

The Nature of Heteroskedasticity

The Nature of Heteroskedasticity, Part I

- Consider our basic linear function:

$$FOOD_EXP_i = \beta_1 + \beta_2 INCOME_i + e_i.$$

- The random error e_i represents the collection of all the factors other than income that affect household expenditure on food.
- The assumption of **strict exogeneity** says that, when using information on household income, our best prediction of the random error is zero.
- If sample values are randomly selected, then the technical expression for this assumption is that given income the conditional expected value of the random error e_i is zero.

The Nature of Heteroskedasticity, Part II

- If the assumption of strict exogeneity holds then the regression function is

$$E(FOOD_EXP_i | INCOME_i) = \beta_1 + \beta_2 INCOME_i$$

- Holding income constant, and given our model, what is the source of the variation in household food expenditures? It must be from the random error.
- Recall that the random error in the regression is the difference between any observation on the outcome variable and its conditional expectation, that is:
 - (8.2) $e_i = FOOD_EXP_i - E(FOOD_EXP_i | INCOME_i)$

Heteroskedasticity in the Multiple Regression Model (cont.)

- Heteroskedasticity often arises when using cross-sectional data.
- The term “cross-sectional data” refers to having data on a number of economic units, such as firms or households, at a *given point in time*.
- Heteroskedasticity is not a property that is necessarily restricted to cross-sectional data.
- With time-series data, it is possible that the conditional error variance will change.

Heteroskedasticity

Consequences for the OLS Estimator

There are two implications of heteroskedasticity.

1. The least squares estimator is still a linear and unbiased estimator, but it is no longer best. There is another estimator with a smaller variance.
2. The standard errors usually computed for the least squares estimator are incorrect. Confidence intervals and hypothesis tests that use these standard errors may be misleading.

FORDHAM
THE JESUIT UNIVERSITY OF NEW YORK

Gabelli School
of Business

Robust Variance Estimators

Heteroskedasticity

Robust Variance Estimator

- Calculation of a correct estimate for the OLS variance is:

$$\widehat{\text{var}}(b_2) = \left[\sum (x_i - \bar{x})^2 \right] \left\{ \sum \left[(x_i - \bar{x})^2 \left(\frac{N}{N-2} \right) \hat{e}_i^2 \right] \right\} \left[\sum (x_i - \bar{x})^2 \right]^{-1}$$

- (8.8)

$$\text{var}(b_2|x) = \left[\sum_{i=1}^N (x_i - \bar{x})^2 \right]^{-1} \sum_{i=1}^N [(x_i - \bar{x})^2 \sigma_i^2] \left[\sum_{i=1}^N (x_i - \bar{x})^2 \right]^{-1}$$

- The **White heteroskedasticity-consistent estimator (HCE)** that is valid in large samples for the simple regression model is:

$$\widehat{\text{var}}(b_2) = \left[\sum (x_i - \bar{x})^2 \right] \left\{ \sum \left[(x_i - \bar{x})^2 \left(\frac{N}{N-2} \right) \hat{e}_i^2 \right] \right\} \left[\sum (x_i - \bar{x})^2 \right]^{-1}$$

Heteroskedasticity

Robust Variance Estimator (cont.)

- Where \hat{e}_i is the least squares residual from the regression model:

$$y_i = \beta_1 + \beta_2 x_i + e_i$$

- This variance estimator is robust because it is valid whether heteroskedasticity is present or not
- If we are not sure whether the random errors are heteroskedastic or homoskedastic, then we can use a robust variance estimator and be confident that our standard errors, t-tests, and interval estimates are valid in large samples
- This does not address the implication that the least square estimator is no longer the best

FORDHAM
THE JESUIT UNIVERSITY OF NEW YORK

Gabelli School
of Business

Generalized Least Squares (GLS)

Known Variance

Generalized Least Squares: Known Form of Variance

- To develop an estimator that is better than the least squares estimator, we need to make a further assumption about how the variances σ_i^2 change with each observation.
- This means making an assumption about the skedastic function $h(x_i)$.
- The further assumption is necessary because the best linear unbiased estimator in the presence of heteroskedasticity, an estimator known as the **generalized least squares (GLS) estimator**, depends on the unknown σ_i^2 .

Transforming the Model: Proportional Heteroskedasticity, Part I

- An estimator known as the **GLS estimator** depends on the unknown σ^2_i .
- To make the GLS estimator operational, some structure is imposed on σ^2_i .
- One possibility is:
 - (8.11) $\text{var}(e_i|x_i) = \sigma_i^2 = \sigma^2 h(x_i) = \sigma^2 x_i, x_i > 0$

Transforming the Model: Proportional Heteroskedasticity, Part II

- We change or transform the model into one with homoskedastic errors.
- Leaving the basic structure of the model intact, we turn the heteroskedastic error model into a homoskedastic error model.
- After the transformation, applying OLS to the transformed model gives a best linear unbiased estimator.

$$(8.12) \quad \frac{y_i}{\sqrt{x_i}} = \beta_1 \left(\frac{1}{\sqrt{x_i}} \right) + \beta_2 \left(\frac{x_i}{\sqrt{x_i}} \right) + \frac{e_i}{\sqrt{x_i}}$$

Transforming the Model: Proportional Heteroskedasticity, Part III

- Define the following transformed variables:

- (8.13) $y_i^* = \frac{y_i}{\sqrt{x_i}}, \quad x_{i1}^* = \frac{1}{\sqrt{x_i}}, \quad x_{i2}^* = \frac{x_i}{\sqrt{x_i}} = \sqrt{x_i}, \quad e_i^* = \frac{e_i}{\sqrt{x_i}}$

- Our model is now:

- (8.14) $y_i^* = \beta_1 x_{i1}^* + \beta_2 x_{i2}^* + e_i^*$

- The beauty of this transformed model is that the new transformed error term e_i^* is homoscedastic

Transforming the Model: Proportional Heteroskedasticity, Part IV

- If X is a random variable and a is a constant, then $\text{var}(aX) = a^2 \text{var}(X)$. Applying that rule here we have:
 - (8.15) $\text{var}(e_i^* | x_i) = \text{var}\left(\frac{e_i}{\sqrt{x_i}} | x_i\right) = \frac{1}{x_i} \text{var}(e_i | x_i) = \frac{1}{x_i} \sigma^2 x_i = \sigma^2$
- The transformed error term will retain the properties of zero mean and zero correlation between different observations.

Transforming the Model: Proportional Heteroskedasticity, Part V

- To obtain the best linear unbiased estimator for a model with heteroskedasticity of the type specified in equation 8.11:
 1. Calculate the transformed variables given in equation 8.13
 2. Use OLS to estimate the transformed model given in equation 8.14, yielding estimates $\widehat{\beta}_1$ and $\widehat{\beta}_2$
- The estimator obtained in this way is called a GLS estimator.

Weighted Least Squares: Proportional Heteroskedasticity

- One way of viewing the GLS estimator is as a **weighted least squares (WLS)** estimator.
- Minimizing the sum of squared transformed errors:

$$\sum_{i=1}^n \frac{(y_i - \beta_1 - \beta_2 x_{i2})^2}{x_i}$$

- The squared errors are weighted by $1/x_i$.

Example Applying GLS/WLS to the Food Expenditure Data

- Applying the generalized (weighted) least squares procedure to our food expenditure problem:

$$\begin{array}{rcccl} \bullet & 8.17 & \widehat{FOOD_EXP}_i & = & 78.68 + 10.45 INCOME_i \\ & & (se) & & (23.79) \quad (1.39) \end{array}$$

- A 95% confidence interval for β_2 is given by:

$$\hat{\beta}_2 \pm t_c se(\hat{\beta}_2) = 10.451 \pm 2.024 \times 1.386 = [7.65, 13.26]$$

FORDHAM
THE JESUIT UNIVERSITY OF NEW YORK

Gabelli School
of Business

Generalized Least Squares (GLS)

Unknown Form of Variance

Generalized Least Squares: Unknown Form of Variance, Part I

- In order to deal with the more general specification, we need a model that is flexible, parsimonious, and for which $\sigma_i^2 > 0$.
- One specification that works well is:

- 8.18
$$\begin{aligned}\sigma_i^2 &= \exp(\alpha_1 + \alpha_2 z_{i2} + \cdots + \alpha_S z_{iS}) \\ &= \exp(\alpha_1) \exp(\alpha_2 z_{i2} + \cdots + \alpha_S z_{iS}) \\ &= \sigma^2 h(z_{i2}, \cdots, z_{iS})\end{aligned}$$

Generalized Least Squares: Unknown Form of Variance, Part II

- Equation 8.18 is called the model of **multiplicative heteroskedasticity**.
- It includes homoskedasticity as a special case; when $\alpha_2 = \dots = \alpha_S = 0$, the error variance is:

$$\sigma_i^2 = \exp(\alpha_1) = \sigma^2$$

- It is called a multiplicative model because:

$$\begin{aligned} & \exp(\alpha_1)\exp(\alpha_2 z_{i2} + \dots + \alpha_S z_{iS}) \\ &= \exp(\alpha_1)\exp(\alpha_2 z_{i2}) \dots \exp(\alpha_S z_{iS}) \end{aligned}$$

Generalized Least Squares: Unknown Form of Variance, Part III

- **Multiplicative heteroskedasticity, special case 1**
 - There are three plausible variance functions; they are special cases of:
 - $\text{var}(e_i|x_i) = \sigma_i^2 = \sigma^2 x_i^{\alpha_2}$
 - Where α_2 is an unknown parameter

Generalized Least Squares: Unknown Form of Variance, Part IV

- **Multiplicative heteroskedasticity, special case 2: grouped heteroskedasticity**
 - Suppose we are considering just two groups.
 - $D_i = 1$ if an observation is in one group and $D_i = 0$ for observations in the other group, Then the variance function is:

$$\begin{aligned} \text{var}(e_i|x_i) &= \exp(\alpha_1 + \alpha_1 D_i) \\ &= \begin{cases} \exp(\alpha_1) = \sigma^2 & D_i = 0 \\ \exp(\alpha_1 + \alpha_2) = \sigma^2 \exp(\alpha_2) & D_i = 1 \end{cases} \end{aligned}$$

Estimating the Multiplicative Model

- How do we proceed with estimation with an assumption like equation 8.18?
- With the model of multiplicative heteroskedasticity, we use several estimation steps
- **Feasible GLS procedure**
 1. Estimate the original model by OLS, saving the OLS residuals $\hat{\epsilon}_i$
 2. Use the least squares residuals and the variables z_{i2}, \dots, z_{iS} to estimate $\alpha_1, \alpha_2, \dots, \alpha_S$

Estimating the Multiplicative Model (cont.)

3. Calculate the estimated skedastic function $\hat{h}(z_{i2}, \dots, z_{iS})$
4. Divide each observation by $\sqrt{\hat{h}(z_{i2}, \dots, z_{iS})}$ and apply OLS to the transformed data, or use WLS regression with weighting factor $1/\hat{h}(z_{i2}, \dots, z_{iS})$
- The resulting estimates are called **feasible generalized least squares (FGLS) estimates** or **estimated generalized least squares (EGLS) estimates**

FORDHAM
THE JESUIT UNIVERSITY OF NEW YORK

Gabelli School
of Business

Detecting Heteroskedasticity

Detecting Heteroskedasticity

- In many applications, there is uncertainty about the presence, or absence, of heteroscedasticity.
- There are two methods we can use to detect heteroskedasticity:
 1. An informal way using residual charts
 2. A formal way using statistical tests

Residual Plots

- If the errors are homoskedastic, there should be no patterns of any sort in the residuals.
- If the errors are heteroskedastic, they may tend to exhibit greater variation in some systematic way.
- We discovered that the absolute values of the residuals do indeed tend to increase as income increases.

Residual Plots (cont.)

- This method of investigating heteroskedasticity can be followed for any simple regression.
- In a regression with more than one explanatory variable, we can plot the least squares residuals against each explanatory variable, or against \hat{y}_i , to see if they vary in a systematic way relative to the specified variable.

The Goldfeld–Quandt Test

- The **Goldfeld–Quandt** test uses the estimated error variances from separate sub-sample regressions as a basis for the test.
- Let the first sub-sample contain N_1 observations.
- Let the regression model in this partition have K_1 parameters.
- The test statistic is 8.22. $GQ = \frac{\hat{\sigma}_1^2}{\hat{\sigma}_2^2} \sim F_{(N_M - K_1, N_2 - K_2)}$

Example The Goldfeld–Quandt Test in the Food Expenditure Model

- With the observations ordered according to income x_i , and the sample split into two equal groups of 20 observations each, it yields:

$$\hat{\sigma}_1^2 = 3574.8 \quad \hat{\sigma}_2^2 = 12921.9$$

- Calculate:

$$F = \frac{\hat{\sigma}_2^2}{\hat{\sigma}_1^2} = \frac{12921.9}{3574.8} = 3.61$$

Example The Goldfeld–Quandt Test in the Food Expenditure Model (cont.)

- Believing that the variances could increase, but not decrease, with income, we use a one-tail test with 5% critical value $F(0.95, 18, 18) = 2.22$.
- Because $3.61 > 2.22$, a null hypothesis of homoskedasticity is rejected in favor of the alternative that the variance increases with income.

A General Test for Conditional Heteroskedasticity, Part I

- In this section, we consider a test for conditional heteroskedasticity that is related to some “explanatory” variables.
- Under assumptions MR1–MR5, the OLS estimator is the best linear unbiased estimator of the parameters $\beta_1, \beta_2, \dots, \beta_k$.
- When conditional heteroskedasticity is a possibility, we suppose that the variance of the random error, e_i , depends on a set of explanatory variables $z_{i2}, z_{i3}, \dots, z_{ik}$ that may include some or all of the explanatory variables x_{i2}, \dots, x_{ik} .

A General Test for Conditional Heteroskedasticity, Part II

- Assume a general expression for the conditional variance
$$\text{var}(e_i|z_i) = \sigma_i^2 = E(e_i^2|z_i) = h(\alpha_1 + \alpha_2 z_{i2} + \cdots + \alpha_s z_{is})$$
- Where $h(\cdot)$ is some smooth function and $\alpha_2, \alpha_3, \dots, \alpha_s$ are nuisance parameters
- We will test for any relationship between the variance of the error term and any function of the selected variables

A General Test for Conditional Heteroskedasticity, Part III

- The null and alternative hypotheses for a test for heteroskedasticity based on the variance function are:
 - Homoskedasticity $\leftrightarrow H_0: \alpha_2 = \alpha_3 = \dots = \alpha_S = 0$
 - Heteroskedasticity $\leftrightarrow H_1: \text{not all the } \alpha_s \text{ in } H_0 \text{ are zero}$
 - If the random errors are homoskedastic, then the sample size multiplied by R^2 , $N \times R^2$ or simply NR^2 , has a chi-square (χ^2) distribution with $S - 1$ degrees of freedom.

A General Test for Conditional Heteroskedasticity, Part IV

- The test statistic is:
 - 8.30 $NR^2 \underset{\sim}{\chi}_{(S-1)}^2$ if the null hypothesis of homoskedasticity is true
- There are several important features of this test.
 1. It is a large sample test. The result in equation 8.30 holds approximately in large samples.
 2. You will often see the test referred to as a Lagrange multiplier test (LM test) or a Breusch–Pagan test for heteroskedasticity.

A General Test for Conditional Heteroskedasticity, Part V

3. One of the amazing features of the Breusch–Pagan/LM test is that the value of the statistic computed from the linear function is valid for testing an alternative hypothesis of heteroskedasticity where the variance function can be of any form given by equation 8.24.
4. The Breusch–Pagan test is for conditional heteroskedasticity. Unconditional heteroskedasticity exists when the error term variance is completely random.

The White Test, Part I

- The variance tests used so far presuppose that we have knowledge of what variables will appear in the variance function if the alternative hypothesis of heteroskedasticity is true.
- We may wish to test for heteroskedasticity without precise knowledge of the relevant variables.

The White Test, Part II

- Suppose that:

$$y_i = \beta_1 + \beta_2 x_{i2} + \beta_3 x_{i3} + e_i$$

- The White test without cross-product terms (interactions) specifies:

$$Z_2 = x_2 \quad Z_3 = x_3 \quad Z_4 = x_2^2 \quad Z_5 = x_3^2$$

- Including interactions adds one further variable:

$$Z_6 = x_2 x_3$$

The White Test, Part III

- The White test is performed using the NR^2 test defined in equation 8.29 or an F-test.
- One difficulty with the White test is that it can detect problems other than heteroskedasticity.
- Thus, while it is a useful diagnostic, be careful about interpreting the result of a significant White test.

Model Specification and Heteroskedasticity, Part I

- As hinted at the end of the previous section, heteroskedasticity can be present because of a model specification error.
- If data partitions are not recognized, or important variables omitted, or an incorrect functional form selected, then heteroskedasticity can appear to be present.
- Don't necessarily believe that a significant heteroskedasticity test means that heteroskedasticity is the problem and that using robust standard errors will be an adequate fix.

Model Specification and Heteroskedasticity, Part II

- Critically examine the model from the point of view of economic reasoning and look for any specification problems.
- One very common specification issue with economic data is the choice of functional form.
- Using a logarithmic transformation of the dependent variable has another feature, **variance stabilization**, that is useful in the context of heteroskedastic data.

Model Specification and Heteroskedasticity, Part III

- Economic variables like wages, incomes, house prices, and expenditures are right-skewed, with a long tail to the right.
- The log-normal probability distribution is useful when modeling such variables.
- If the random variable y has a log-normal probability density function, then $\ln(y)$ has a normal distribution, which is symmetrical and bell-shaped, and not skewed.

Model Specification and Heteroskedasticity, Part IV

- The feature of the log-normal random variable that we are now interested in is that its variance increases when its mean and median increase.
- By choosing a log-linear or log-log model, we are implicitly assuming a curvilinear and heteroskedastic relationship between the variables y and x .
- However, there is a linear and homoskedastic relation between $\ln(y)$ and x .

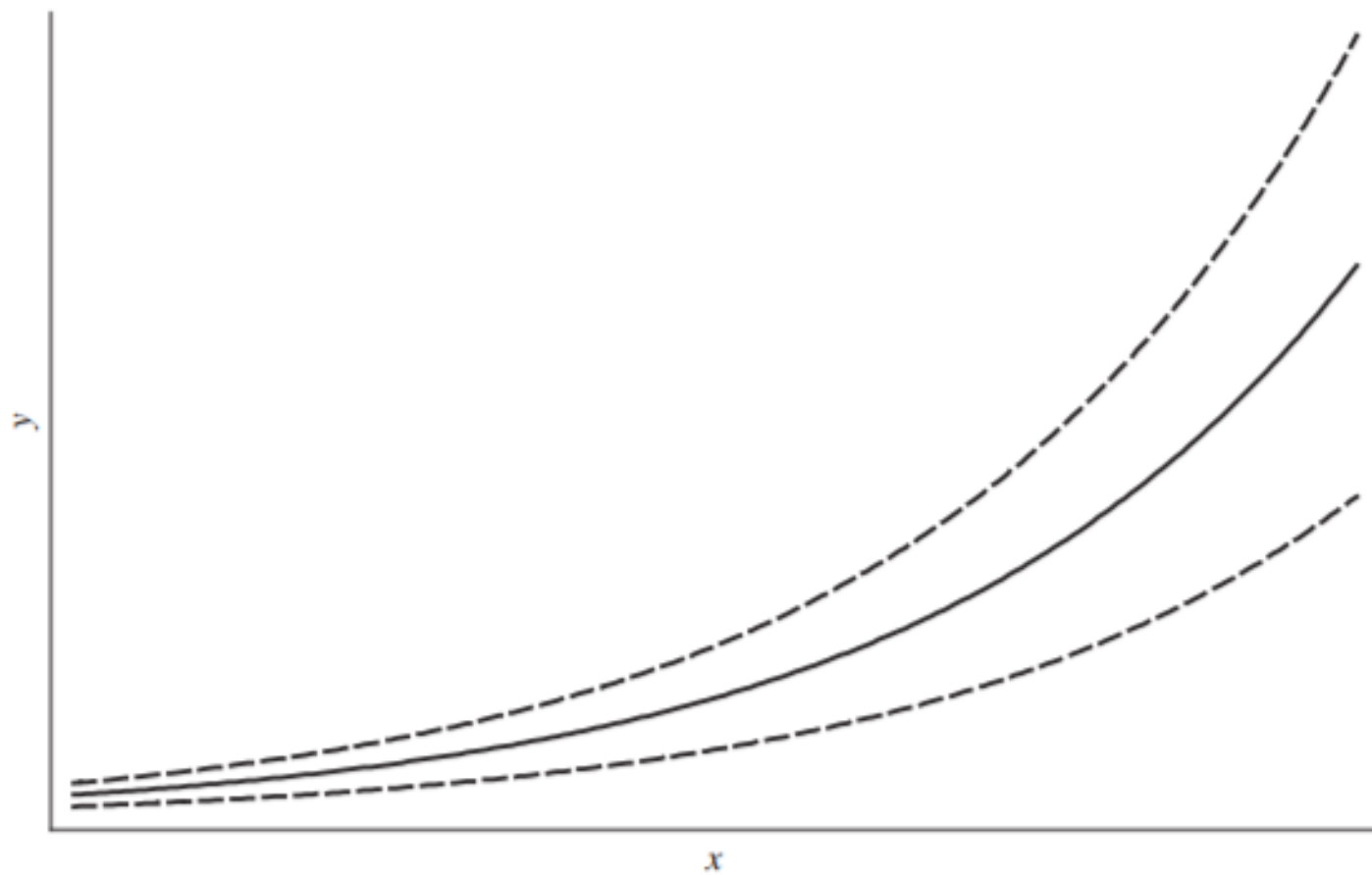


FIGURE 8.6 A log-linear relationship.

Example: Variance Stabilizing Log-Transformation, Part I

- Consider the data file `cex5_small`.
- Figure 8.7(a) shows a histogram of household expenditures on entertainment per person, `ENTERT`, for those households who have positive spending, and Figure 8.7(b) is the histogram for $\ln(\text{ENTERT})$.
- The extremely skewed distribution of entertainment expenditures shows the effect of the log-transformation.

Example: Variance Stabilizing Log-Transformation, Part II

- The variation in ENTERT about the fitted line increases as INCOME increases.
- Estimating the model $ENTERT = \beta_1 + \beta_2 INCOME + \beta_3 COLLEGE + \beta_4 ADVANCED + e$.
- Obtain the least squares residuals and then estimate by OLS the model.

$$\widehat{e_i^2} = \alpha_1 + \alpha_2 INCOME_i + v_i$$

Example: Variance Stabilizing Log-Transformation, Part III

- From this regression, $NR^2 = 31.34$. The critical value for a 1% level of significance.
- Heteroskedasticity test is 6.635; thus, we conclude that heteroskedasticity is present.
- There is little if any visual evidence of heteroskedasticity, and the value of the heteroskedasticity test statistic is $NR^2 = 0.36$.
- So we do not reject the null hypothesis of homoskedasticity. The log-transformation has “cured” the heteroskedasticity problem.

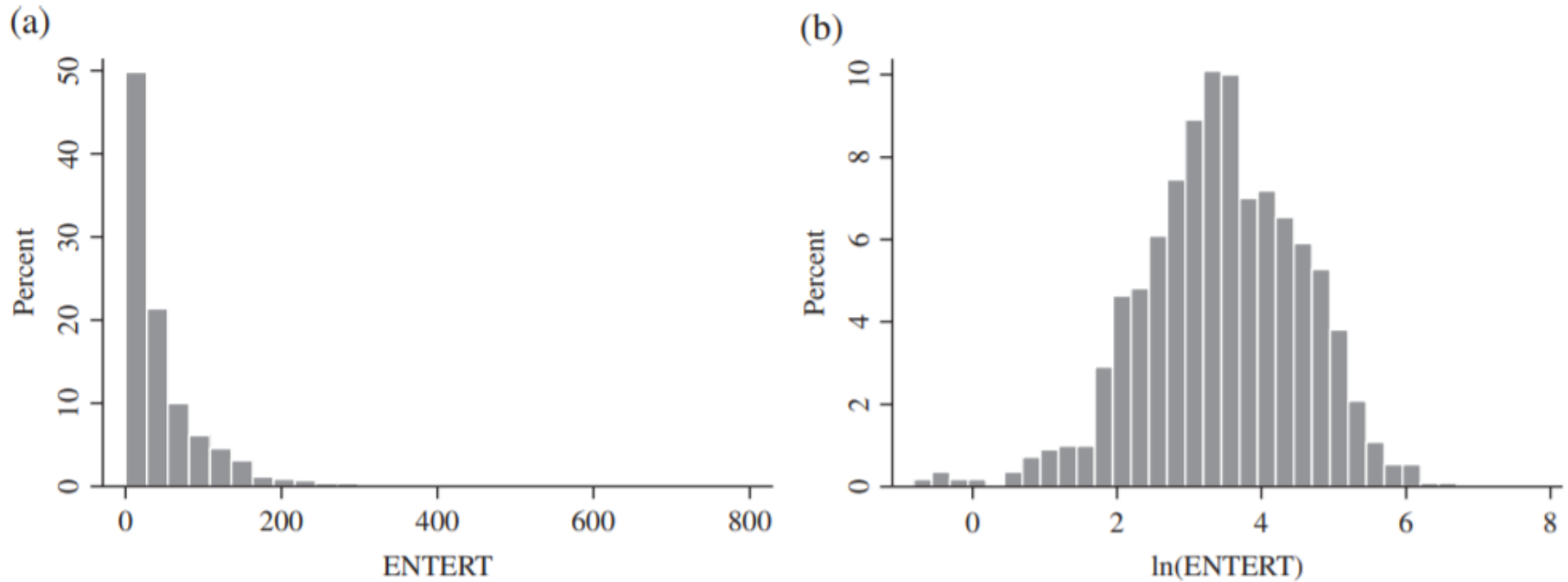


FIGURE 8.7 Histograms of entertainment expenditures.

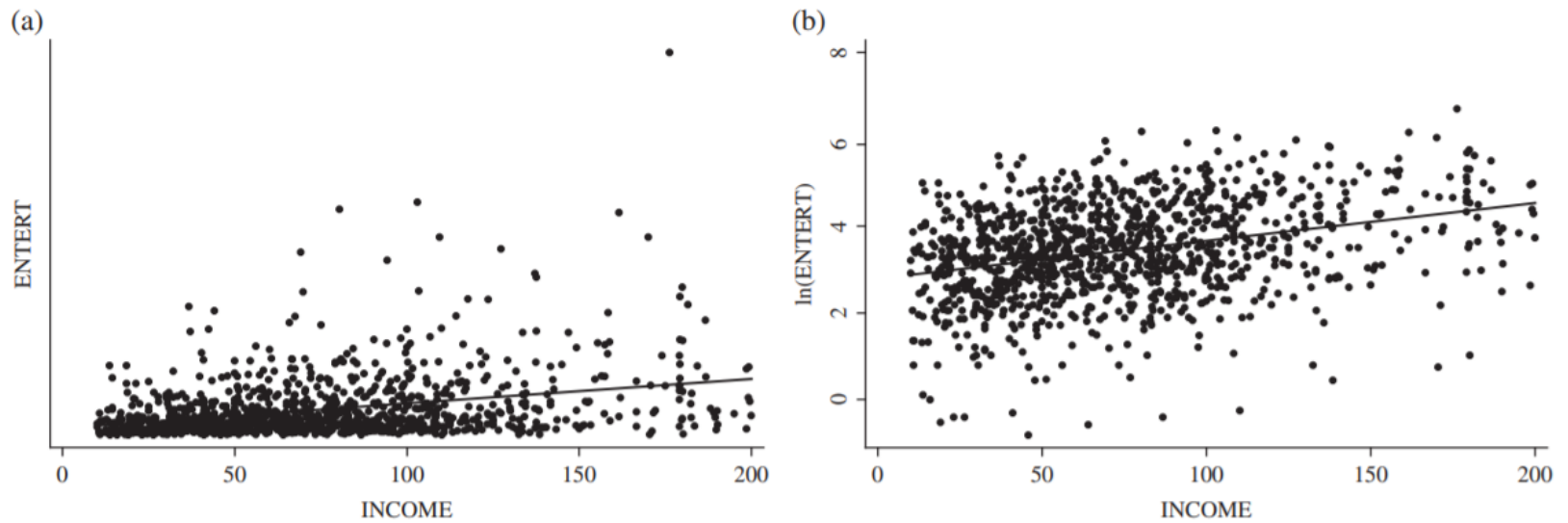


FIGURE 8.8 Linear and log-linear models for entertainment expenditures.

FORDHAM
THE JESUIT UNIVERSITY OF NEW YORK

Gabelli School
of Business

Key Words

- Breusch–Pagan test
- Generalized least squares (GLS)
- Goldfeld–Quandt test
- Grouped heteroskedasticity
- Heteroskedasticity
- Heteroskedasticity-consistent standard errors
- Homoskedasticity
- Lagrange multiplier test
- Linear probability model
- Mean function
- Residual plot
- Robust standard errors
- Skedastic function
- Transformed model
- Variance function
- Weighted least squares (WLS)
- White test

Copyright

Copyright © 2018 John Wiley & Sons, Inc.

All rights reserved. Reproduction or translation of this work beyond that permitted in Section 117 of the 1976 United States Act without the express written permission of the copyright owner is unlawful. Request for further information should be addressed to the Permissions Department, John Wiley & Sons, Inc. The purchaser may make back-up copies for his/her own use only and not for distribution or resale. The Publisher assumes no responsibility for errors, omissions, or damages, caused by the use of these programs or from the use of the information contained herein.

FORDHAM
THE JESUIT UNIVERSITY OF NEW YORK

Gabelli School
of Business