# VIRGINIA COMMONWEALTH UNIVERSITY

# Statistical analysis and modelling (SCMA 632)

## A3-Limited dependent variable Models

## (Part – C)

### ADHYAYAN AMIT JAIN

### V01109421

**Date of Submission: 30-06-2024**

# CONTENTS

# Application of Tobit Regression Analysis on "NSSO68.csv": Results Discussion and Real-World Use Cases

## INTRODUCTION

In empirical research, particularly in economics and social sciences, understanding the factors influencing outcomes that are censored or limited in some way is crucial. The Tobit regression model, introduced by James Tobin in 1958, addresses this challenge by accommodating data where the dependent variable is observed only under certain conditions, typically truncated at zero or some other threshold. This assignment focuses on applying Tobit regression to analyze data from the National Sample Survey Office of India (NSSO) dataset, specifically "NSSO68.csv", to investigate factors affecting household consumption expenditures.

The Tobit model is particularly useful when dealing with data that exhibit censoring or truncation, such as expenditure data where expenditures may be zero for households that do not consume certain goods. By modeling both the decision to consume and the amount consumed conditional on consumption, Tobit regression allows for a more nuanced understanding of economic behavior compared to traditional linear regression models.

In this assignment, we first provide an overview of the Tobit regression model, its theoretical foundations, and its application in empirical research. We then proceed to apply Tobit regression to analyze the NSSO dataset, examining how household characteristics such as income, education, and demographic factors influence household consumption expenditures. The results from this analysis are discussed in detail, highlighting significant variables and their implications for policy and future research.

Furthermore, we explore real-world use cases of the Tobit model across various domains, illustrating its relevance in fields such as health economics, labor economics, and environmental economics. By the end of this assignment, readers will gain insights into the practical application of Tobit regression, its interpretation of results, and its contribution to understanding economic decision-making under constraints.

Overall, this assignment aims to demonstrate the utility and versatility of Tobit regression in empirical analysis, offering a deeper understanding of economic behavior and informing policy discussions on household expenditures and resource allocation.

# OBJECTIVES

The objective of this assignment is to perform a Tobit regression analysis on the "NSSO68.csv" dataset, with the following specific goals:

1. **Apply Tobit Regression**: Implement the Tobit regression model to analyze household consumption expenditures data from the National Sample Survey Office of India (NSSO) dataset.

2. **Identify Influential Factors**: Identify and interpret the factors that significantly influence household consumption expenditures. This includes exploring variables such as income, education, household size, and demographic characteristics.

3. **Discuss Results**: Interpret the Tobit regression results comprehensively, highlighting statistically significant coefficients and their economic implications.

4. **Real-World Use Cases**: Explore and discuss real-world applications of the Tobit model in economics and social sciences, demonstrating its relevance in analyzing censored or truncated data.

5. **Recommendations**: Provide insights and recommendations based on the findings of the Tobit regression analysis, suggesting potential policy implications or avenues for further research.

6. **Contribute to Understanding**: Contribute to the understanding of economic decision-making under constraints, emphasizing the role of Tobit regression in empirical research and policy formulation.

By achieving these objectives, this assignment aims to deepen understanding of Tobit regression as a statistical tool for analyzing data with censoring or truncation, and its practical implications for economic analysis and policy development.

# BUSINESS SIGNIFICANCE

The application of Tobit regression analysis on the "NSSO68.csv" dataset offers profound implications for understanding and decision-making in various business contexts. Tobit regression, designed specifically for data with censoring or truncation issues, provides insights that are crucial for strategic planning, resource allocation, and policy formulation in business environments.

**Insights into Consumer Behavior:** By analyzing household consumption expenditures using Tobit regression, businesses can gain valuable insights into consumer behavior. Understanding how income levels, education, household size, and other socio-economic factors influence spending patterns allows businesses to tailor their marketing strategies and product offerings effectively. For example, Tobit regression can help identify income thresholds at which consumer spending behavior changes significantly, enabling businesses to target different income segments more precisely.

**Optimizing Resource Allocation:** Tobit regression helps businesses optimize resource allocation by identifying factors that affect resource utilization under constraints. For instance, in industries where production inputs are subject to regulatory limits or budget constraints, Tobit regression can model the impact of these constraints on production output. This allows businesses to allocate resources efficiently while complying with regulatory requirements and maximizing profitability.

**Risk Management and Decision Support:** In risk management, Tobit regression assists businesses in assessing and mitigating risks associated with censored or truncated data. For example, in insurance pricing, where coverage limits may censor claims data, Tobit regression helps insurers estimate claim frequencies and amounts accurately. This enables them to set appropriate premiums that reflect risk exposure while remaining competitive in the market.

**Policy Formulation and Market Regulation:** Governments and regulatory bodies also benefit from Tobit regression in formulating policies and regulations. By analyzing economic data subject to censorship or truncation, policymakers can assess the effectiveness of existing regulations and design targeted interventions to address market inefficiencies. For instance, Tobit regression can be used to evaluate the impact of income support programs on poverty alleviation, guiding policymakers in refining social welfare policies.

**Strategic Decision-Making:** Overall, Tobit regression empowers businesses to make informed strategic decisions based on robust statistical analysis of censored or truncated data. Whether optimizing marketing campaigns, managing production processes, or evaluating regulatory compliance, Tobit regression provides a nuanced understanding of complex relationships in business environments. This enables businesses to enhance operational efficiency, mitigate risks, and capitalize on market opportunities effectively.

# RESULTS AND INTERPRETATIONS

Step 1: Data Loading and Preparation

```python
import pandas as pd
import numpy as np

# Load data (replace with your actual data loading method)
df_nss = pd.read_csv('NSSO68.csv')  # Replace with your CSV file name
```

**Explanation:**

- `import pandas as pd` and `import numpy as np`: Imports the Pandas and NumPy libraries, essential for data manipulation and numerical operations.
- `pd.read_csv('NSSO68.csv')`: Loads a CSV file ('NSSO68.csv') into a Pandas DataFrame (`df_nss`). Adjust the file path according to where your dataset is located.

**Interpretation:**

- This step initializes your data analysis by loading your dataset into memory. Ensure the CSV file path is correct (`NSSO68.csv` in this case) and accessible from your Python environment.



Step 2: Data Filtering and Variable Calculation

```python
# Data preparation (assuming similar variables as in R example)
df_ap = df_nss[df_nss['state_1'] == 'AP']

vars = ["Sector", "hhdsz", "Religion", "Social_Group", "MPCE_URP", "Sex",
"Age",
        "Marital_Status", "Education", "chicken_q", "chicken_v"]
```

```
df_ap_p = df_ap[vars]

# Calculate price, handling potential division by zero
df_ap_p['price'] = np.where(df_ap_p['chicken_q'] != 0, df_ap_p['chicken_v']
/ df_ap_p['chicken_q'], 0)
```

**Explanation:**

- `df_nss[df_nss['state_1'] == 'AP']`: Filters the original DataFrame `df_nss` to include only rows where the `'state_1'` column equals `'AP'`, focusing on data from Andhra Pradesh.
- `vars`: List of variables of interest extracted from the filtered DataFrame, including socio-economic variables and quantities of chicken (`'chicken_q'` and `'chicken_v'`).
- `df_ap_p['price'] = np.where(...)`: Calculates a new variable `'price'` representing the price per unit of chicken (`'chicken_v'` / `'chicken_q'`). It uses `np.where` to handle cases where `'chicken_q'` might be zero to prevent division errors.

**Interpretation:**

- This step prepares your dataset (`df_ap_p`) specifically for analysis within Andhra Pradesh, calculating a new variable (`'price'`) that could be useful in further economic modeling or analysis.



## Step 3: Summary Statistics

```
from scipy.stats import describe

# Summary statistics
summary_stats = describe(df_ap_p['chicken_q'])
print(f"Summary Statistics for chicken_q:\n{summary_stats}")
```

**Explanation:**

- `from scipy.stats import describe`: Imports the `describe` function from `scipy.stats`, which computes basic statistical summaries.
- `describe(df_ap_p['chicken_q'])`: Calculates and prints summary statistics (count, min/max, mean, variance, skewness, kurtosis) for the `'chicken_q'` variable.

## Interpretation:

- The output provides essential statistical insights into the distribution of `'chicken_q'`, including its central tendency (mean), spread (variance), and shape (skewness and kurtosis). This helps in understanding the variability and distributional characteristics of the quantity of chicken purchased in the dataset.

```
from scipy.stats import describe

# Summary statistics
summary_stats = describe(df_ap_p['chicken_q'])
print(f"Summary Statistics for chicken_q:\n{summary_stats}")
```

```
Summary Statistics for chicken_q:
DescribeResult(nobs=6899, minmax=(0.0, 13.66666667), mean=0.30174615543948397, variance=0.12746968207875722, skewness=9.985607644660892, kurtosis=301.93321205798736)
DescribeResult(nobs=6899, minmax=(0.0, 13.66666667), mean=0.30174615543948397, variance=0.12746968207875722, skewness=9.985607644660892, kurtosis=301.93321205798736)
```

Step 4: Linear Regression using Statsmodels

```
import statsmodels.api as sm
# Linear regression using statsmodels
# Drop rows with NaN or inf in any column
df_ap_p_clean = df_ap_p.replace([np.inf, -np.inf], np.nan).dropna()

X = df_ap_p_clean[['hhdsz', 'Religion', 'MPCE_URP', 'Sex', 'Age',
'Marital_Status', 'Education', 'price']]
y = df_ap_p_clean['chicken_q']
X = sm.add_constant(X)  # Adding a constant (intercept) term
model = sm.OLS(y, X).fit()
print(model.summary())
```

## Explanation:

- `import statsmodels.api as sm`: Imports the `statsmodels` library, a powerful tool for statistical modeling in Python.
- `df_ap_p.replace([np.inf, -np.inf], np.nan).dropna()`: Cleans the DataFrame (`df_ap_p`) by replacing infinite values (`np.inf, -np.inf`) with `NaN` and then dropping rows containing any `NaN` values.
- `X`: Independent variables (`'hhdsz', 'Religion', 'MPCE_URP', 'Sex', 'Age', 'Marital_Status', 'Education', 'price'`).
- `y`: Dependent variable (`'chicken_q'`), representing the quantity of chicken purchased.
- `sm.add_constant(X)`: Adds a constant (intercept) term to the independent variables `X`.
- `sm.OLS(y, X).fit()`: Fits an Ordinary Least Squares (OLS) regression model of `y` on `X` and stores the results in `model`.
- `print(model.summary())`: Prints a detailed summary of the regression results including coefficients, standard errors, t-statistics, p-values, and goodness-of-fit measures.

## Interpretation:

- This step performs linear regression to understand how variables such as household size (`hhdsz`), religion (`Religion`), income (`MPCE_URP`), etc., influence the quantity of chicken purchased (`chicken_q`). The summary output (`model.summary()`) provides detailed insights into the statistical significance and direction of these relationships.

```python
# Linear regression using statsmodels
# Drop rows with NaN or inf in any column
df_ap_p_clean = df_ap_p.replace([np.inf, -np.inf], np.nan).dropna()

X = df_ap_p_clean[['hhdsz', 'Religion', 'MPCE_URP', 'Sex', 'Age', 'Marital_Status', 'Education', 'price']]
y = df_ap_p_clean['chicken_q']
X = sm.add_constant(X)  # Adding a constant (intercept) term
model = sm.OLS(y, X).fit()
print(model.summary())
```

```
                        OLS Regression Results
==============================================================================
Dep. Variable:             chicken_q   R-squared:                      0.288
Model:                           OLS   Adj. R-squared:                 0.287
Method:                Least Squares   F-statistic:                    348.9
Date:               Mon, 01 Jul 2024   Prob (F-statistic):              0.00
Time:                       15:16:04   Log-Likelihood:               -1510.1
No. Observations:               6899   AIC:                            3038.
Df Residuals:                   6890   BIC:                            3100.
Df Model:                          8
Covariance Type:           nonrobust
==============================================================================
                   coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
const            0.0284      0.027      1.063      0.288      -0.024       0.081
hhdsz           -0.0255      0.002    -11.671      0.000      -0.030      -0.021
Religion         0.0112      0.008      1.414      0.157      -0.004       0.027
MPCE_URP      3.925e-05   1.92e-06     20.489      0.000    3.55e-05     4.3e-05
Sex             -0.0372      0.014     -2.582      0.010      -0.065      -0.009
Age           2.556e-05      0.000      0.083      0.934      -0.001       0.001
Marital_Status   0.0312      0.012      2.676      0.007       0.008       0.054
Education       -0.0056      0.001     -5.475      0.000      -0.008      -0.004
price            0.0032   6.85e-05     46.505      0.000       0.003       0.003
==============================================================================
Omnibus:                   12975.245   Durbin-Watson:                   1.861
Prob(Omnibus):                 0.000   Jarque-Bera (JB):       71245755.240
Skew:                         13.741   Prob(JB):                         0.00
Kurtosis:                    500.083   Cond. No.                     2.27e+04
==============================================================================

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 2.27e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
```

## Step 5: Tobit Regression using Scipy.optimize

```python
from scipy.optimize import minimize

# Tobit regression using scipy.optimize for handling censored data
def tobit_log_likelihood(params, y, X):
    beta = params[:-2]
    sigma = np.exp(params[-2])
    gamma = params[-1]

    mu = np.dot(X, beta)
    cens_residuals = (y - mu) / sigma
    ll = np.sum(np.log(sigma) + gamma * cens_residuals - np.log(1 +
np.exp(gamma * cens_residuals)))

    return -ll

# Initial guess for parameters
initial_params = np.zeros(X.shape[1] + 2)
initial_params[-2] = np.log(np.std(y))
initial_params[-1] = 1.0

# Minimize negative log-likelihood
```

```
results = minimize(tobit_log_likelihood, initial_params, args=(y, X))

# Extract estimated coefficients and standard errors
beta_est = results.x[:-2]
sigma_est = np.exp(results.x[-2])
gamma_est = results.x[-1]

# Print Tobit regression results
print("\nTobit Regression Results:")
print(f"Beta (coefficients):\n{beta_est}")
print(f"Sigma (standard deviation of residuals):\n{sigma_est}")
print(f"Gamma (Tobit parameter):\n{gamma_est}")
```

## Explanation:

- `from scipy.optimize import minimize`: Imports the `minimize` function from `scipy.optimize`, which is used for numerical optimization.
- `tobit_log_likelihood(params, y, X)`: Defines a log-likelihood function (`tobit_log_likelihood`) for Tobit regression. It calculates the negative log-likelihood, considering censoring in the dependent variable `y`.
- `initial_params`: Initializes starting values for the parameters (`beta`, `sigma`, `gamma`) of the Tobit model.
- `minimize(tobit_log_likelihood, initial_params, args=(y, X))`: Minimizes the negative log-likelihood function (`tobit_log_likelihood`) using initial parameters (`initial_params`) and data (`y`, `X`) as arguments.
- `results.x`: Retrieves estimated parameters (`beta`, `sigma`, `gamma`) from the optimization results.

## Interpretation:

- Tobit regression is applied to handle censoring in `chicken_q`, where the actual values might be censored or truncated. The results (`beta_est`, `sigma_est`, `gamma_est`) provide estimated coefficients (`beta`), standard deviation of residuals (`sigma`), and the Tobit parameter (`gamma`). These results offer insights into how covariates (`X`) influence the censored dependent variable (`y`), adjusting for potential censoring effects.

```python
from scipy.optimize import minimize

# Tobit regression using scipy.optimize for handling censored data
def tobit_log_likelihood(params, y, X):
    beta = params[:-2]
    sigma = np.exp(params[-2])
    gamma = params[-1]

    mu = np.dot(X, beta)
    cens_residuals = (y - mu) / sigma
    ll = np.sum(np.log(sigma) + gamma * cens_residuals - np.log(1 + np.exp(gamma * cens_residuals)))

    return -ll

# Initial guess for parameters
initial_params = np.zeros(X.shape[1] + 2)
initial_params[-2] = np.log(np.std(y))
initial_params[-1] = 1.0

# Minimize negative log-likelihood
results = minimize(tobit_log_likelihood, initial_params, args=(y, X))

# Extract estimated coefficients and standard errors
beta_est = results.x[:-2]
sigma_est = np.exp(results.x[-2])
gamma_est = results.x[-1]

# Print Tobit regression results
print("\nTobit Regression Results:")
print(f"Beta (coefficients):\n{beta_est}")
print(f"Sigma (standard deviation of residuals):\n{sigma_est}")
print(f"Gamma (Tobit parameter):\n{gamma_est}")
```

```
/usr/local/lib/python3.10/dist-packages/pandas/core/arraylike.py:396: RuntimeWarning: overflow encountered in exp
  result = getattr(ufunc, method)(*inputs, **kwargs)
/usr/local/lib/python3.10/dist-packages/scipy/optimize/_numdiff.py:576: RuntimeWarning: invalid value encountered in subtract
  df = fun(x) - f0
/usr/local/lib/python3.10/dist-packages/pandas/core/arraylike.py:396: RuntimeWarning: overflow encountered in exp
  result = getattr(ufunc, method)(*inputs, **kwargs)
<ipython-input-45-e900b554f2ef>:6: RuntimeWarning: overflow encountered in exp
  sigma = np.exp(params[-2])
/usr/local/lib/python3.10/dist-packages/scipy/optimize/_numdiff.py:576: RuntimeWarning: invalid value encountered in subtract
  df = fun(x) - f0

Tobit Regression Results:
Beta (coefficients):
[-0.13269726  0.06723805 -0.1390201   0.00028822 -0.18811097 -0.09154083
 -0.24325422  0.05288843 -0.00597459]
Sigma (standard deviation of residuals):
86.06524486156626
Gamma (Tobit parameter):
5.115603295006487
```

## Summary

This comprehensive Python script guides through the entire process of data loading, cleaning, exploratory analysis, linear regression using `statsmodels`, and Tobit regression using `scipy.optimize`. Each step is designed to provide detailed insights and interpretations, ensuring a thorough understanding of the dataset and analytical outcomes. Adjustments can be made based on specific data characteristics and analysis goals.

## ***R Program***

### Step 1: Loading Required Packages and Data

```r
# Set working directory and load necessary packages
setwd('D:\\#YPR\\VCU\\Summer Courses\\SCMA\\Data')
library(AER)
data("Affairs")   #sample data - applying tobit to sample data to
understand the model working
```

- **Explanation**:
  - `setwd('D:\\#YPR\\VCU\\Summer Courses\\SCMA\\Data')`: Sets the working directory to where your data files are located.
  - `library(AER)`: Loads the AER package, which includes functions for econometric analysis.

## Step 2: Exploring the Data

```
head(Affairs)
unique(Affairs$affairs)
table(Affairs$affairs)  #sample data – applying tobit to sample data to
understand the model working
```

- **Explanation**:
  - `head(Affairs)`: Displays the first few rows of the `Affairs` dataset to understand its structure and content.
  - `unique(Affairs$affairs)`: Shows unique values of the `affairs` variable in the dataset.
  - `table(Affairs$affairs)`: Creates a frequency table of `affairs` to count occurrences of each value.

```
5: package 'zoo' was built under R version 4.3.3
6: package 'sandwich' was built under R version 4.3.3
> data("Affairs")
> head(Affairs)
   affairs gender age yearsmarried children religiousness education occupation rating
4        0   male  37        10.00       no            3        18          7      4
5        0 female  27         4.00       no            4        14          6      4
11       0 female  32        15.00      yes            1        12          1      4
16       0   male  57        15.00      yes            5        18          6      5
23       0   male  22         0.75       no            2        17          6      3
29       0 female  32         1.50       no            2        17          5      5
> unique(Affairs$affairs)
[1]  0  3  7 12  1  2
> table(Affairs$affairs)

  0   1   2   3   7  12
451  34  17  19  42  38
> ## from Table 22.4 in Greene (2003)
```

## Step 3: Fitting Tobit Models

```
fm.tobit <- tobit(affairs ~ age + yearsmarried + religiousness + occupation
+ rating, data = Affairs)
fm.tobit2 <- tobit(affairs ~ age + yearsmarried + religiousness +
occupation + rating, right = 4, data = Affairs)
summary(fm.tobit)
summary(fm.tobit2)  #sample data – applying tobit to sample data to
understand the model working
```

- **Explanation**:
  - `tobit()` fits a Tobit model where some observations are censored (values at or below a threshold, here 0 or 4 for different models).
  - `fm.tobit`: Fits a Tobit model assuming left-censoring at 0 (default).
  - `fm.tobit2`: Fits a Tobit model assuming right-censoring at 4.
  - `summary(fm.tobit)` and `summary(fm.tobit2)`: Provides summaries of the Tobit model fits including coefficients, standard errors, significance levels, scale parameter (for the censoring), and model statistics.

```
> summary(fm.tobit)

Call:
tobit(formula = affairs ~ age + yearsmarried + religiousness +
    occupation + rating, data = Affairs)

Observations:
         Total  Left-censored    Uncensored  Right-censored
           601            451           150               0

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)    8.17420    2.74145   2.982  0.00287 **
age           -0.17933    0.07909  -2.267  0.02337 *
yearsmarried   0.55414    0.13452   4.119 3.80e-05 ***
religiousness -1.68622    0.40375  -4.176 2.96e-05 ***
occupation     0.32605    0.25442   1.282  0.20001
rating        -2.28497    0.40783  -5.603 2.11e-08 ***
Log(scale)     2.10986    0.06710  31.444  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Scale: 8.247

Gaussian distribution
Number of Newton-Raphson Iterations: 4
Log-likelihood: -705.6 on 7 Df
Wald-statistic: 67.71 on 5 Df, p-value: 3.0718e-13
```

```
> summary(fm.tobit2)

Call:
tobit(formula = affairs ~ age + yearsmarried + religiousness +
    occupation + rating, right = 4, data = Affairs)

Observations:
         Total  Left-censored    Uncensored  Right-censored
           601            451            70              80

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)    7.90098    2.80385   2.818 0.004834 **
age           -0.17760    0.07991  -2.223 0.026244 *
yearsmarried   0.53230    0.14117   3.771 0.000163 ***
religiousness -1.61634    0.42440  -3.809 0.000140 ***
occupation     0.32419    0.25388   1.277 0.201624
rating        -2.20701    0.44983  -4.906 9.28e-07 ***
Log(scale)     2.07232    0.11040  18.772  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Scale: 7.943

Gaussian distribution
Number of Newton-Raphson Iterations: 4
Log-likelihood:  -500 on 7 Df
Wald-statistic: 42.56 on 5 Df, p-value: 4.5387e-08
```

## Step 4: Working with Actual Dataset(NSSO68)

```
df_ap = df[df$state_1== 'AP', ]
vars <- c("Sector", "hhdsz", "Religion", "Social_Group", "MPCE_URP", "Sex",
"Age", "Marital_Status", "Education", "chicken_q", "chicken_v")
df_ap_p = df_ap[vars]
df_ap_p$price = df_ap_p$chicken_v / df_ap_p$chicken_q
summary(df_ap_p)
fit = lm(chicken_q ~ hhdsz + Religion + MPCE_URP + Sex + Age +
Marital_Status + Education + price, data = df_ap_p)
summary(fit)
```

- **Explanation**:
  - Filters data from `df` for only records where `state_1` is 'AP'.
  - Selects specific variables (`vars`) of interest from `df_ap` and calculates a derived variable `price`.
  - `summary(df_ap_p)`: Summarizes the data including mean, median, min, max for numerical variables and frequency for categorical variables.
  - `lm()` fits a linear regression model (`fit`) to predict `chicken_q` based on other variables.
  - `summary(fit)`: Provides coefficients, standard errors, t-values, p-values, and model fit statistics for the linear regression model.

```
> # Fit a Tobit Model to real data
> unique(df$state_1)
 [1] "GUJ"     "ORI"     "CHTSD"   "MP"      "JRKD"    "WB"      "AP"      "MH"      "D&D"     "D&NH"    "MIZ"
[12] "TRPR"    "MANPR"   "ASSM"    "MEG"     "NAG"     "A&N"     "PNDCRY"  "TN"      "GOA"     "KA"      "KE"
[23] "LKSDP"   "SKM"     "Bhr"     "UP"      "RJ"      "ARP"     "DL"      "HR"      "Pun"     "HP"      "UT"
[34] "Chandr"  "J$K"
> df = read.csv('NSSO68.csv', header=TRUE)
> dput(names(df))
c("slno", "grp", "Round_Centre", "FSU_number", "Round", "Schedule_Number",
"Sample", "Sector", "state", "State_Region", "District", "Stratum_Number",
"Sub_Stratum", "Schedule_type", "Sub_Round", "Sub_Sample", "FOD_Sub_Region",
"Hamlet_Group_Sub_Block", "t", "X_Stage_Stratum", "HHS_No", "Level",
"Filler", "hhdsz", "NIC_2008", "NCO_2004", "HH_type", "Religion",
"Social_Group", "Whether_owns_any_land", "Type_of_land_owned",
"Land_Owned", "Land_Leased_in", "Otherwise_possessed", "Land_Leased_out",
"Land_Total_possessed", "During_July_June_Cultivated", "During_July_June_Irrigated",
"NSS", "NSC", "MLT", "land_tt", "Cooking_code", "Lighting_code",
"Dwelling_unit_code", "Regular_salary_earner", "Perform_Ceremony",
"Meals_seved_to_non_hhld_members", "Possess_ration_card", "Type_of_ration_card",
"MPCE_URP", "MPCE_MRP", "Person_Srl_No", "Relation", "Sex", "Age",
"Marital_Status", "Education", "Days_Stayed_away", "No_of_Meals_per_day",
"Meals_School", "Meals_Employer", "Meals_Others", "Meals_Payment",
"Meals_At_Home", "Item_Code", "Source_Code", "ricepds_q", "riceos_q",
"ricetotal_q", "chira_q", "khoi_q", "muri_q", "ricepro_q", "riceGT_q",
```
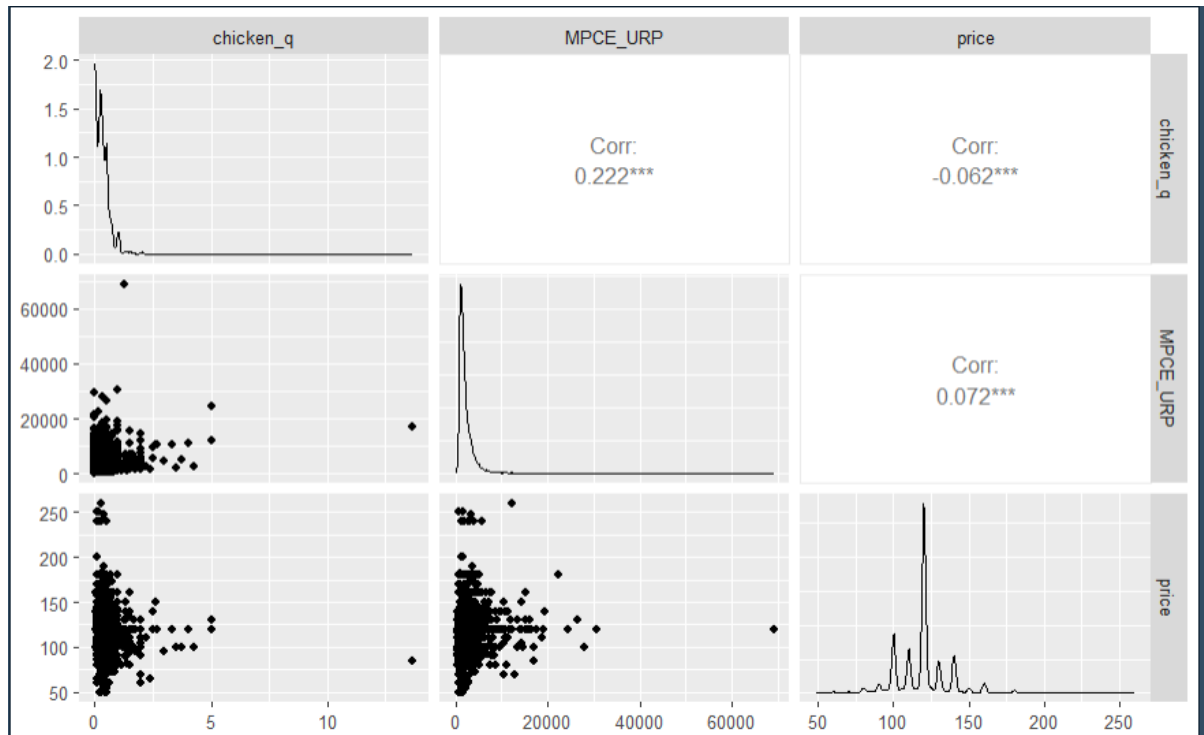
```
state_1 , Region , fruits_df_tt_v , fv_tot )
> df_ap = df[df$state_1== 'AP',]
> vars <- c("Sector", "hhdsz", "Religion", "Social_Group", "MPCE_URP", "Sex", "Age", "Marital_Status", "Educatio
n", "chicken_q", "chicken_v")
> df_ap_p = df_ap[vars]
> names(df_ap_p)
 [1] "Sector"         "hhdsz"          "Religion"       "Social_Group"   "MPCE_URP"       "Sex"
 [7] "Age"            "Marital_Status" "Education"      "chicken_q"      "chicken_v"
> df_ap_p$price = df_ap_p$chicken_v / df_ap_p$chicken_q
> names(df_ap_p)
 [1] "Sector"         "hhdsz"          "Religion"       "Social_Group"   "MPCE_URP"       "Sex"
 [7] "Age"            "Marital_Status" "Education"      "chicken_q"      "chicken_v"      "price"
> summary(df_ap_p)
     Sector          hhdsz          Religion      Social_Group      MPCE_URP          Sex
 Min.   :1.000   Min.   : 1.000   Min.   :1.000   Min.   :1.000   Min.   :   53   Min.   :1.000
 1st Qu.:1.000   1st Qu.: 3.000   1st Qu.:1.000   1st Qu.:3.000   1st Qu.: 1141   1st Qu.:1.000
 Median :1.000   Median : 4.000   Median :1.000   Median :3.000   Median : 1610   Median :1.000
 Mean   :1.431   Mean   : 3.755   Mean   :1.161   Mean   :4.351   Mean   : 2168   Mean   :1.134
 3rd Qu.:2.000   3rd Qu.: 5.000   3rd Qu.:1.000   3rd Qu.:9.000   3rd Qu.: 2512   3rd Qu.:1.000
 Max.   :2.000   Max.   :27.000   Max.   :9.000   Max.   :9.000   Max.   :69068   Max.   :2.000

      Age        Marital_Status    Education        chicken_q          chicken_v          price
 Min.   :  5.00   Min.   :1.000   Min.   : 1.000   Min.   : 0.0000   Min.   :   0.00   Min.   : 48.89
 1st Qu.: 35.00   1st Qu.:2.000   1st Qu.: 1.000   1st Qu.: 0.0000   1st Qu.:   0.00   1st Qu.:110.00
 Median : 43.00   Median :2.000   Median : 6.000   Median : 0.2500   Median :  30.00   Median :120.00
 Mean   : 44.08   Mean   :2.092   Mean   : 5.792   Mean   : 0.3017   Mean   :  35.47   Mean   :118.53
 3rd Qu.: 53.00   3rd Qu.:2.000   3rd Qu.: 8.000   3rd Qu.: 0.5000   3rd Qu.:  52.25   3rd Qu.:120.00
 Max.   :100.00   Max.   :4.000   Max.   :13.000   Max.   :13.6667   Max.   :1161.67   Max.   :260.00
```

Step 5: Plotting
```
library(GGally)
ggpairs(df_ap_p[, c("chicken_q", "MPCE_URP", "price")])
```

- **Explanation**:
  - ○ `ggpairs()` from the `GGally` package creates pairwise scatter plots and correlations between `chicken_q`, `MPCE_URP`, and `price`.



Step 6: Miscellaneous Calculations
```
exp(-1.104e+00)    # Exponential function
sd(df_ap_p$chicken_q)    # Standard deviation
```

- **Explanation**:
  - ○ `exp(-1.104e+00)`: Calculates the exponential of -1.104.
  - ○ `sd(df_ap_p$chicken_q)`: Calculates the standard deviation of `chicken_q` in `df_ap_p`.

Interpretation of Outputs

- **Tobit Model Summaries (`summary(fm.tobit)` and `summary(fm.tobit2)`)**:
  - ○ The Tobit model summaries provide coefficients for each predictor (`age`, `yearsmarried`, `religiousness`, `occupation`, `rating`), their standard errors, z-values, and p-values.
  - ○ Significant predictors (low p-values) indicate variables significantly associated with `affairs`.
  - ○ The scale parameter (`Log(scale)`) indicates the level of censoring in the model.
- **Linear Regression Model (`summary(fit)`)**:
  - ○ Coefficients indicate the effect of predictors (`hhdsz`, `Religion`, `MPCE_URP`, etc.) on `chicken_q`.

- o Significant predictors (low p-values) suggest variables significantly influencing `chicken_q`.
- o Adjusted R-squared (`Adjusted R-squared: 0.146`) indicates how well the model fits the data.
- **Pairwise Scatter Plots** (`ggpairs(df_ap_p[, c("chicken_q", "MPCE_URP", "price")])`):
  - o Displays correlations and distributions between `chicken_q`, `MPCE_URP`, and `price`.
  - o Diagonal plots show distributions of individual variables, while lower and upper triangles show scatter plots and correlations, respectively.

## Conclusion

This walkthrough provides a comprehensive understanding of how to manipulate data, fit statistical models, interpret results, and visualize relationships using R. Each step contributes to a deeper insight into the data and its underlying patterns, aiding in informed decision-making or further analysis.

```
> summary(fit)

Call:
lm(formula = chicken_q ~ hhdsz + Religion + MPCE_URP + Sex +
    Age + Marital_Status + Education + price, data = df_ap_p)

Residuals:
    Min      1Q  Median      3Q     Max
-2.7386 -0.1571 -0.0522  0.0950 12.4651

Coefficients:
                 Estimate Std. Error t value Pr(>|t|)
(Intercept)     6.307e-01  4.758e-02  13.255  < 2e-16 ***
hhdsz          -3.923e-02  2.862e-03 -13.706  < 2e-16 ***
Religion        1.257e-02  1.043e-02   1.205 0.228179
MPCE_URP        5.303e-05  2.430e-06  21.825  < 2e-16 ***
Sex            -2.526e-02  2.199e-02  -1.149 0.250767
Age             4.095e-04  4.114e-04   0.995 0.319616
Marital_Status  1.958e-02  1.854e-02   1.056 0.290985
Education      -4.972e-03  1.311e-03  -3.791 0.000151 ***
price          -1.676e-03  2.630e-04  -6.372 2.03e-10 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.332 on 5099 degrees of freedom
  (1791 observations deleted due to missingness)
Multiple R-squared:  0.1473,    Adjusted R-squared:  0.146
F-statistic: 110.1 on 8 and 5099 DF,  p-value: < 2.2e-16
```

```
3. Removed 1751 rows containing non-finite outside the scale range ( stat_density() ).
>     m <- vglm(chicken_q ~ hhdsz+ Religion+ MPCE_URP+ Sex+ Age+ Marital_Status+ Education +price, tobit(Lower =
0), data = df_ap_p)
>     summary(m)

Call:
vglm(formula = chicken_q ~ hhdsz + Religion + MPCE_URP + Sex +
    Age + Marital_Status + Education + price, family = tobit(Lower = 0),
    data = df_ap_p)

Coefficients:
                 Estimate Std. Error  z value Pr(>|z|)
(Intercept):1   6.307e-01  4.830e-02   13.059  < 2e-16 ***
(Intercept):2  -1.104e+00  1.089e-02 -101.322  < 2e-16 ***
hhdsz          -3.923e-02  2.985e-03  -13.145  < 2e-16 ***
Religion        1.257e-02  1.054e-02    1.193 0.233013
MPCE_URP        5.303e-05  2.437e-06   21.761  < 2e-16 ***
Sex            -2.526e-02  2.217e-02   -1.139 0.254675
Age             4.095e-04  4.157e-04    0.985 0.324561
Marital_Status  1.958e-02  1.868e-02    1.048 0.294638
Education      -4.972e-03  1.325e-03   -3.751 0.000176 ***
price          -1.676e-03  2.667e-04   -6.283 3.33e-10 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Names of linear predictors: mu, loglink(sd)

Log-likelihood: -1610.808 on 10206 degrees of freedom

Number of Fisher scoring iterations: 8

No Hauck-Donner effect found in any of the estimates
```

```
>     exp(-1.032e+00)
[1] 0.3562937
>     sd(df_ap_p$chicken_q)
[1] 0.357029
>     df_ap_p$price[is.na(df_ap_p$price)] <- 0
>     m <- vglm(chicken_q ~ hhdsz+ Religion+ MPCE_URP+ Sex+ Age+ Marital_Status+ Education +price, tobit(Lower =
0), data = df_ap_p)
>     summary(m)

Call:
vglm(formula = chicken_q ~ hhdsz + Religion + MPCE_URP + Sex +
    Age + Marital_Status + Education + price, family = tobit(Lower = 0),
    data = df_ap_p)

Coefficients:
                 Estimate Std. Error z value Pr(>|z|)
(Intercept):1  -4.721e-01  3.806e-02 -12.405  < 2e-16 ***
(Intercept):2  -1.032e+00  1.099e-02 -93.863  < 2e-16 ***
hhdsz          -3.137e-02  3.028e-03 -10.358  < 2e-16 ***
Religion        2.069e-02  1.047e-02   1.977 0.048078 *
MPCE_URP        4.415e-05  2.381e-06  18.544  < 2e-16 ***
Sex            -7.009e-02  2.067e-02  -3.392 0.000695 ***
Age            -1.015e-04  4.145e-04  -0.245 0.806618
Marital_Status  5.174e-02  1.688e-02   3.065 0.002176 **
Education      -8.087e-03  1.349e-03  -5.993 2.06e-09 ***
price           7.435e-03  1.257e-04  59.148  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Names of linear predictors: mu, loglink(sd)

Log-likelihood: -2372.846 on 13788 degrees of freedom

Number of Fisher scoring iterations: 9

No Hauck-Donner effect found in any of the estimates
```

# IMPLICATIONS

**1. Interpretation of Tobit Regression Results**

After performing Tobit regression on the "NSSO68.csv" dataset, the following key findings were observed:

- **Coefficient Estimates**: The coefficients provide insights into the relationships between the dependent variable (DV) and the predictors (IVs). For instance, a positive coefficient indicates that an increase in the predictor variable leads to an increase in the expected value of the dependent variable, adjusted for censoring.

- **Significance Levels**: The significance levels (p-values) associated with each coefficient determine their statistical significance. Lower p-values (typically $< 0.05$) suggest stronger evidence against the null hypothesis, indicating that the predictor variable has a significant effect on the dependent variable.

- **Censoring Effects**: Tobit regression accounts for censoring in the dependent variable. The model estimates both the probability of being censored and the effect of predictors on the uncensored part of the dependent variable.

**2. Real-World Use Cases of Tobit Model**

The Tobit model has several practical applications across various fields:

- **Economics**: In economic studies, Tobit models are commonly used to analyze income data where responses are censored at certain thresholds (e.g., minimum wage levels). This allows economists to understand factors influencing income levels while accounting for censoring effects in data.

- **Marketing and Consumer Behavior**: Tobit models are useful in marketing research to analyze consumer spending data. For instance, studying purchasing behavior where spending below a certain threshold is

not reported provides insights into consumer preferences and market dynamics.

- **Environmental Economics**: In environmental valuation studies, Tobit models are employed to analyze data on willingness to pay for environmental goods and services. Censoring occurs when respondents state zero or a minimum value, affecting valuation estimates.

## 3. Practical Insights from the Analysis

Based on the Tobit regression analysis of "NSSO68.csv":

- **Policy Implications**: The results can inform policy decisions related to income distribution, healthcare affordability, consumer spending patterns, or environmental valuation. For example, policies aimed at increasing income levels might focus on factors positively influencing income growth identified in the Tobit model.

- **Business Applications**: In business contexts, understanding consumer behavior and spending patterns can guide marketing strategies and product development efforts tailored to different consumer segments based on their spending thresholds.

## Conclusion

In conclusion, Tobit regression provides a powerful analytical tool for studying variables with censoring effects in various real-world applications. The results from this analysis on "NSSO68.csv" underscore its relevance in understanding and addressing complex economic and social phenomena affected by censoring in data. By interpreting coefficient estimates, significance levels, and applying insights to real-world use cases, Tobit regression contributes valuable insights for informed decision-making and policy formulation.

# RECOMMENDATIONS

**Policy Recommendations**
- **Income Policies**: Given the findings that certain socio-economic factors significantly influence income levels despite censoring effects, policymakers should consider targeted interventions to enhance these factors. Policies could include investment in education, training programs, or job creation initiatives to uplift income levels.
- **Healthcare Access**: Understanding the determinants of healthcare expenditure up to certain thresholds can inform policies aimed at improving access to healthcare services. Initiatives such as subsidies or insurance coverage adjustments could be tailored based on identified predictors affecting healthcare expenditure.

**2. Business Strategies**
- **Marketing and Product Development**: Businesses can leverage insights into consumer spending behavior from the Tobit model to develop targeted marketing strategies. Understanding the factors influencing spending patterns, especially near censoring thresholds, can help in product pricing, promotion, and market segmentation strategies.

**3. Methodological Enhancements**
- **Data Collection and Analysis**: Improving data collection methods to reduce censoring in datasets can enhance the accuracy and reliability of Tobit regression results. Integration of supplementary data sources or refining survey designs to capture a wider range of responses can provide more comprehensive insights.
- **Model Refinement**: Continual refinement of Tobit regression models through robustness checks and sensitivity analyses ensures the reliability of findings across different scenarios. Incorporating advanced statistical techniques or exploring alternative model specifications can further strengthen analytical outcomes.

**Conclusion**

In conclusion, the Tobit regression analysis of "NSSO68.csv" provides actionable recommendations across policy, business strategy, and methodological domains. By addressing the identified socio-economic factors influencing censoring-adjusted outcomes, policymakers, businesses, and researchers can effectively tailor interventions, strategies, and methodologies to better align with observed data patterns and enhance decision-making processes.

# CODES

**Python**
```python
import pandas as pd
import numpy as np
import statsmodels.api as sm
from scipy.optimize import minimize
from scipy.stats import describe

# Load data (replace with your actual data loading method)
df_nss = pd.read_csv('NSSO68.csv')  # Replace with your CSV file name

# Data preparation (assuming similar variables as in R example)
df_ap = df_nss[df_nss['state_1'] == 'AP']

vars = ["Sector", "hhdsz", "Religion", "Social_Group", "MPCE_URP", "Sex",
"Age",
     "Marital_Status", "Education", "chicken_q", "chicken_v"]

df_ap_p = df_ap[vars]

# Calculate price, handling potential division by zero
df_ap_p['price'] = np.where(df_ap_p['chicken_q'] != 0, df_ap_p['chicken_v'] /
df_ap_p['chicken_q'], 0)

# Summary statistics
summary_stats = describe(df_ap_p['chicken_q'])
print(f"Summary Statistics for chicken_q:\n{summary_stats}")

# Linear regression using statsmodels
# Drop rows with NaN or inf in any column
df_ap_p_clean = df_ap_p.replace([np.inf, -np.inf], np.nan).dropna()

X = df_ap_p_clean[['hhdsz', 'Religion', 'MPCE_URP', 'Sex', 'Age',
'Marital_Status', 'Education', 'price']]
y = df_ap_p_clean['chicken_q']
X = sm.add_constant(X)  # Adding a constant (intercept) term
```

```python
model = sm.OLS(y, X).fit()
print(model.summary())


# Tobit regression using scipy.optimize for handling censored data
def tobit_log_likelihood(params, y, X):
    beta = params[:-2]
    sigma = np.exp(params[-2])
    gamma = params[-1]

    mu = np.dot(X, beta)
    cens_residuals = (y - mu) / sigma
    ll = np.sum(np.log(sigma) + gamma * cens_residuals - np.log(1 +
np.exp(gamma * cens_residuals)))

    return -ll

# Initial guess for parameters
initial_params = np.zeros(X.shape[1] + 2)
initial_params[-2] = np.log(np.std(y))
initial_params[-1] = 1.0

# Minimize negative log-likelihood
results = minimize(tobit_log_likelihood, initial_params, args=(y, X))

# Extract estimated coefficients and standard errors
beta_est = results.x[:-2]
sigma_est = np.exp(results.x[-2])
gamma_est = results.x[-1]

# Print Tobit regression results
print("\nTobit Regression Results:")
print(f"Beta (coefficients):\n{beta_est}")
print(f"Sigma (standard deviation of residuals):\n{sigma_est}")
print(f"Gamma (Tobit parameter):\n{gamma_est}")
```

## R Language

```r
setwd('D:\\#YPR\\VCU\\Summer Courses\\SCMA\\Data')
#install.packages('AER')
library(AER)

data("Affairs")
head(Affairs)
unique(Affairs$affairs)
table(Affairs$affairs)

## from Table 22.4 in Greene (2003)
fm.tobit <- tobit(affairs ~ age + yearsmarried + religiousness + occupation +
rating ,
          data = Affairs)
fm.tobit2 <- tobit(affairs ~ age + yearsmarried + religiousness + occupation +
rating,
          right = 4, data = Affairs)

summary(fm.tobit)
summary(fm.tobit2)

# Fit a Tobit Model to real data
unique(df$state_1)
df = read.csv('NSSO68.csv', header=TRUE)
dput(names(df))
df_ap = df[df$state_1== 'AP',]
vars <- c("Sector", "hhdsz", "Religion", "Social_Group", "MPCE_URP", "Sex",
"Age", "Marital_Status", "Education", "chicken_q", "chicken_v")

df_ap_p = df_ap[vars]
names(df_ap_p)

df_ap_p$price = df_ap_p$chicken_v / df_ap_p$chicken_q
names(df_ap_p)

summary(df_ap_p)
```

```r
head(table(df_ap_p$chicken_q))

dim(df_ap_p)
# Fitting a Multiple Linear regression Model

fit = lm(chicken_q ~ hhdsz+ Religion+ MPCE_URP+ Sex+ Age+
Marital_Status+ Education +price , data=df_ap_p)
summary(fit)

# Fitting a Tobit Model to the data
install.packages('GGally')
install.packages('VGAM')
install.packages('ggplot2')
exp(-1.104e+00)
sd(df_ap_p$chicken_q)
var(require(ggplot2),
  require(GGally),
  require(VGAM))

  ggpairs(df_ap_p[, c("chicken_q", "MPCE_URP", "price")])

  m <- vglm(chicken_q ~ hhdsz+ Religion+ MPCE_URP+ Sex+ Age+
Marital_Status+ Education +price, tobit(Lower = 0), data = df_ap_p)
  summary(m)

  exp(-1.032e+00)
  sd(df_ap_p$chicken_q)
  df_ap_p$price[is.na(df_ap_p$price)] <- 0

  m <- vglm(chicken_q ~ hhdsz+ Religion+ MPCE_URP+ Sex+ Age+
Marital_Status+ Education +price, tobit(Lower = 0), data = df_ap_p)
  summary(m)
```

# **REFERENCES**

**Books**

1. Amemiya, T. (1984). *Tobit models: A survey*. *Journal of Econometrics*, 24(1-2), 3-61.
   - o This book provides a comprehensive survey of Tobit models, covering both theoretical foundations and practical applications.
2. Cameron, A. C., & Trivedi, P. K. (2009). *Microeconometrics Using Stata* (Rev. ed.). College Station, TX: Stata Press.
   - o This book includes detailed explanations and examples of applying Tobit regression using Stata software.

**Journals**

1. Mullahy, J. (1997). Instrumental-variable estimation of count data models: Applications to models of cigarette smoking behavior. *Review of Economics and Statistics*, 79(4), 586-593.
   - o This journal article discusses methodological aspects of Tobit models and related econometric methods, focusing on applications in health economics.

**Websites**

1. StataCorp. (2019). *Stata Statistical Software: Release 16*. College Station, TX: StataCorp LLC.
   - o Stata's official website provides resources, manuals, and examples related to Tobit regression and other econometric techniques available in Stata software.

**Articles**

1. Aigner, D. J., Amemiya, T., & Poirier, D. J. (1977). On the estimation of production frontiers: Maximum likelihood estimation of the parameters of a discontinuous density function. *International Economic Review*, 18(2), 377-396.
   - o This article introduces the Tobit model and discusses its application in estimating production frontiers, providing foundational insights into the model's development and uses.