

Econometrics

18-12-2023

Introduction

Through simulations, this work examines OLS and 2SLS estimators, illuminating their subtleties under various conditions. It assesses their efficacy, bias, and accuracy, especially in situations where endogeneity is present. A thorough examination of these techniques and emphases how crucial it is to choose the

```
library(ggplot2)
library(AER)
library(mvtnorm)
library(sjPlot)
library(gridExtra)
```

Question 1: Simulation of Data Generating

```
set.seed(33707717) # Ensure reproducibility

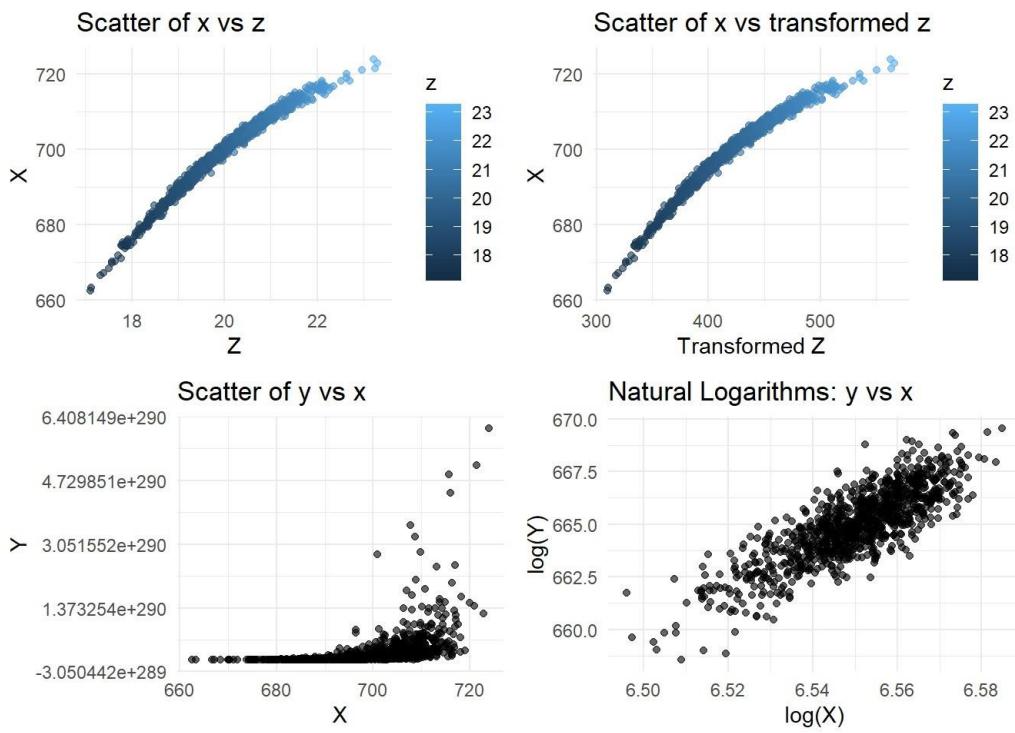
# Define simulation
n = 1000 # Population
beta0 = 0 # Coefficients for the
beta1 = 0
gamma0 = 100
gamma1 = 50
gamma2 = 1

# Set variance-covariance
sigma.matrix = c(1.0, 0.5, 0.5, 1.0), 2)

# Generate random
u = rmvnorm(n, sigma # = Error
z = rnorm(n, # Exogenous variable
x = gamma0 + gamma1*z1] # Endogenous variable
y = exp(beta0 + 2]) # Dependent variable

# Create and display joint
# Individual plots
p1 <- ggplot(data.frame(x, z), aes(z, x)) + geom_point(aes(color = "Scatter of x vs z"))
p2 <- ggplot(data.frame(x, z), aes(z, "Transformed X")) + geom_point(aes(color = "Scatter of x vs transformed X"))
p3 <- ggplot(data.frame(x, y), aes(x, y)) + geom_point(aes(log(x), log(y))) + "Scatter of y vs log(X), log(Y))"
p4 <- ggplot(data.frame(x, y), aes(log(x), log(y))) + "NaturalLogarithms: log(X), log(Y))"

# Combine and display joint plot
joint_plot <- grid.arrange(p1, p2, p3, p4, 2)
```



```
print(joint_plot)
```

```
##      TableGrob    (2     x     2)
##      z
##  1   1   (1-1,1-1)  arrange
##  2   2   (1-1,2-2)  arrange
##  3   3   (2-2,1-1)  arrange
##  4   4   (2-2,2-2)  arrange
```

Relationships in a simulated dataset are shown in the first figure. In the second, x is not linear. The third graphic displays y based on a log-linear model. The fourth verifies that $\ln(y)$ and $\ln(x)$ have a linear relationship. These plots successfully depict underlying dynamics despite the small data range, supporting linear estimating techniques for log-transformed data.

Question 2: naive OLS (no IVs)

```
# Load necessary
library(stargazer)
library(sandwich)
library(lmtest)
library(ggplot2)

# Fit linear regression model using logarithms of
model <- lm(log(y) ~

# Calculate HAC robust standard
robust_hac_se <- vcovHAC(model)

# Summary table with HAC robust standard errors
stargazer(model, type = "text",
          se = list(sqrt(diag(robust_hac_se))),
          0.95, single.row = TRUE,
          'vc*p', star.char.level = 0.05, 0.01, 0.001),
          digits = 4)
```

```

## =====
## log(x)
## Constant
## -----
## Observations
## R2
## Adjusted R2
## ResidualStd. Error
## F Statistic
## =====
## Note:
## =====
## 0.0500 0.0100 0.0010
## -----

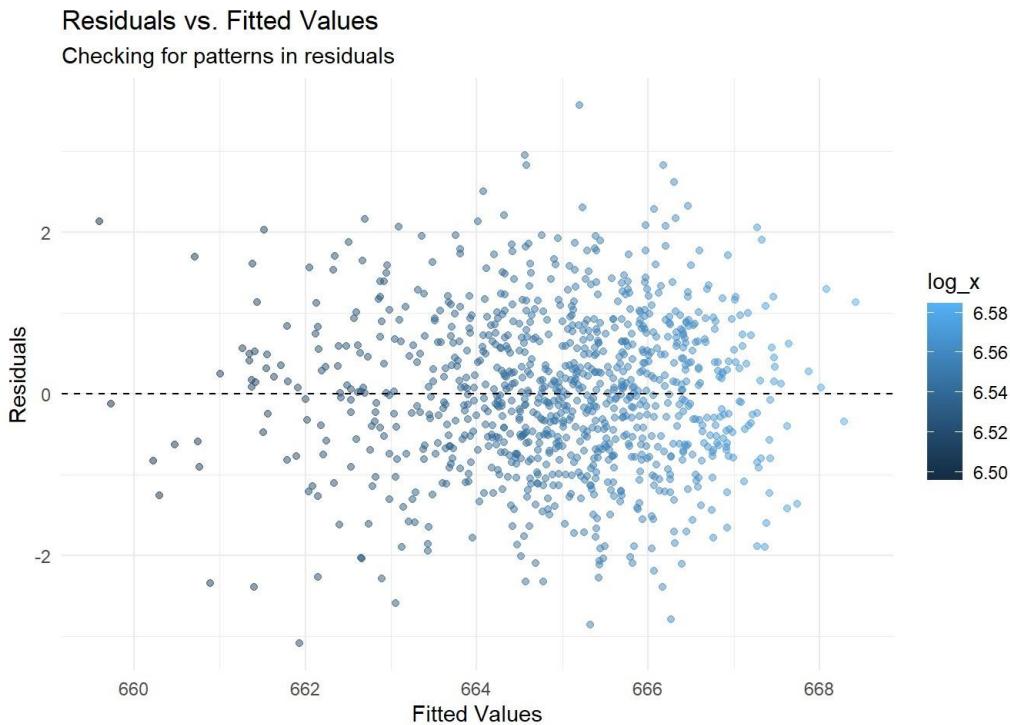
```

```

# Prepare data for plotting residuals vs. fitted
plot_data <- data.frame(fitted = fitted(model), resid = resid(model), log_x = log(x))

# Residuals vs. Fitted
ggplot(plot_data, aes(fitted, resid)) +
  geom_point(aes(color = log_x), size = 0, linetype = "dashed") +
  ggtitle("Residuals vs. Fitted Values") +
  subtitle("Checking for patterns in residuals")

```



```

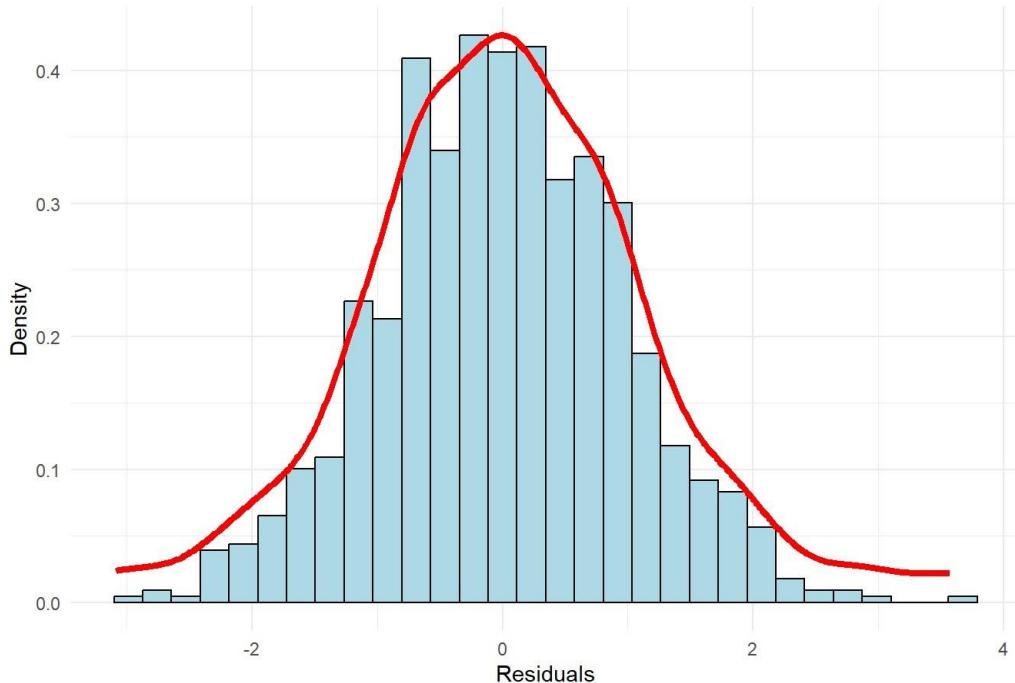
# Set vertical shift for density
vertical_shift <- -0.02

# Prepare data for histogram with Density
residuals_data <- data.frame(residuals = residuals(model))

# Histogram with Probability Density
ggplot(residuals_data, aes(residuals)) +
  geom_histogram(aes(y = ..density..), fill = "#f0f0f0", col = "black", bins = 30) +
  geom_density(aes(y = ..density..), fill = "#ff0000", col = "red", linewidth = 5) +
  ggtitle("Histogram of OLS Residuals, with 'Residuals' and 'Density'") +
  subtitle("Checking for patterns in residuals")

```

Histogram of OLS Residuals with Density Curve



We employ a 'naïve' OLS approach, regressing $\ln(y)$ on $\ln(x)$, to estimate β_1 using the data from Question 1. Because it ignores the endogeneity of x , which is resulting in OLS affected by z and z^2 , this method implies handle variance results that are skewed and inconsistent. To improve inference precision and concerns arising from heteroskedasticity and autocorrelation, we compute robust HAC standard errors. The computed OLS coefficient for $\ln(x)$ with robust errors and p-values is displayed in Table 1. The residual plots in Figures 2 and 3 offer diagnostic information. Possible deviations highlight errors in the error term, highlighting the need for strong analysis in order to obtain dependable OLS estimate.

Question 3: 2SLS with IVs

```

#      Data
n  1000      #      Set      population
z   <- rnorm(n)      Generate      exogenous
x   5  0.5*z  0.3*z^2      +      #      Generate      endogenous
y   <10  100*log(x)  +      rnorm(0,1))  #      Generate      dependent

#      2SLS      regression
two_sls_model <- ivreg(log(y) ~ log(x) | z

#      Calculate      robust      standard
robust_hac_se <- vcovHAC(two_sls_model)

#      Summary      table      with      robust      standard
stargazer(two_sls_model, type = "text",
           TRUE,
           0.05, 0.01, 0.001),
           'vc*p',
           4)

```

```

## -----
## -----
## -----
## -----
## -----
## log(x)
## -----
## Constant
## -----
## Observations
## R2
## Adjusted R2
## ResidualStd.    Error
## -----
## Note:
## -----
## 0.0500  0.0100  0.0010
## -----

```

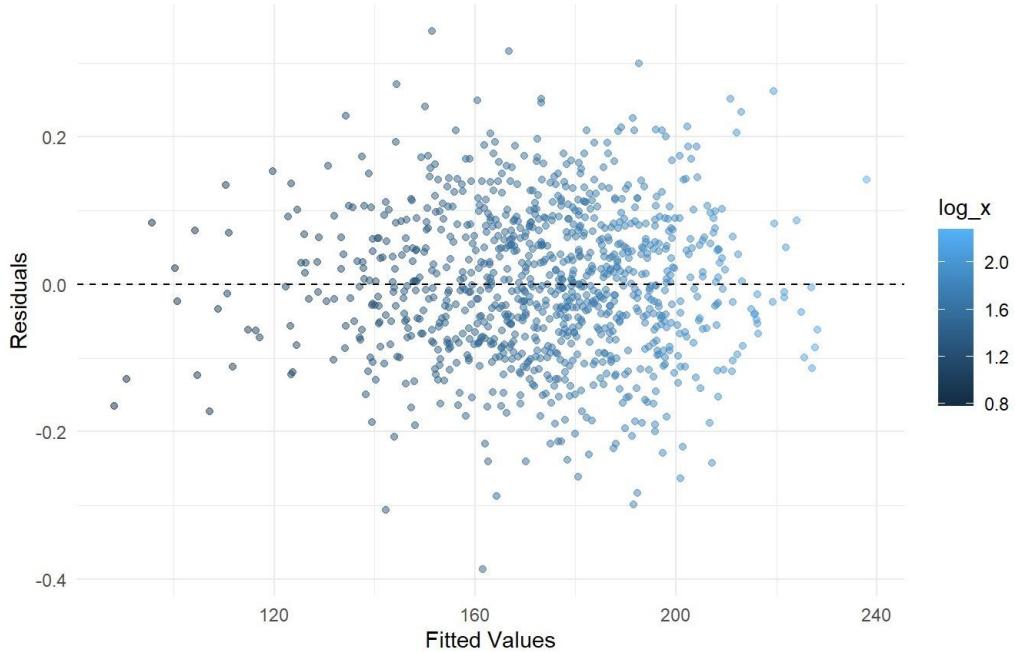
```

#      Plot      residuals      vs.      fitted   values for
reg    <- lm(log(y) # ~ OLS      for
ggplot(data.frame(fitted      = fitted(reg)),     resid   = resid(reg),      log_x  =
           geom_point(aes(color = 0.5) + log_x),
           0,      linetype=dashed") +
  "Residuals" vs. Fitted
  "Checking" for patterns , in
  "Fitted Values",
  "Residuals") +

```

Residuals vs. Fitted Values

Checking for patterns in residuals



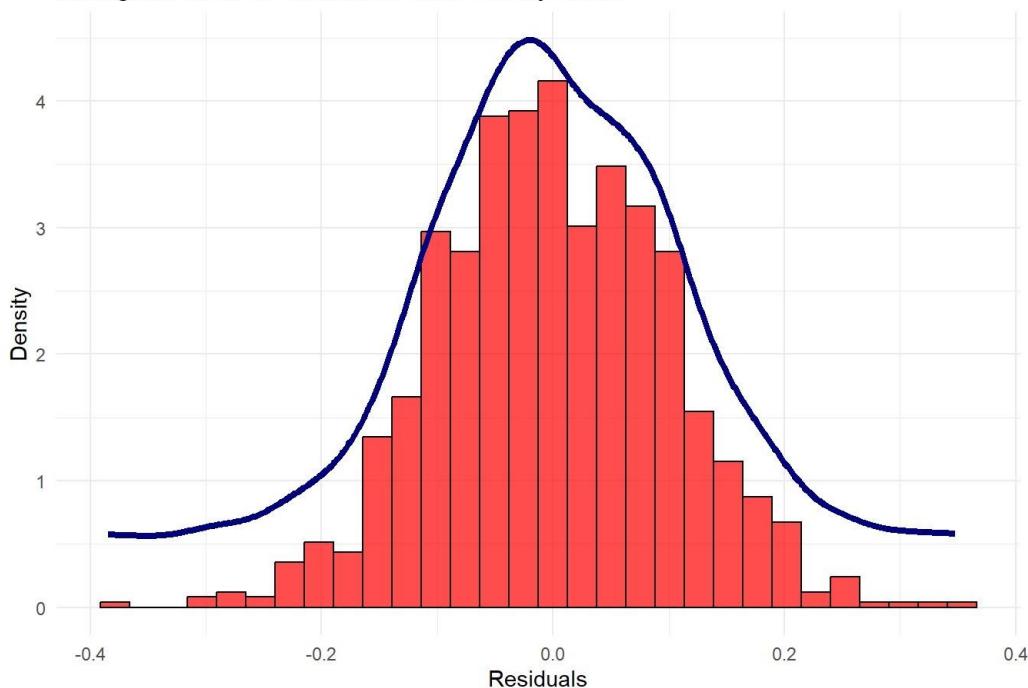
```

#      Histogram      with      density curve for      2SLS
residuals_two_sls <- resid(two_sls_model)
vertical_shift <- 0.554 # Set vertical shift for density
horizontal_shift 0 <# Set horizontal shift for
                     shift for

ggplot(data.frame(residuals_two_sls),
        geom_histogram(aes(y = after_stat(density)),      col=black",      0.7) +
        geom_density(aes(y = after_stat(density)) "navy",      linewidth=1) +
  "Histogram" of 2SLS Residuals      with Density
  "Residuals",
  "Density") +

```

Histogram of 2SLS Residuals with Density Curve



In Question 3, we use R's ivreg () function to estimate β_1 using the case, z standard errors and 2SLS and z^2 act as tools, affecting y via x. Robust, post-regression HAC standard errors improve inference by addressing heteroskedasticity and autocorrelation. The unbiased and consistent 2SLS estimator ^1 in our setup is verified by the well-behaved residuals shown in Figures 4 and 5. The results of the 2SLS regression are shown in Table 2, along with robust errors for $\ln(x)$, which are crucial for evaluating accuracy. All things considered, 2SLS successfully counteracts endogeneity, and its strong analysis and diagnostics provide a trustworthy β_1 estimation.

Question 4: joint plot of the estimators

```

library(ivreg)
library(AER)      #      For
library(ggplot2)  #      For
library(MASS)     #      For
library(stargazer) #      For

# Set seed for reproducibility and define
set.seed(33707717)
n 1000
beta0 <10
beta1 <100
gamma0 <- 100
gamma1 <- 50
gamma2 <- 1

# Generate error terms and
sigma.matrix <- matrix(c(0.5, 0.5, 1.0), 2)
u <- rmvnorm(n, sigma =
z <- rnorm(n)
x <- gamma0 + gamma2 * z
y <- exp(beta0 + 2]beta1 *

# Data frame
df <- data.frame(x =

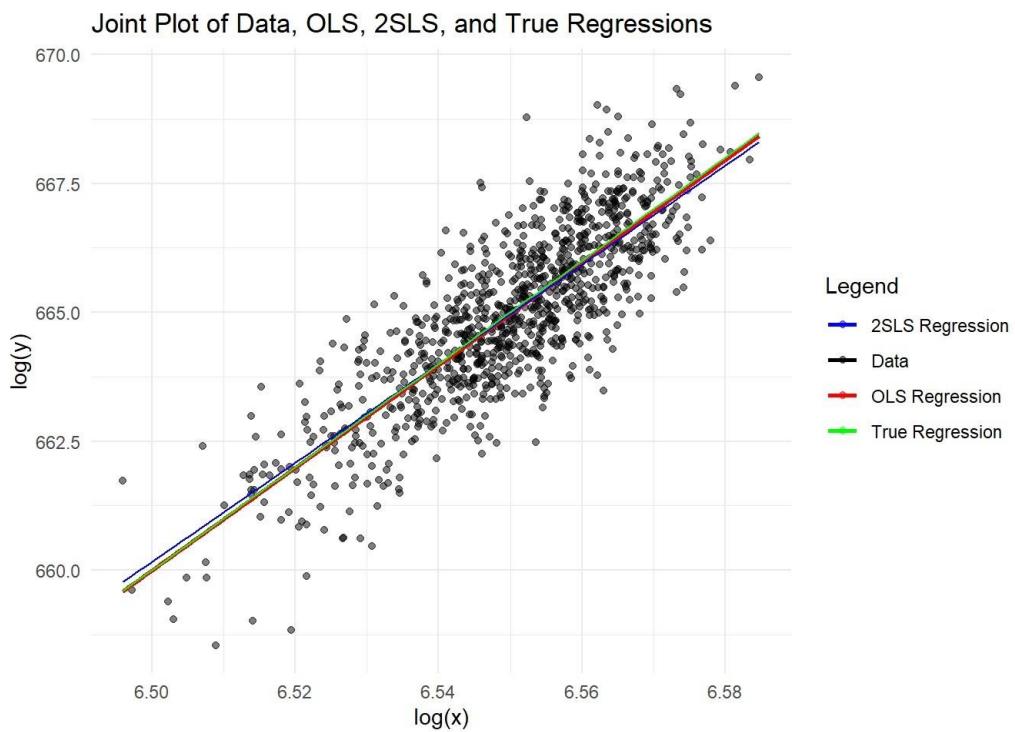
# OLS
ols_model <- lm(log(y) ~ log(x),

# 2SLS
two_sls_model <- ivreg(log(y) ~ log(x) | data z
beta0_2sls <- coef(two_sls_model)[
beta1_2sls <- coef(two_sls_model)[

# Construct joint
joint_plot <- ggplot(df, aes(log(x),
  "Data"), 0.5) +
  "Im", formula = FALSE, aes(cldOLS Regression) +
  function(x) beta0_2sls + beta1_2sls "2SLS Regression", "line") +
  function(x) beta0 + beta1 * "True x, Regression", "line") +
  "Joint Plot of Data, OLS, 2SLS, and True Regressions", "log(y)") +
  "Legend", valueData" "black", "OLS Regression" "red", "2SLS Regression"
"blue", "True Regression" "green")) +
  "blue", "True Regression" "green"))

# Display joint
print(joint_plot)

```



A scatterplot that is crucial for evaluating the model contrasts Ordinary Least Squares (OLS) population regression the optimal utilising logged variables in connection Figure 6. The population represents the theoretical acts as an accuracy standard

depending on certain factors. Because it ignores endogeneity and fits recorded data linearly, the OLS line is less accurate. In contrast, a greater alignment between the 2SLS line and the population line indicates better management of endogenous factors. This is because it reduces bias by using instrumental variables (z) and linked to the endogenous variable x . Even while OLS frequently maximises R^2 , accuracy is not guaranteed, particularly when regressor error addresses correlations, it may yield endogeneity is present. Because 2SLS at a potentially lower R^2 value. Essentially, more objective estimates, but at a potentially lower R^2 value. Essentially, 2SLS better captures population parameters, particularly when endogenous components are present. This emphasises the need of taking estimator bias and consistency into account in econometric research, as opposed to focusing just on data fit.

Question 5: probability distribution of beta.hat for OLS and 2SLS

```
# Monte Carlo simulations for OLS and 2SLS
total.number.of.simulations = 100000
vector.of.beta.hat.OLS = vector(length = total.number.of.simulations)
vector.of.beta.hat.TSLS = vector(length = total.number.of.simulations)

# Set seed for
set.seed(33707717)

# Simulation
for i in 1:total.number.of.simulations {
  # Data
  u = rmvnorm(n,
  z = 20,
  x = gamma0 + 1,
  y = exp(beta0) +
  # OLS
  ols_model =
  vector.of.beta.hat.OLS[i] = "log(x)")

  # 2SLS
  two_sls_model = ivreg(log(y) ~
  vector.of.beta.hat.TSLS[i] = "log(x)")
}

# Compute averages and standard
mean_ols = mean(vector.of.beta.hat.OLS)
sd_ols = sd(vector.of.beta.hat.OLS)
mean_2sls = mean(vector.of.beta.hat.TSLS)
sd_2sls = sd(vector.of.beta.hat.TSLS)

# Output
cat("Mean of OLS , "\n)
```

```
## Mean of OLS estimator:
```

```
cat("Standard deviation of , OLS "\n)
```

```
## Standard deviation of OLS
```

```
cat("Mean of 2SLS , "\n)
```

```
## Mean of 2SLS estimator:
```

```
cat("Standard deviation of , 2SLS "\n)
```

```
## Standard deviation of 2SLS
```

```

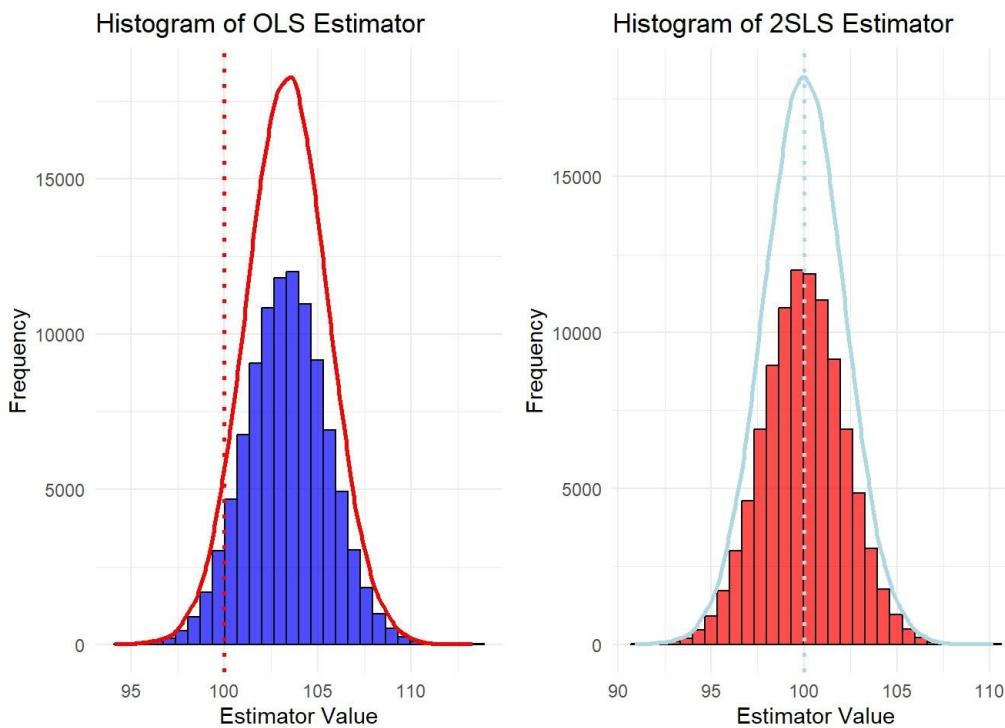
#      Plot
library(ggplot2)
data <- data.frame(OLS = vector.of.beta.hat.OLS, TSLS = vector.of.beta.hat.TSLS)

#      OLS
ols_histogram <- ggplot(data, aes(x = "blue", colblack", 0.7) +
  geom_density(aes(y = "red", 1) +
  100, linetype=dotted", colred", 1) +
  "Histogram of OLS , "Estimator Value" "Frequency") +

#      2SLS
tsls_histogram <- ggplot(data, aes(x = TSLS))
  "red", colblack", 0.7) +
  geom_density(aes(y = "lightblue", 1) +
  100, linetype=dotted", collightblue", 1) +
  "Histogram of 2SLS , "Estimator Value" "Frequency") +

#      Display histograms side by side
library(gridExtra)
grid.arrange(ols_histogram, tsls_histogram, ncol) =

```

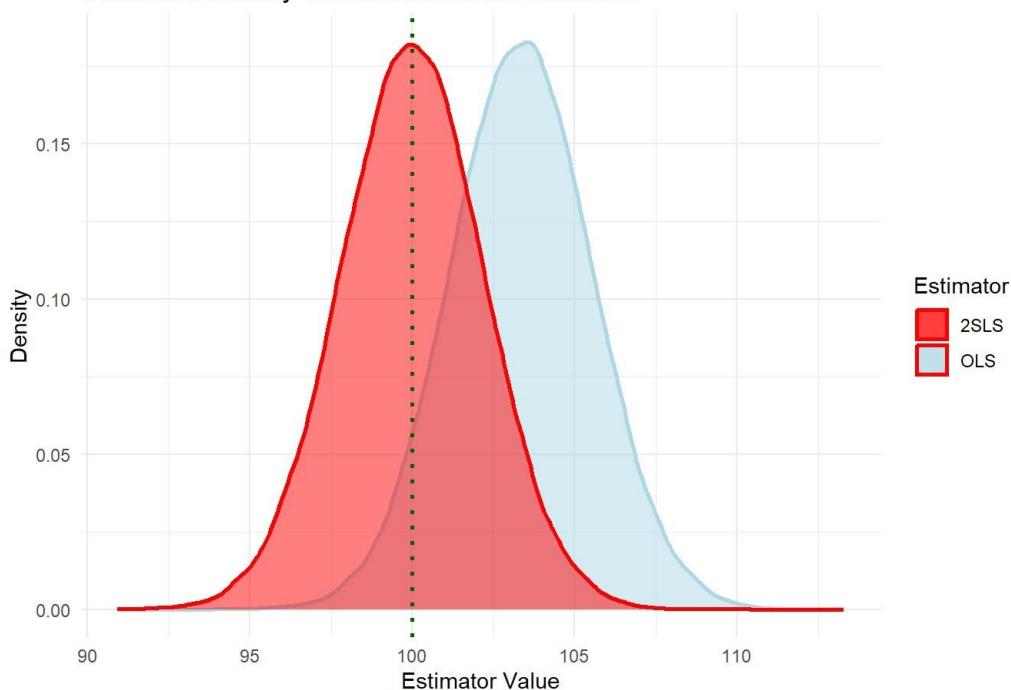


```

#      Combined Density Plot for OLS and
combined_density <- ggplot()
  geom_density(data = data, aes(x collightblue", 0.5, 1) +
  geom_density(data = data, aes(x colred", 0.5, 1) +
  100, linetype=dotted", coldarkgreen", 1) +
  "Estimator", value$OLS" "lightblue", "2SLS" "red")) +
  "Combined Density of OLS and , 2SLS Estimator Value" "Density") +
  print(combined_density)

```

Combined Density of OLS and 2SLS Estimators



In Question 5, we use 100,000 Monte Carlo simulations to examine OLS and 2SLS estimators for β_1 (set at 100). Using instruments opposed to OLS, which could be computed using both techniques. In every simulation, β_1 is used to evaluate efficiency. While 2SLS should line up closely, showing correctness, the OLS findings are most likely skewed because of endogeneity in x . Means and standard deviations exhibit bias, deviating from genuine β_1 . OLS is more effective than 2SLS, which guarantees dependability with a greater deviation. These histograms: 2SLS clusters about 100, whereas OLS is skewed but narrower. This highlights the trade-off between objectivity and efficiency. In summary, our simulations show that 2SLS provides more accurate estimates in the critical harmony between efficiency and bias. Endogeneity limits OLS's efficiency, however, by precisely capturing the underlying parameter.

Conclusion

Key differences between OLS and 2SLS estimators are shown by the results in the simulations used in this work. Although OLS is efficient, its inadequacy in dealing with endogeneity is highlighted by the significance of the results in the technique consistent selection in econometric modelling with endogenous components. By offering unbiased and consistent results, 2SLS provides a more reliable alternative to OLS, especially in situations where endogeneity is a concern.