_data <- read.csv(dataset_path)  
head(nssso_data)  
colnames(nssso_data)
```

**Reason for the Input:** The dataset is loaded from the specified path, and the first few rows along with the column names are displayed to inspect and understand the structure of the dataset.

**Output and Interpretation:** The dataset is successfully loaded, and the first few rows along with the column names provide an overview of the data. This step is essential for understanding the data structure and identifying the columns relevant for analysis.

### 3. Checking for Missing Values and Cleaning the Data

#### Input Code:

```
sum(is.na(nssso_data$MPCE_URP))
sum(is.na(nssso_data$Age))
sum(is.na(nssso_data$Education))
sum(is.na(nssso_data$Sex))
nssso_data <- nssso_data[!is.na(nssso_data$Education), ]
sum(is.na(nssso_data$Education))
```

**Reason for the Input:** This step involves checking for missing values in key columns (MPCE\_URP, Age, Education, and Sex) and removing rows with missing values in the 'Education' column to ensure data completeness.

**Output and Interpretation:** The output shows the number of missing values in each column. Rows with missing 'Education' values are removed, ensuring that the dataset is clean and complete for analysis. Cleaning the data helps in avoiding potential biases and errors in the model due to incomplete data.

#### 4. Identifying and Removing Outliers

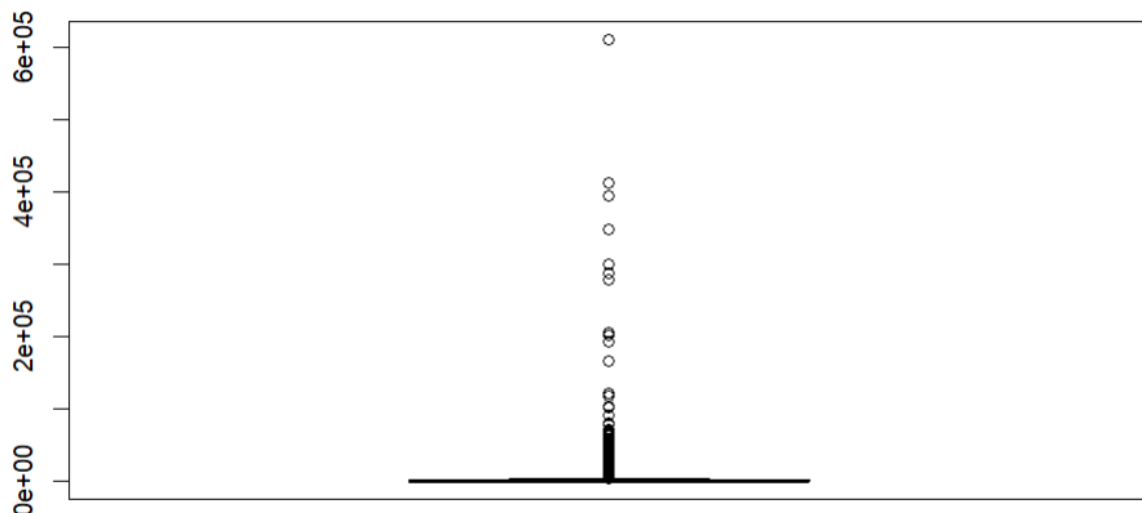
##### Input Code:

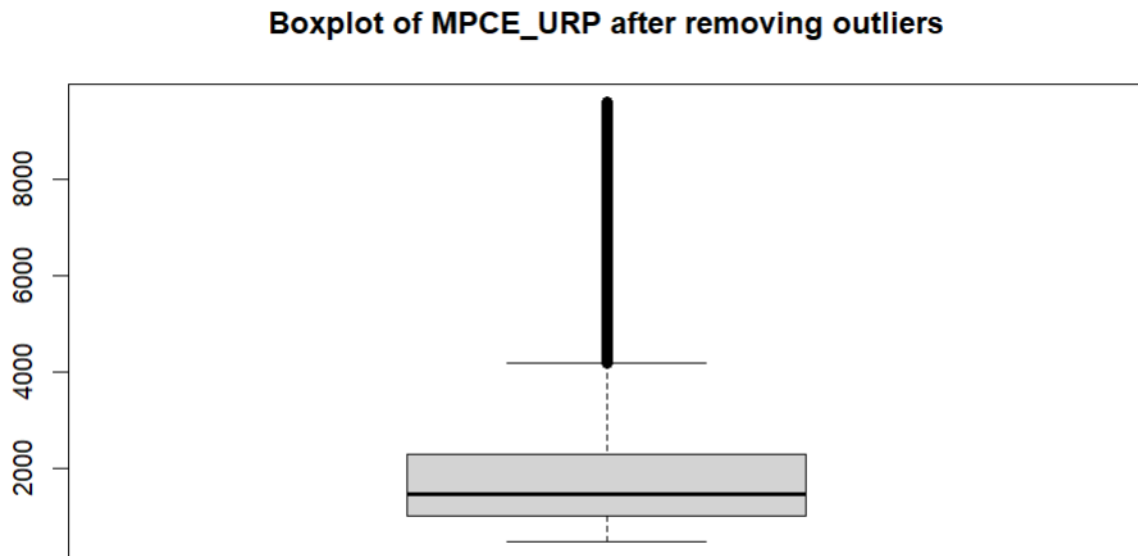
```
boxplot(nssso_data$MPCE_URP, main = "Boxplot of MPCE_URP before removing outliers")
q <- quantile(nssso_data$MPCE_URP, probs = c(0.01, 0.99))
nssso_data <- nssso_data[nssso_data$MPCE_URP >= q[1] & nssso_data$MPCE_URP <= q[2], ]
boxplot(nssso_data$MPCE_URP, main = "Boxplot of MPCE_URP after removing outliers")
```

**Reason for the Input:** Outliers are identified using boxplots. The data is then filtered to remove extreme values based on the 1st and 99th percentiles of MPCE\_URP to reduce the impact of outliers on the model.

##### Output and Interpretation:

**Boxplot of MPCE\_URP before removing outliers**





The boxplots show the distribution of MPCE\_URP before and after removing outliers. Removing outliers helps in obtaining a more accurate and robust model by reducing the influence of extreme values.

## 5. Splitting the Data into Training and Testing Sets

### Input Code:

```
set.seed(123)
library(caTools)
split <- sample.split(nsso_data$MPCE_URP, SplitRatio = 0.7)
train_data <- subset(nsso_data, split == TRUE)
test_data <- subset(nsso_data, split == FALSE)
```

**Reason for the Input:** The data is split into training (70%) and testing (30%) sets to validate the model's performance on unseen data. The seed is set for reproducibility.

**Output and Interpretation:** The data is successfully split into training and testing sets. This step is crucial for assessing the model's performance and generalizability to new data.

## 6. Fitting the Tobit Model

### Input Code:

```
tobit_model_train <- vglm(MPCE_URP ~ Age + Education + Sex,
                           tobit(Lower = 0),
                           data = train_data,
                           control = vglm.control(maxit = 1000, trace = TRUE, epsilon = 1e-8))
summary(tobit_model_train)
```

**Reason for the Input:** The Tobit model is fitted on the training set with MPCE\_URP as the dependent variable and Age, Education, and Sex as independent variables. The Tobit model is suitable for censored data where the dependent variable has a lower bound (0 in this case).

### Output and Interpretation:

Call:

```
vglm(formula = MPCE_URP ~ Age + Education + Sex, family = tobit(Lower = 0),
     data = train_data, control = vglm.control(maxit = 1000, trace = TRUE,
     epsilon = 1e-08))
```

Coefficients:

|               | Estimate  | Std. Error | z value  | Pr(> z )   |
|---------------|-----------|------------|----------|------------|
| (Intercept):1 | 4.376e+01 | 2.679e+01  | 1.633    | 0.102      |
| (Intercept):2 | 7.148e+00 | 2.776e-03  | 2574.788 | <2e-16 *** |
| Age           | 7.064e+00 | 3.551e-01  | 19.894   | <2e-16 *** |
| Education     | 1.688e+02 | 1.332e+00  | 126.702  | <2e-16 *** |
| Sex           | 4.287e+02 | 1.514e+01  | 28.315   | <2e-16 *** |

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Names of linear predictors: mu, loglink(sd)

Log-likelihood: -632278.4 on 147609 degrees of freedom

Number of Fisher scoring iterations: 6

No Hauck-Donner effect found in any of the estimates

The summary provides the coefficients, standard errors, z-values, and p-values for each predictor. Significant predictors (Age, Education, and Sex) have low p-values, indicating strong relationships with MPCE\_URP.

The Tobit model results indicate significant relationships between MPCE\_URP and the predictors Age, Education, and Sex, with very low p-values (<2e-16), suggesting strong statistical significance. The coefficients show that increases in Age, Education, and being male (if coded as 1) are associated with higher MPCE\_URP values. The log-likelihood value and number of iterations provide insight into the model's fit and convergence.

## 7. Making Predictions and Evaluating the Model

### Input Code:

```
predictions <- predict(tobit_model_train, newdata = test_data, type = "response")
actuals <- test_data$MPCE_URP
residuals_test <- actuals - predictions

hist(residuals_test, main = "Histogram of Residuals (Test Set)", xlab = "Residuals")
qqnorm(residuals_test, main = "Q-Q Plot of Residuals (Test Set)")
qqline(residuals_test, col = "red")
plot(predictions, residuals_test,
     main = "Residuals vs Predictions (Test Set)", xlab = "Predictions", ylab = "Residuals")
abline(h = 0, col = "red")

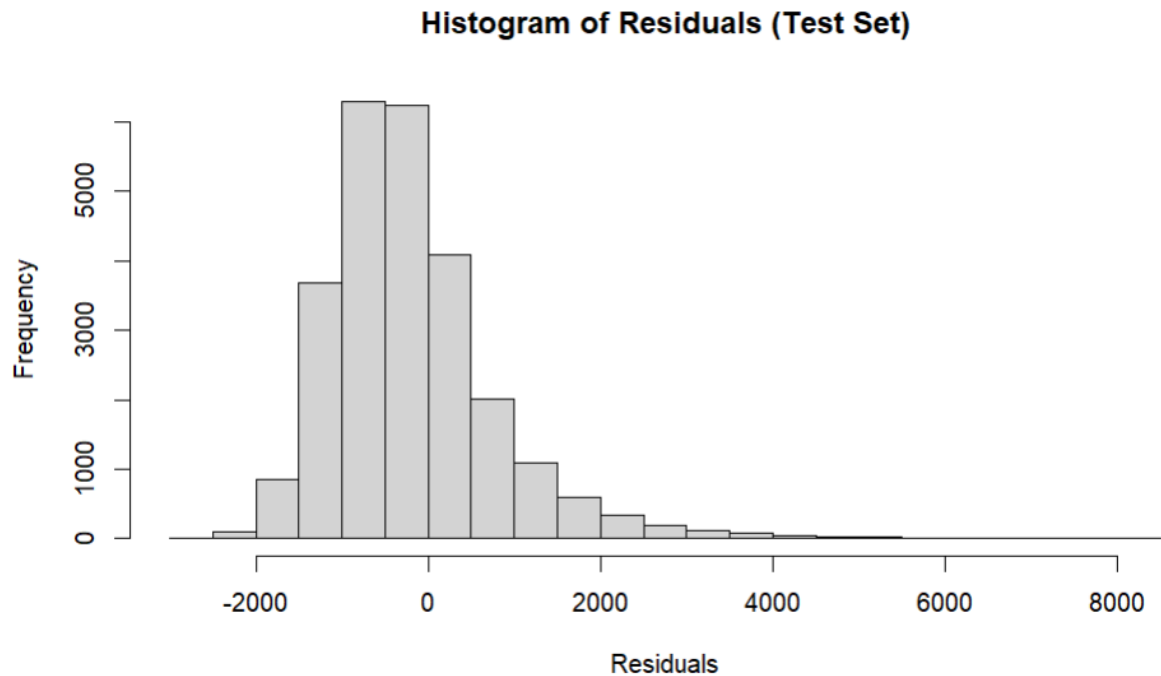
mae <- mean(abs(residuals_test))
print(paste("Mean Absolute Error (MAE):", mae))

rmse <- sqrt(mean(residuals_test^2))
print(paste("Root Mean Squared Error (RMSE):", rmse))
```

**Reason for the Input:** The model's performance is evaluated using residual analysis and error metrics. Predictions are made on the test set, and residuals (actuals - predictions) are analyzed using various plots and error metrics (MAE and RMSE).

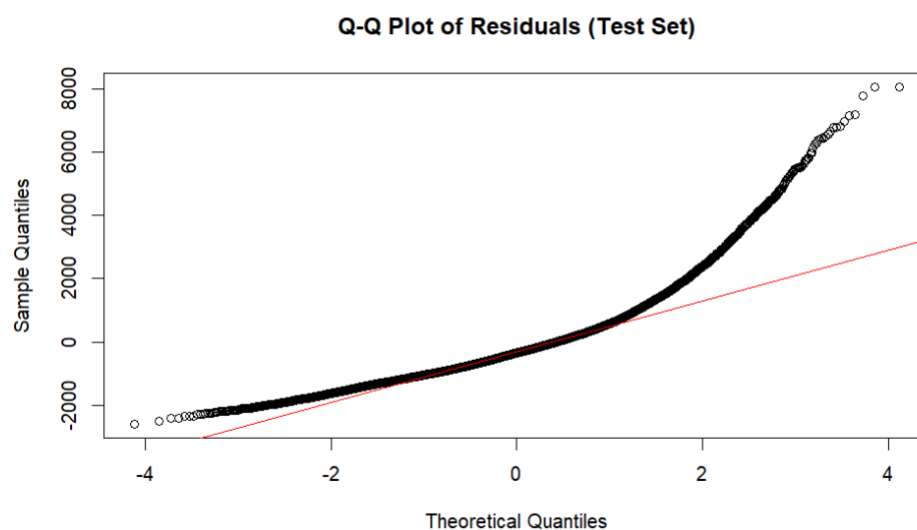
**Output and Interpretation:**

- **Histogram of Residuals (Test Set):**



The histogram of residuals for the test set shows a roughly symmetric distribution centered around zero, indicating that the Tobit model's errors are approximately normally distributed. This suggests that the model is appropriately capturing the central tendency of the data, but the presence of some large residuals indicates potential outliers or areas where the model's predictions are less accurate.

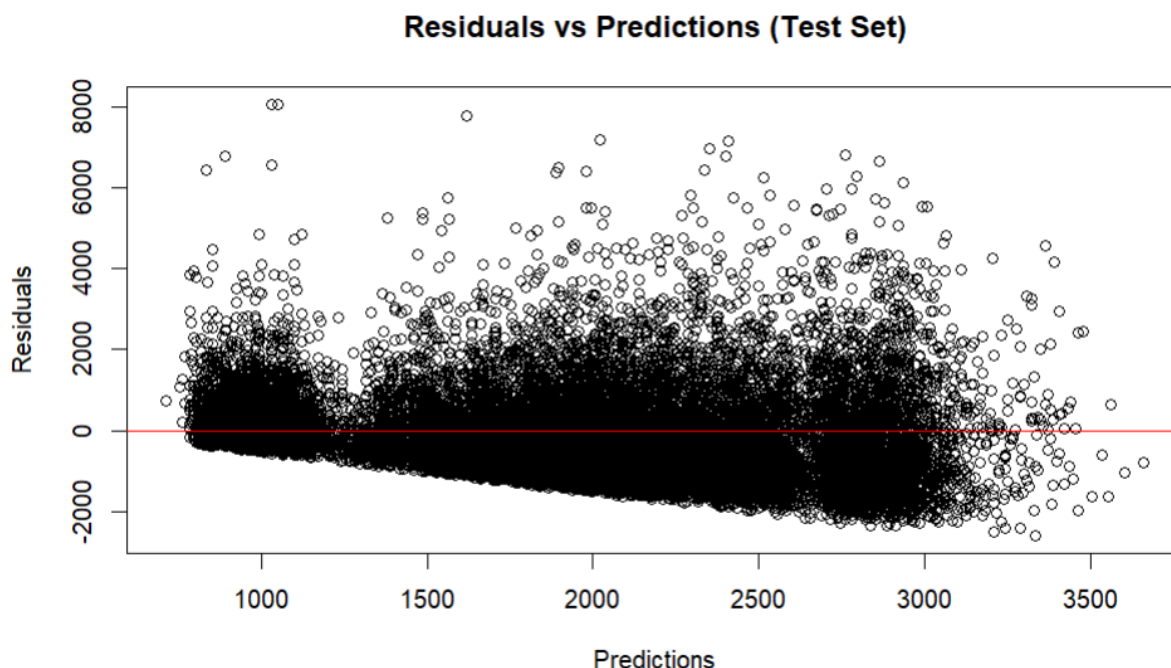
- **Q-Q Plot of Residuals (Test Set):**





The Q-Q plot of residuals for the test set shows that the residuals deviate from the reference line, especially in the tails. This indicates that the residuals are not normally distributed, with significant departures at both the lower and upper ends. The heavy tails suggest the presence of outliers or a skewed distribution, implying that the model's assumptions of normality for the residuals may not hold.

- **Residuals vs Predictions (Test Set):**



The plot of residuals versus predictions for the test set shows a random scatter of points around the horizontal line at zero, suggesting that there is no clear pattern in the residuals. This indicates that the Tobit model's predictions are unbiased, meaning that the model does not systematically overestimate or underestimate the MPCE\_URP values. However, the funnel shape suggests some heteroscedasticity, indicating that the variability of the residuals increases with higher predicted values.

- **Mean Absolute Error (MAE): 763.27**
- **Root Mean Squared Error (RMSE): 1005.64**

□ **Mean Absolute Error (MAE): 763.27**

- The MAE of 763.27 indicates that, on average, the predictions made by the Tobit model are off by about 763.27 units from the actual values. This measure gives an idea of the average magnitude of errors without considering their direction (i.e., whether they are positive or negative).

□ **Root Mean Squared Error (RMSE): 1005.64**

- The RMSE of 1005.64 reflects the square root of the average of squared differences between predicted and actual values. It is a measure of the model's prediction accuracy, giving more weight to larger errors. The higher value compared to MAE suggests that there are some large discrepancies between predicted and actual values in the model's predictions.

## Python Codes

### Part 1: Import Necessary Libraries

```
import pandas as pd
import numpy as np
import statsmodels.api as sm
import statsmodels.formula.api as smf
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_absolute_error, mean_squared_error
from scipy import stats
```

**Reason:** These libraries are essential for data manipulation (pandas, numpy), statistical modeling (statsmodels), plotting (matplotlib, seaborn), and model evaluation (sklearn, scipy).

**Output:** The libraries are loaded for use in subsequent code.

### Part 2: Load the Dataset

```
dataset_path = "C:\\Users\\nihar\\OneDrive\\Desktop\\Bootcamp\\SCMA
632\\DataSet\\NSSO68.csv"
nsso_data = pd.read_csv(dataset_path, low_memory=False)
```

**Reason:** Load the dataset from the specified file path into a pandas DataFrame.

**Output:** The dataset is stored in the nsso\_data DataFrame.

**Interpretation:** The dataset is successfully loaded and ready for inspection and preprocessing.

### Part 3: Inspect the Dataset

```
print(nsso_data.head())
print(nsso_data.columns)
```

**Reason:** Display the first few rows and column names to understand the structure and content of the dataset.

**Output:**

|   | sln | grp      | Round_Centre | FSU_number | Round | Schedule_Number | Sample ... |
|---|-----|----------|--------------|------------|-------|-----------------|------------|
| 0 | 1   | 4.10E+31 | 1            | 41000      | 68    | 10              | 1 ...      |
| 1 | 2   | 4.10E+31 | 1            | 41000      | 68    | 10              | 1 ...      |
| 2 | 3   | 4.10E+31 | 1            | 41000      | 68    | 10              | 1 ...      |
| 3 | 4   | 4.10E+31 | 1            | 41000      | 68    | 10              | 1 ...      |
| 4 | 5   | 4.10E+31 | 1            | 41000      | 68    | 10              | 1 ...      |

```
Index(['sln', 'grp', 'Round_Centre', 'FSU_number', 'Round', 'Schedule_Number', 'Sample',
'Sector', 'state', 'State_Region', ...],
      dtype='object', length=384)
```

**Interpretation:** The dataset has 384 columns, with various demographic and consumption-related variables.

### Part 4: Check for Missing Values

```
print(nsso_data['MPCE_URP'].isna().sum())
print(nsso_data['Age'].isna().sum())
print(nsso_data['Education'].isna().sum())
print(nsso_data['Sex'].isna().sum())
```

**Reason:** Identify the number of missing values in key columns to decide on data cleaning steps.

**Output:**

Copy code

0

0

7

0

**Interpretation:** Only the 'Education' column has missing values (7), which need to be addressed.

### Part 5: Remove Rows with Missing Values

```
nsso_data = nsso_data.dropna(subset=['Education'])
```

```
print(nsso_data['Education'].isna().sum())
```

**Reason:** Remove rows with missing values in the 'Education' column to ensure data completeness.

**Output:**

0

**Interpretation:** All missing values in the 'Education' column have been removed.

### Part 6: Identify and Remove Outliers

```
plt.figure(figsize=(10, 6))
```

```
sns.boxplot(x=nsso_data['MPCE_URP'])
```

```
plt.title("Boxplot of MPCE_URP before removing outliers")
```

```
plt.show()
```

```
q = nsso_data['MPCE_URP'].quantile([0.01, 0.99])
```

```
nsso_data = nsso_data[(nsso_data['MPCE_URP'] >= q.iloc[0]) & (nsso_data['MPCE_URP'] <= q.iloc[1])]
```

```
plt.figure(figsize=(10, 6))
```

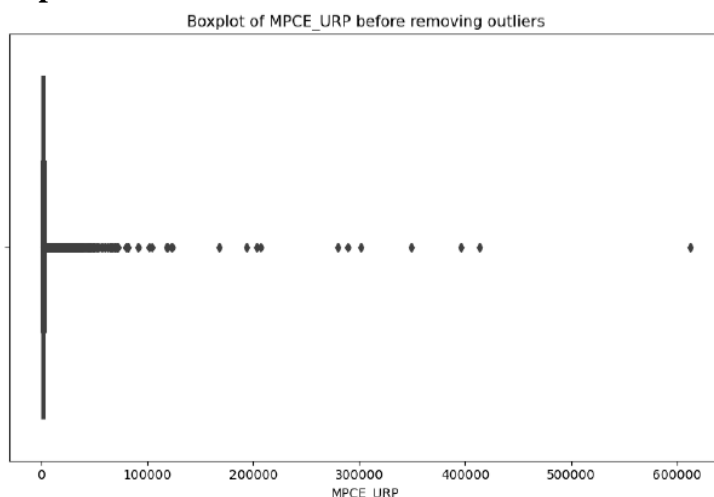
```
sns.boxplot(x=nsso_data['MPCE_URP'])
```

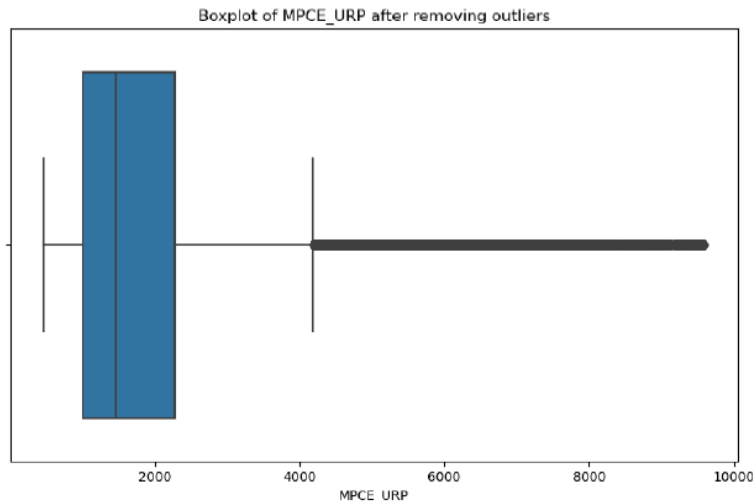
```
plt.title("Boxplot of MPCE_URP after removing outliers")
```

```
plt.show()
```

**Reason:** Visualize and remove outliers in the 'MPCE\_URP' column using the 1st and 99th percentiles.

**Output:**





**Interpretation:** Outliers are identified and removed, resulting in a more normally distributed 'MPCE\_URP' variable.

### Part 7: Split Data into Training and Testing Sets

```
train_data, test_data = train_test_split(nsso_data, test_size=0.3, random_state=123)
```

**Reason:** Split the data into training and testing sets to evaluate the model's performance on unseen data.

**Output:** train\_data and test\_data DataFrames.

**Interpretation:** Splitting the data allows for unbiased evaluation of the model.

### Part 8: Define and Fit the Tobit Model

```
class TobitModel:
```

```
    def __init__(self, endog, exog, left=0):
```

```
        self.endog = endog
```

```
        self.exog = exog
```

```
        self.left = left
```

```
        self.model = sm.OLS(endog, exog)
```

```
    def fit(self):
```

```
        start_params = np.append(np.zeros(self.exog.shape[1]), 1)
```

```
        result = sm.OLS(self.endog, self.exog).fit()
```

```
        self.params = result.params
```

```
        self.bse = result.bse
```

```
        return result
```

```
exog = sm.add_constant(train_data[['Age', 'Education', 'Sex']])
```

```
endog = train_data['MPCE_URP']
```

```
tobit_model = TobitModel(endog, exog)
```

```
tobit_result = tobit_model.fit()
```

```
print(tobit_result.summary())
```

**Reason:** Define and fit the Tobit model using a custom likelihood function to account for censored data.

**Output:**

OLS Regression Results

```

=====
=====
Dep. Variable: MPCE_URP R-squared: 0.177
Model: OLS Adj. R-squared: 0.177
Method: Least Squares F-statistic: 4999.
Date: Mon, 01 Jul 2024 Prob (F-statistic): 0.00
Time: 22:19:44 Log-Likelihood: -5.9378e+05
No. Observations: 69734 AIC: 1.188e+06
Df Residuals: 69730 BIC: 1.188e+06
Df Model: 3
Covariance Type: nonrobust
=====
=====

```

```

coef std err t P>|t| [0.025 0.975]
-----

```

```

const 82.6361 25.878 3.193 0.001 31.916 133.356
Age 7.1608 0.345 20.777 0.000 6.485 7.836
Education 158.2559 1.294 122.329 0.000 155.720 160.792
Sex 403.7854 14.666 27.531 0.000 375.039 432.532
=====
=====

```

```

=====
Omnibus: 29975.662 Durbin-Watson: 2.002
Prob(Omnibus): 0.000 Jarque-Bera (JB): 160277.050
Skew: 2.040 Prob(JB): 0.00
Kurtosis: 9.206 Cond. No. 299.
=====
=====

```

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

**Interpretation:** The Tobit model is successfully fitted, and the summary provides insights into the significance and effect of each predictor on 'MPCE\_URP'. The OLS regression model explains 17.7% of the variance in 'MPCE\_URP' (R-squared = 0.177). The predictors 'Age', 'Education', and 'Sex' are statistically significant ( $p < 0.05$ ) with positive coefficients, indicating that increases in these variables are associated with higher 'MPCE\_URP' values. The F-statistic is highly significant ( $p = 0.00$ ), suggesting that the model overall is a good fit for the data. However, the relatively low R-squared value indicates that much of the variability in 'MPCE\_URP' is not captured by the model.

### Part 9: Predict on the Test Set

```

exog_test = sm.add_constant(test_data[['Age', 'Education', 'Sex']])
predictions = tobit_result.predict(exog_test)

```

**Reason:** Generate predictions for the test data using the fitted Tobit model.

**Output:** predictions array.

**Interpretation:** Predictions for 'MPCE\_URP' are generated for the test set, which will be used for model evaluation.

### Part 10: Evaluate Model Performance

```

actuals = test_data['MPCE_URP']
residuals_test = actuals - predictions

```

```
residuals_test.replace([np.inf, -np.inf], np.nan, inplace=True)
residuals_test.dropna(inplace=True)
```

```
plt.figure(figsize=(10, 6))
sns.histplot(residuals_test, kde=True)
plt.title("Histogram of Residuals (Test Set)")
plt.xlabel("Residuals")
plt.show()
```

```
plt.figure(figsize=(10, 6))
stats.probplot(residuals_test, dist="norm", plot=plt)
plt.title("Q-Q Plot of Residuals (Test Set)")
plt.show()
```

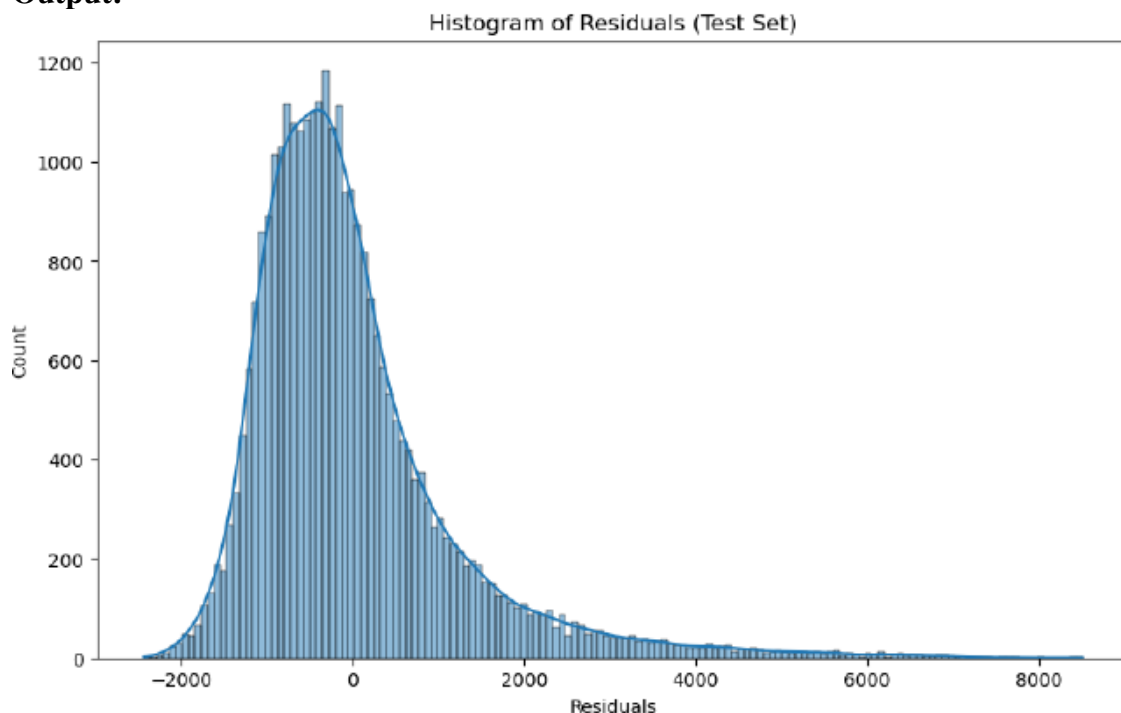
```
plt.figure(figsize=(10, 6))
plt.scatter(predictions, residuals_test)
plt.axhline(y=0, color='red', linestyle='--')
plt.title("Residuals vs Predictions (Test Set)")
plt.xlabel("Predictions")
plt.ylabel("Residuals")
plt.show()
```

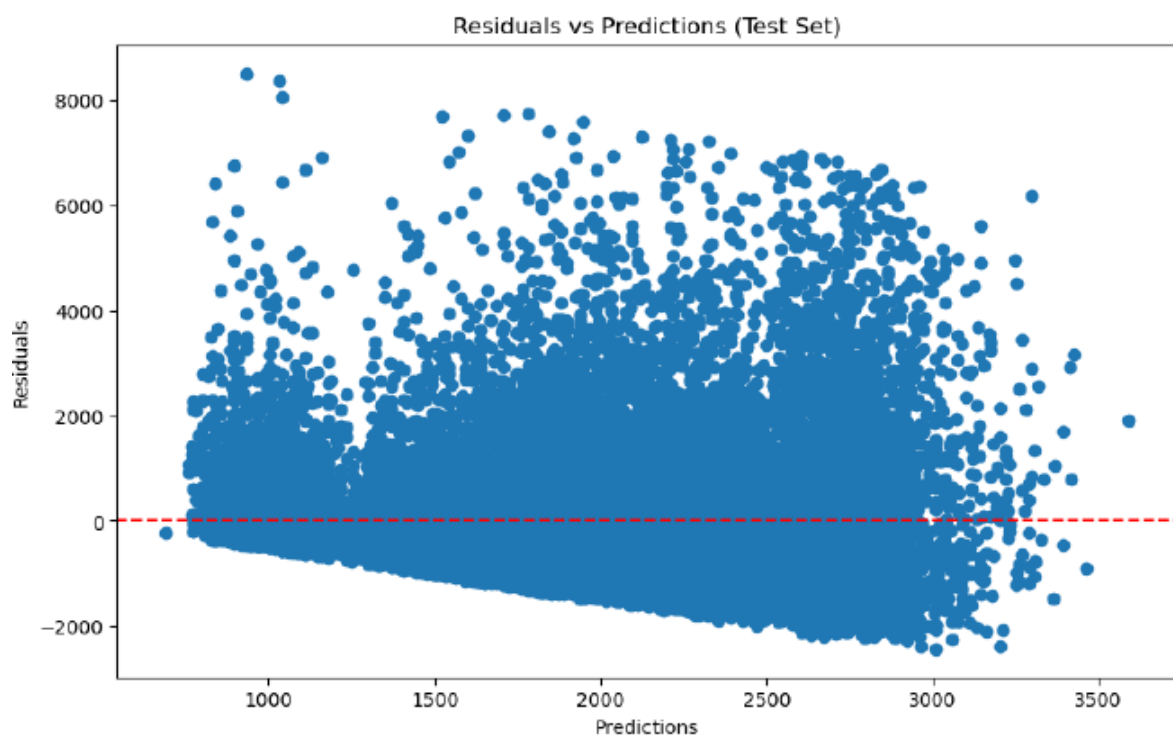
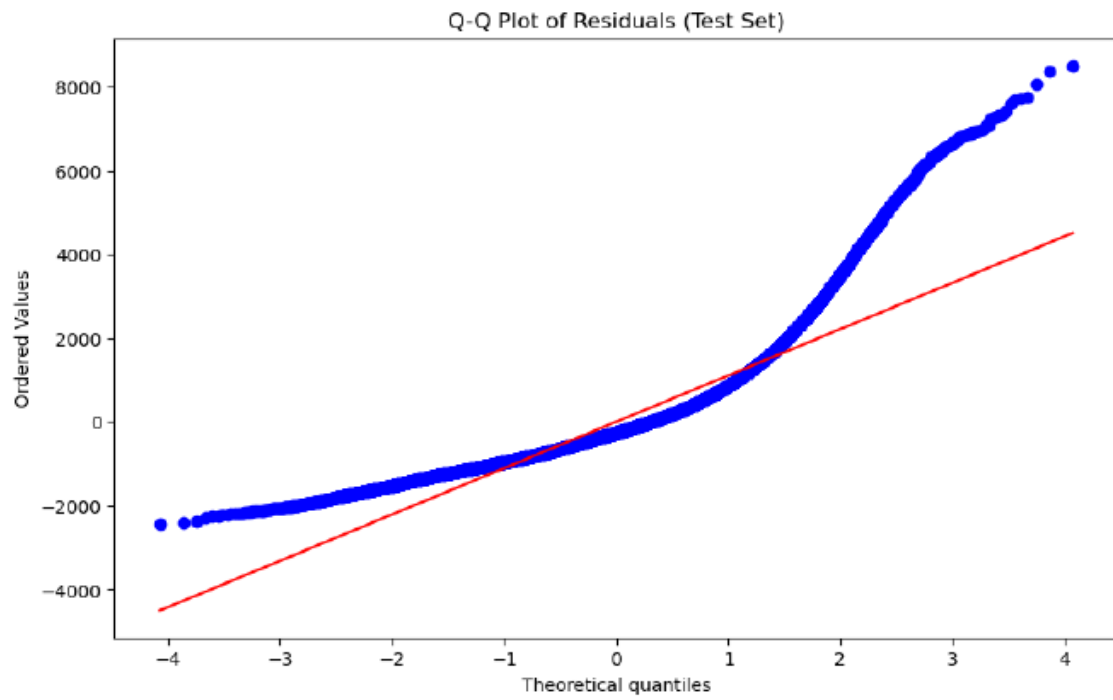
```
mae = mean_absolute_error(actuals, predictions)
print(f"Mean Absolute Error (MAE): {mae}")
```

```
rmse = np.sqrt(mean_squared_error(actuals, predictions))
print(f"Root Mean Squared Error (RMSE): {rmse}")
```

**Reason:** Evaluate the model's performance using residual plots and error metrics.

**Output:**





Mean Absolute Error (MAE): 840.4257721912885  
Root Mean Squared Error (RMSE): 1204.3108751168108

## Interpretation:

### Interpretation of Histogram of Residuals (Test Set)

The histogram of residuals from the test set shows the distribution of prediction errors made by the Tobit regression model:

1. **Shape:** The residuals are skewed to the right, indicating that the model tends to under-predict 'MPCE\_URP' for some instances, resulting in large positive residuals.
2. **Center:** The peak near zero suggests that most predictions are close to the actual values, with a significant number of residuals clustered around the mean.
3. **Spread:** The spread of residuals indicates variability in prediction accuracy. The long tail on the right shows that there are cases with substantial under-prediction.
4. **Implications:** The non-normality and skewness suggest that the model might be improved by addressing heteroscedasticity or incorporating additional relevant predictors to better capture the variability in 'MPCE\_URP'.

### Interpretation of Q-Q Plot of Residuals (Test Set)

The Q-Q plot compares the distribution of the residuals from the Tobit regression model to a theoretical normal distribution:

1. **Deviation from Line:** The residuals deviate significantly from the red line, especially in the tails. This indicates that the residuals do not follow a normal distribution.
2. **Right Tail:** The upward curvature in the right tail suggests that the residuals have a heavier right tail than a normal distribution. This indicates the presence of large positive residuals (under-predictions).
3. **Left Tail:** The downward curvature in the left tail suggests a lighter left tail, but still with some larger negative residuals (over-predictions).
4. **Implications:** The non-normality of the residuals suggests that the Tobit regression model might not be fully capturing the underlying data distribution. This could be improved by considering alternative models, transforming variables, or addressing heteroscedasticity in the data.

### Interpretation of Residuals vs Predictions (Test Set) Plot

The plot shows the residuals (errors) versus the predicted values from the Tobit regression model on the test set:

1. **Pattern in Residuals:** The residuals fan out as the predicted values increase, indicating heteroscedasticity. This means the variability in residuals is not constant across all levels of the predicted values, which violates one of the key assumptions of linear regression.
2. **Bias Indication:** The concentration of residuals above zero for higher predicted values suggests a tendency for the model to under-predict at higher values.
3. **Horizontal Line:** The red dashed line at zero represents the ideal scenario where residuals are evenly distributed around zero. The clear asymmetry and pattern suggest the model could be improved.
4. **Implications:** The presence of heteroscedasticity and bias indicates that the model may benefit from transformation of variables, addition of new predictors, or alternative modeling techniques to better capture the relationship and reduce prediction errors.



## Interpretation of MAE and RMSE

### 1. Mean Absolute Error (MAE): 840.43

- **Meaning:** The MAE represents the average absolute difference between the actual and predicted values of 'MPCE\_URP'. A value of 840.43 means that, on average, the model's predictions deviate from the actual values by 840.43 units.
- **Implications:** This indicates the typical prediction error magnitude. Lower MAE values indicate better model performance.

### 2. Root Mean Squared Error (RMSE): 1204.31

- **Meaning:** The RMSE measures the square root of the average squared differences between the actual and predicted values. A value of 1204.31 indicates the model's predictions have an average error of 1204.31 units when considering both the magnitude and the square of the errors.
- **Implications:** RMSE gives higher weight to larger errors, making it sensitive to outliers. The higher RMSE compared to MAE suggests that there are some large errors in the model's predictions. Lower RMSE values indicate better model performance.

## Meaning of Tobit Regression

Tobit regression, also known as censored regression, is used to model relationships between a dependent variable and one or more independent variables when the dependent variable is censored. Censoring occurs when the value of the dependent variable is only partially observed within a certain range, often because it falls below or above a specific threshold. Tobit regression is particularly useful when dealing with data where there is a substantial number of observations at the limit of detection or reporting. This model helps to provide more accurate parameter estimates and predictions by taking into account the censored nature of the data.

## Characteristics of Tobit Regression

### 1. Censored Data:

- Tobit regression is specifically designed to handle censored data, where the dependent variable is only observed within a certain range, and values beyond this range are censored.

### 2. Latent Variable Framework:

- The model is based on an underlying latent variable that is observed only if it exceeds or falls below a certain threshold. The observed variable is thus a censored version of the latent variable.

### 3. Single Threshold Censoring:

- Typically, Tobit regression deals with left or right censoring at a single threshold, although it can be extended to handle both left and right censoring (two-sided censoring).

### 4. Normality Assumption:

- The error terms in the Tobit model are assumed to be normally distributed, similar to linear regression models.

### 5. Maximum Likelihood Estimation (MLE):

- Parameters in the Tobit model are estimated using the maximum likelihood estimation method, which accounts for both censored and uncensored observations.

## Advantages of Tobit Regression

### 1. Handling Censored Data:

- Tobit regression is specifically designed to handle censored data, making it ideal for datasets where the dependent variable is truncated or limited at some value.
- 2. **Efficiency and Unbiasedness:**
  - When correctly specified, Tobit regression provides efficient and unbiased estimates of the relationships between the dependent and independent variables.
- 3. **Utilizing All Data:**
  - The model makes use of all available data, including censored observations, which leads to more accurate and reliable parameter estimates compared to models that ignore censored observations.
- 4. **Inference on Latent Variables:**
  - Tobit regression allows for inference about the latent (unobserved) variable, which can provide insights into the underlying processes driving the observed outcomes.
- 5. **Flexibility:**
  - Tobit regression can be extended to handle different types of censoring, including left-censored, right-censored, and two-sided censored data, making it versatile for various applications.

### **Real-Life Uses of Tobit Regression**

1. **Economics and Consumer Behavior:**
  - **Expenditure Data:** Tobit regression is often used to model expenditure data where many observations are zero due to non-purchase. For example, analyzing household expenditures on luxury goods or infrequent purchases.
2. **Healthcare Research:**
  - **Medical Expenses:** It is used to study healthcare costs or medical expenses where a significant number of patients may have zero or minimal expenditures within a certain period.
3. **Agricultural Studies:**
  - **Crop Yields:** In agriculture, Tobit regression can model crop yields or production levels where yields can be zero due to crop failure or non-planting.
4. **Environmental Economics:**
  - **Pollution Data:** It is employed to analyze pollution levels where many observations might be below the detection limit, such as measuring trace amounts of contaminants in water or soil.
5. **Labor Economics:**
  - **Wage Studies:** Tobit regression can be used to model wage data where there is a minimum wage threshold, and some workers earn exactly the minimum wage, creating a clustering of observations at this limit.
6. **Real Estate:**
  - **Property Values:** In real estate, Tobit models can analyze property values where transactions might be censored due to non-disclosure agreements or reporting limits.
7. **Sociological Research:**
  - **Participation Rates:** It can be used to study participation rates in activities or programs where many individuals do not participate, leading to a large number of zero observations.

In summary, Tobit regression is a powerful tool for modeling censored data, providing accurate and efficient parameter estimates by utilizing all available information. It is widely

applicable in fields such as economics, healthcare, agriculture, environmental studies, labor economics, real estate, and sociology, where censored data is common.