

Ch06. Multiple Regression Analysis: Further Issues

Ping Yu

HKU Business School
The University of Hong Kong

Effects of Data Scaling on OLS Statistics

Section 2.4a: SLR

- The fitted regression line in the example of CEO salary and return on equity is

$$\widehat{\text{salary}} = 963.191 + 18.501 \text{roe},$$

$$n = 209, R^2 = .0132$$

- How will the intercept and slope estimates change when the units of measurement of the dependent and independent variables changes?
- Suppose the salary is measured in dollars rather than thousands of dollars, the intercept should be 963,191 and the slope should be 18,501. (*why?*)
- Solution:** convert a new problem to an old problem whose solution is known.

A Brute-Force Solution

- Suppose $y_i^* = w_1 y_i$ and $x_i^* = w_2 x_i$.
- Then

$$\hat{\beta}_1^* = \frac{\sum_{i=1}^n (x_i^* - \bar{x}^*) y_i^*}{\sum_{i=1}^n (x_i^* - \bar{x}^*)^2} = \frac{w_1 w_2 \sum_{i=1}^n (x_i - \bar{x}) y_i}{w_2^2 \sum_{i=1}^n (x_i - \bar{x})^2} = \frac{w_1}{w_2} \hat{\beta}_1,$$

and

$$\hat{\beta}_0^* = \bar{y}^* - \bar{x}^* \hat{\beta}_1^* = w_1 \bar{y} - w_2 \bar{x} \frac{w_1}{w_2} \hat{\beta}_1 = w_1 (\bar{y} - \bar{x} \hat{\beta}_1) = w_1 \hat{\beta}_0,$$

where $\bar{x}^* = w_2 \bar{x}$ and $\bar{y}^* = w_1 \bar{y}$.

- But when the number of regressors is large, the formula of $\hat{\beta}$ is complicated.
- The textbook derives the relationship between $\hat{\beta}^*$ and $\hat{\beta}$ by FOCs, while we employ the objective functions.

Relationship Between Objective Functions

- Recall that

$$\left(\widehat{\beta}_0, \widehat{\beta}_1\right) = \arg \min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2, \quad (1)$$

where arg means arguments.

- Now, for the rescaled data,

$$\begin{aligned} & \min_{\beta_0^*, \beta_1^*} \sum_{i=1}^n (y_i^* - \beta_0^* - \beta_1^* x_i^*)^2 \\ &= \min_{\beta_0^*, \beta_1^*} \sum_{i=1}^n (w_1 y_i - \beta_0^* - \beta_1^* w_2 x_i)^2 \\ &= \min_{\beta_0^*, \beta_1^*} w_1^2 \sum_{i=1}^n \left(y_i - \frac{\beta_0^*}{w_1} - \beta_1^* \frac{w_2}{w_1} x_i \right)^2, \end{aligned}$$

where note that since (β_0^*, β_1^*) can be freely chosen, $\left(\frac{\beta_0^*}{w_1}, \beta_1^* \frac{w_2}{w_1}\right)$ can also be freely chosen although (w_1, w_2) are fixed, just as (β_0, β_1) in (1).

Relationship Between $\hat{\beta}$'s and $\hat{\sigma}^2$'s

- So (why?)

$$\begin{aligned}\hat{\beta}_0 &= \frac{\hat{\beta}_0^*}{w_1} \text{ and } \hat{\beta}_1 = \hat{\beta}_1^* \frac{w_2}{w_1} \\ \Rightarrow \hat{\beta}_0^* &= w_1 \hat{\beta}_0 \text{ and } \hat{\beta}_1^* = \frac{w_1}{w_2} \hat{\beta}_1\end{aligned}$$

- Also,

$$\begin{aligned}\hat{\sigma}^{*2} &= \frac{1}{n-2} \sum_{i=1}^n \hat{u}_i^{*2} = \frac{1}{n-2} \sum_{i=1}^n (y_i^* - \hat{\beta}_0^* - \hat{\beta}_1^* x_i^*)^2 \\ &= \frac{w_1^2}{n-2} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \\ &= \frac{w_1^2}{n-2} \sum_{i=1}^n \hat{u}_i^2 \\ &= w_1^2 \hat{\sigma}^2.\end{aligned}$$

- Or, the standard error of the regression (SER) $\hat{\sigma}^* = w_1 \hat{\sigma}$.

Relationship Between R^{*2} 's and Standard Errors

- It turns out that $R^{*2} = R^2$:

$$\begin{aligned}
 R^{*2} &= 1 - \frac{SSR^*}{SST^*} \\
 &= 1 - \frac{\sum_{i=1}^n (y_i^* - \hat{\beta}_0^* - \hat{\beta}_1^* x_i^*)^2}{\sum_{i=1}^n (y_i^* - \bar{y}^*)^2} \\
 &= 1 - \frac{w_1^2 \sum_{i=1}^n \hat{u}_i^2}{w_1^2 \sum_{i=1}^n (y_i - \bar{y})^2} \\
 &= 1 - \frac{SSR}{SST} = R^2
 \end{aligned}$$

- Standard Errors:

$$\begin{aligned}
 \hat{\beta}_0^* &= w_1 \hat{\beta}_0 \text{ and } \hat{\beta}_1^* = \frac{w_1}{w_2} \hat{\beta}_1 \\
 \implies \text{se}(\hat{\beta}_0^*) &= w_1 \cdot \text{se}(\hat{\beta}_0) \text{ and } \text{se}(\hat{\beta}_1^*) = \frac{w_1}{w_2} \cdot \text{se}(\hat{\beta}_1)
 \end{aligned}$$

Relationship Between t Statistics and CIs (Section 6.1: MLR)

- It is natural to predict that t statistic is the same as before:

$$t_{\hat{\beta}_1^*} = \frac{\hat{\beta}_1^*}{\text{se}(\hat{\beta}_1^*)} = \frac{\frac{w_1}{w_2} \hat{\beta}_1}{\frac{w_1}{w_2} \cdot \text{se}(\hat{\beta}_1)} = \frac{\hat{\beta}_1}{\text{se}(\hat{\beta}_1)} = t_{\hat{\beta}_1}.$$

- The CI for β_1 is the original CI multiplied by $\frac{w_1}{w_2}$:

$$\begin{aligned} & \left[\hat{\beta}_1^* - 1.96 \cdot \text{se}(\hat{\beta}_1^*), \hat{\beta}_1^* + 1.96 \cdot \text{se}(\hat{\beta}_1^*) \right] \\ &= \frac{w_1}{w_2} \left[\hat{\beta}_1 - 1.96 \cdot \text{se}(\hat{\beta}_1), \hat{\beta}_1 + 1.96 \cdot \text{se}(\hat{\beta}_1) \right]. \end{aligned}$$

- The results for β_0 are similar; the only difference is to replace $\frac{w_1}{w_2}$ by w_1 .

Unit Change in Logarithmic Form

- Changing the unit of measurement of y and x , when they appear in logarithmic form, does not affect any of the slope estimates, but may affect the intercept estimate.
- why?

$$\begin{aligned}
 & \min_{\beta_0^*, \beta_1^*} \sum_{i=1}^n [\log(y_i^*) - \beta_0^* - \beta_1^* \log(x_i^*)]^2 \\
 = & \min_{\beta_0^*, \beta_1^*} \sum_{i=1}^n [\log(y_i) + \log(w_1) - \beta_0^* - \beta_1^* \log(x_i) - \beta_1^* \log(w_2)]^2 \\
 = & \min_{\beta_0^*, \beta_1^*} \sum_{i=1}^n [\log(y_i) - (\beta_0^* + \beta_1^* \log(w_2) - \log(w_1)) - \beta_1^* \log(x_i)]^2,
 \end{aligned}$$

so

$$\begin{aligned}
 \hat{\beta}_0 &= \hat{\beta}_0^* + \hat{\beta}_1^* \log(w_2) - \log(w_1) \quad \text{and} \quad \hat{\beta}_1 = \hat{\beta}_1^* \\
 \implies \hat{\beta}_0^* &= \hat{\beta}_0 - \hat{\beta}_1 \log(w_2) + \log(w_1) \quad \text{and} \quad \hat{\beta}_1^* = \hat{\beta}_1.
 \end{aligned}$$

- I.e., the elasticity is invariant to the units of measurement of either y or x , and the intercept is related to both the original intercept and slope.

Example: Japanese Learning

- Suppose we want to study the SLR,

$$y_i = \beta_0 + \beta_1 x_i + u_i,$$

where

y_i = exam mark in Japanese language course

x_i = hours of study per day during the semester (average hours)

- The fitted regression line is

$$\hat{y}_i = 20 + 14x_i$$

(5.6) (3.5)

- If you do not study at all, the predicted mark is 20. One additional hour of study per day increases exam mark by 14 marks. At the mean value $\bar{x} = 3$ hours of study per day is expected to result in a mark of $\bar{y} = 62$.

continue

- Now if we report hours of study per week (x_i^*) rather than per day, the variable has been scaled:

$$x_i^* = w_2 x_i = 7x_i \text{ and } y_i^* = w_1 y_i = y_i.$$

- If we run the regression based on x_i^* and y_i , we get

$$\hat{y}_i = 20 + 2x_i^* \quad (5.6) \quad (0.5)$$

- Each additional hour of study per week increases the exam mark by 2 marks [[intuition here](#)].
- t statistic remains the same: $\frac{2}{0.5} = \frac{14}{3.5}$.
- The CI for β_1^* , $[2 - 1.96 \times 0.5, 2 + 1.96 \times 0.5] = [1.02, 2.98]$, is 1/7 of the CI for β_1 , which is $[14 - 1.96 \times 3.5, 14 + 1.96 \times 3.5] = [7.14, 20.98]$.
- Note that the means \bar{x}^* and \bar{y}^* will still be on the newly estimated regression line:

$$\begin{aligned} \bar{y}^* &= 20 + 2\bar{x}^* \\ 62 &= 20 + 2 \times 21 \end{aligned}$$

More on Functional Form

a: More on Using Logarithmic Functional Forms

- Logarithmic transformations have the convenient percentage/elasticity interpretation.
- Slope coefficients of logged variables are invariant to rescalings.
- Taking logs often eliminates/mitigates problems with outliers. ([why? figure here](#))
- Taking logs often helps to secure normality (e.g., $\log(\text{wage})$ vs. wage) and homoskedasticity (see Chapter 8 for an example).
- Variables measured in units such as years (e.g., education, experience, tenure, age, etc) should not be logged.
- Variables measured in percentage points (e.g., unemployment rate, participation rate of a pension plan, etc.) should also not be logged.
- Logs must not be used if variables take on zero or negative values (e.g., hours of work during a month).
- It is hard to reverse the log-operation when constructing predictions. (we will discuss more on this point later in this chapter)

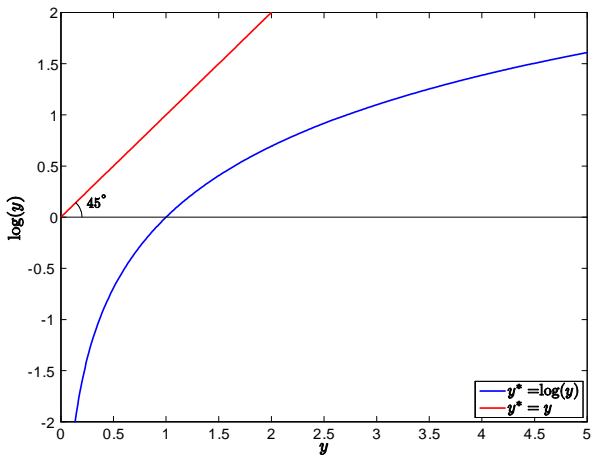


Figure: $\log y < y$ and $\lim_{y \rightarrow \infty} \frac{\log y}{y} = 0$

b: Models with Quadratics

- **Example:** Suppose the fitted regression line for the wage equation is

$$\widehat{wage} = 3.73 + .298 \text{exper} - .0061 \text{exper}^2$$

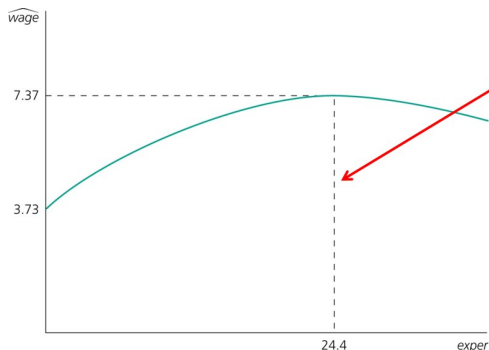
$$(.35) \quad (.041) \quad (.0009)$$

$$n = 526, R^2 = .093 \text{ (quite small)}$$

- The predicted *wage* is a concave function of *exper*. [[figure here](#)]
- The marginal effect of *exper* on wage is

$$\frac{\partial \widehat{wage}}{\partial \text{exper}} = \hat{\beta}_1 + 2\hat{\beta}_2 \text{exper} = .298 - 2 \times .0061 \text{exper}.$$

- The first year of experience increases the wage by some \$.30, the second year by $.298 - 2(.0061)(1) = \$.29 < \$.30$ etc.



$$x^* = \left| \frac{\hat{\beta}_1}{2\hat{\beta}_2} \right| = \left| \frac{.298}{2(.0061)} \right| \approx 24.4$$

Does this mean the return to experience becomes negative after 24.4 years?

Not necessarily. It depends on how many observations in the sample lie right of the turnaround point.

In the given example, these are about 28% of the observations. There may be a specification problem (e.g. omitted variables).

Figure: Wage Maximum with Respect to Work Experience

- Mincer (1974, *Schooling, Experience and Earnings*) assumed $\log(\text{wage}) = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{exper}^2 + u$. [[photo here](#)]

History of the Wage Equation



Jacob Mincer (1922-2006), Columbia,
father of modern labor economics

Example: Effects of Pollution on Housing Prices

- The fitted regression line is

$$\begin{aligned} \widehat{\log(\text{price})} &= 13.39 - .902 \log(\text{nox}) - .087 \log(\text{dist}) \\ &\quad (.57) \quad (.115) \quad (.043) \\ &\quad - .545 \text{rooms} + .062 \text{rooms}^2 - .048 \text{stratio} \\ &\quad (.165) \quad (.013) \quad (.006) \\ n &= 506, R^2 = .603 \end{aligned}$$

where

nox = nitrogen oxide in air

dist = distance from employment centers, in miles

stratio = student/teacher ratio

- The predicted $\log(\text{price})$ is a convex function of *rooms*. [\[figure here\]](#)
- The coefficient of *rooms* is negative. Does this mean that, at a low number of rooms, more rooms are associated with lower prices?
- The marginal effect of *rooms* on $\log(\text{price})$ is

$$\frac{\partial \log(\text{price})}{\partial \text{rooms}} = \frac{\partial \text{price} / \text{price}}{\partial \text{rooms}} = -.545 + 2 \times .062 \text{rooms}.$$

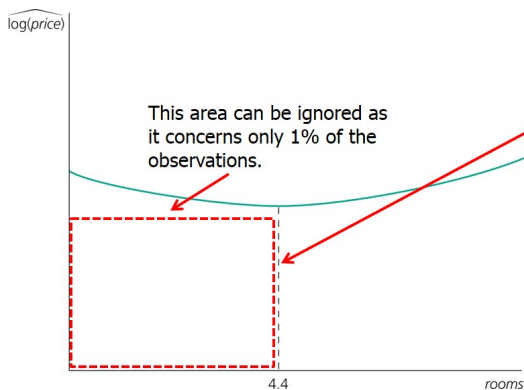


Figure: Calculation of the Turnaround Point

Other Possibilities

- Using Quadratics Along with Logarithms:

$$\log(\text{price}) = \beta_0 + \beta_1 \log(\text{nox}) + \beta_2 \log(\text{nox})^2 + \beta_3 \text{crime} + \beta_4 \text{rooms} + \beta_5 \text{rooms}^2 + \beta_6 \text{stratio} + u,$$

which implies

$$\frac{\partial \log(\text{price})}{\partial \log(\text{nox})} = \frac{\% \partial \text{price}}{\% \partial \text{nox}} = \beta_1 + 2\beta_2 \log(\text{nox}).$$

- Higher Order Polynomials:** It is often assumed that the total cost takes the following form,

$$\text{cost} = \beta_0 + \beta_1 \text{quantity} + \beta_2 \text{quantity}^2 + \beta_3 \text{quantity}^3 + u,$$

which implies a U-shaped marginal cost (MC), where β_0 is the total fixed cost. [\[figure here\]](#)

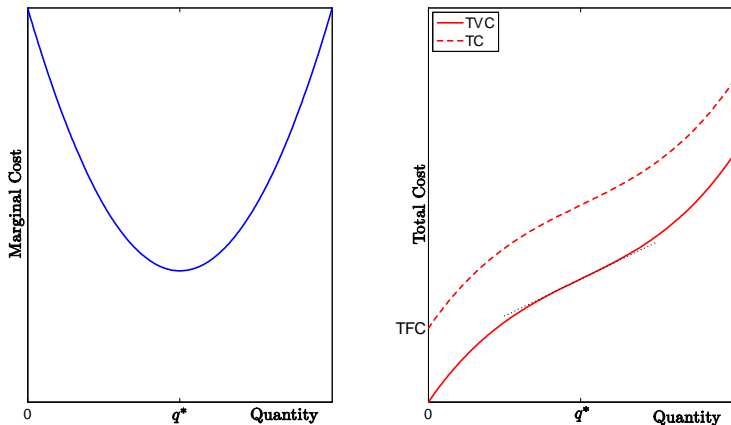


Figure: Quadratic MC Implies Cubic TC: q^* is the inflection point

c: Models with Interaction Terms

- In the model

$$price = \beta_0 + \beta_1 sqft + \beta_2 bdrms + \beta_3 sqft \cdot bdrms + \beta_4 bthrms + u,$$

$sqft \cdot bdrms$ is the **interaction term**.

- The marginal effect of $bdrms$ on $price$ is

$$\frac{\partial price}{\partial bdrms} = \beta_2 + \beta_3 sqft.$$

- The effect of the number of bedrooms depends on the level of square footage.
- Interaction effects complicate interpretation of parameters: β_2 is the effect of number of bedrooms, but for a square footage of zero.
- How to avoid this interpretation difficulty?

Reparametrization of Interaction Effects

- The model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + u$$

can be reparametrized as

$$y = \alpha_0 + \delta_1 x_1 + \delta_2 x_2 + \beta_3 (x_1 - \mu_1)(x_2 - \mu_2) + u,$$

where $\mu_1 = E[x_1]$ and $\mu_2 = E[x_2]$ are population means of x_1 and x_2 , and can be replaced by their sample means.

- What is the relationship between $(\hat{\alpha}_0, \hat{\delta}_1, \hat{\delta}_2)$ and $(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2)$? (Exercise)

- Now,

$$\frac{\partial y}{\partial x_2} = \delta_2 + \beta_3 (x_1 - \mu_1),$$

i.e., δ_2 is the effect of x_2 if all other variables take on their mean values.

- Advantages** of reparametrization:
 - It is easy to interpret all parameters.
 - Standard errors for partial effects at the mean values are available.
 - If necessary, interaction may be centered at other interesting values.

More on Goodness-of-Fit and Selection of Regressors

a: Adjusted R -Squared

- General remarks on R -squared:
 - A high R -squared does not imply that there is a causal interpretation.
 - A low R -squared does not preclude precise estimation of partial effects.
- Recall that

$$R^2 = 1 - \frac{SSR/n}{SST/n} = 1 - \frac{\tilde{\sigma}_u^2}{\tilde{\sigma}_y^2},$$

so R^2 is estimating the **population R -squared**

$$\rho^2 = 1 - \frac{\sigma_u^2}{\sigma_y^2},$$

the proportion of the variation in y in the population explained by the independent variables.

- **Adjusted R -Squared:**

$$\bar{R}^2 = 1 - \frac{SSR / (n - k - 1)}{SST / (n - 1)} = 1 - \frac{\hat{\sigma}_u^2}{\hat{\sigma}_y^2},$$

is sometimes also called **R -bar squared** [[photo here](#)], where $\hat{\sigma}_u^2$ and $\hat{\sigma}_y^2$ are unbiased estimators of σ_u^2 and σ_y^2 due to the correction of dfs.

History of \bar{R}^2 

Henri Theil (1924-2000)¹, Chicago and Florida

¹He is a Dutch econometrician. Two other Dutch econometricians, Jan Tinbergen (1903-1994) and Tjalling Koopmans (1910-1985) won the Nobel Prize in economics in 1969 and 1975, respectively.

continue

- \bar{R}^2 takes into account degrees of freedom of the numerator and denominator, so is generally a better measure of goodness-of-fit.
- \bar{R}^2 imposes a penalty for adding new regressors: $k \uparrow \implies \bar{R}^2 \downarrow$
- \bar{R}^2 increases if and only if the t statistic of a newly added regressor is greater than one in absolute value. [proof not required]
 - Compared with $y = \beta_0 + \beta_1 x_1 + u$, the regression $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$ has a larger \bar{R}^2 if and only if

$$\left| t_{\hat{\beta}_2} \right| > 1.$$

- Relationship between R^2 and \bar{R}^2 : by

$$1 - R^2 = \frac{SSR}{SST} = \frac{n-k-1}{n-1} \frac{SSR/(n-k-1)}{SST/(n-1)} = \frac{n-k-1}{n-1} (1 - \bar{R}^2),$$

we have

$$\bar{R}^2 = 1 - (1 - R^2) \frac{n-1}{n-k-1} < R^2$$

unless $k = 0$ or $R^2 = 1$. [figure here]

- Note that \bar{R}^2 even gets negative if $R^2 < \frac{k}{n-1}$.

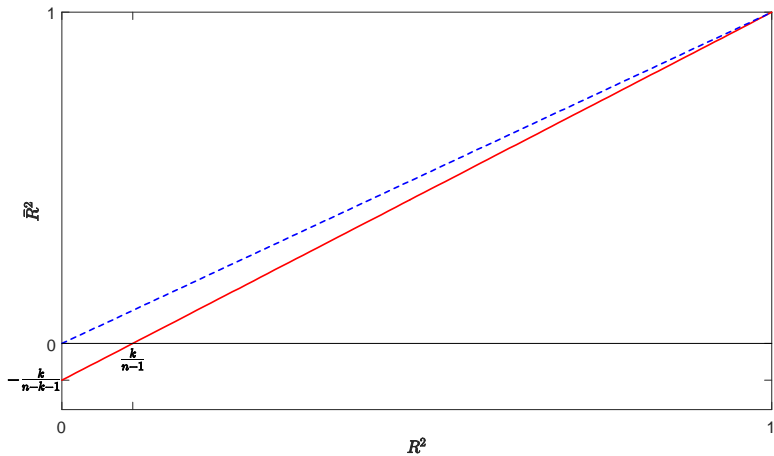


Figure: Relationship Between \bar{R}^2 and R^2

b: Using Adjusted R -squared to Choose between Nonnested Models

- Models are **nonnested** if neither model is a special case of the other.
- For example, to incorporate diminishing return of *sales* to R&D, we consider two models:

$$\begin{aligned} rdintens &= \beta_0 + \beta_1 \log(sales) + u, \\ rdintens &= \beta_0 + \beta_1 sales + \beta_2 sales^2 + u, \end{aligned}$$

where

$rdintens$ = R&D intensity.

- $R^2 = .061$ and $\bar{R}^2 = .030$ in model 1 and $R^2 = .148$ and $\bar{R}^2 = .090$ in model 2.
- A comparison between the R -squared of both models would be unfair to the first model because the first model contains fewer parameters.
- In the given example, even after adjusting for the difference in degrees of freedom, the quadratic model is preferred.

Comparing Models with Different Dependent Variables

- R -squared or adjusted R -squared must not be used to compare models which differ in their definition of the dependent variable.
- **Example** (CEO Compensation and Firm Performance):

$$\widehat{salary} = 830.63 + .0163sales + 19.63roe$$

$$(223.90)(.0089) \quad (11.08)$$

$$n = 209, R^2 = .029, \bar{R}^2 = .020, SST = 391,732,982$$

and

$$\widehat{lsalary} = 4.36 + .275lsales + .0179roe$$

$$(0.29)(.033) \quad (.0040)$$

$$n = 209, R^2 = .282, \bar{R}^2 = .275, SST = 66.72$$

- There is much less variation in $\log(salary)$ that needs to be explained than in $salary$, so it is not fair to compare R^2 and \bar{R}^2 of the two models. (we will discuss how to compare the fitting of these two models later in this chapter)

c: Controlling for Too Many Factors in Regression Analysis

- In some cases, certain variables should not be held fixed:
 - In a regression of traffic fatalities on state beer taxes (and other factors) one should not directly control for beer consumption.
 - **why?** Beer taxes influence traffic fatalities only through beer consumption; if beer consumption is controlled, then the coefficient of beer taxes measures the indirect effect of beer taxes, which is hardly interesting.
 - In a regression of family health expenditures on pesticide usage among farmers one should not control for doctor visits.
 - **why?** Health expenditures include doctor visits, and we would like to pick up all effects of pesticide use on health expenditure.
- Different regressions may serve different purposes:
 - In a regression of house prices on house characteristics, one would include price assessments and also housing attributes if the purpose of the regression is to study the validity of assessments; one should not include price assessments if the purpose of the regression is to estimate a **hedonic price model**,² which measures the marginal values of various housing attributes.

²What consumers are seeking to acquire is not goods themselves (e.g. cars or train journeys) but the characteristics they contain (e.g., display of fashion sense, transport from A to B).

History of Hedonic Price Model

- **Hedonic Utility:** Lancaster, Kelvin J., 1966, A New Approach to Consumer Theory, *Journal of Political Economy*, 74, 132-157.
- **Hedonic Pricing:** Rosen, S., 1974, Hedonic Prices and Implicit Markets: Product Differentiation in Pure Competition, *Journal of Political Economy*, 82, 34-55.



Kelvin J. Lancaster (1924-1999), Columbia



Sherwin Rosen (1938-2001)³, Chicago

³His student Robert H. Thaler (1945-) at the University of Chicago won the Nobel Prize in economics in 2017.

d: Adding Regressors to Reduce the Error Variance

- Recall that

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{SST_j(1 - R_j^2)}.$$

- Adding regressors may exacerbate multicollinearity problems ($R_j^2 \uparrow$).
- On the other hand, adding regressors reduces the error variance ($\sigma^2 \downarrow$).
- Variables that are uncorrelated with other regressors should be added because they reduce error variance ($\sigma^2 \downarrow$) without increasing multicollinearity (R_j^2 remains the same).
- However, such uncorrelated variables may be hard to find.
- Example** (Individual Beer Consumption and Beer Prices): Including individual characteristics in a regression of beer consumption on beer prices leads to more precise estimates of the price elasticity if individual characteristics are uncorrelated with beer prices.

$$\log(\text{cons}) = \beta_0 + \beta_1 \log(\text{price}) + \underbrace{\text{indchar}}_{\text{uncorrelated with } \log(\text{price})} + u.$$

(**) Prediction and Residual Analysis

c: Predicting y When $\log(y)$ is the Dependent Variable

- We study only this prediction problem as promised.
- Note that

$$\log(y) = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k + u$$

implies

$$y = \exp(\beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k) \exp(u) = m(\mathbf{x}) \exp(u).$$

- Under the additional assumption that u is independent of (x_1, \dots, x_k) , we have

$$\begin{aligned} E[y|\mathbf{x}] &= \exp(\beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k) E[\exp(u)|\mathbf{x}] \\ &= \exp(\beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k) E[\exp(u)] \\ &\equiv m(\mathbf{x}) \alpha_0, \end{aligned}$$

where the second equality is due to the independence between u and \mathbf{x} , so the predicted y is

$$\hat{y} = \hat{m}(\mathbf{x}) \hat{\alpha}_0$$

where

$$\hat{m}(\mathbf{x}) = \exp(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \cdots + \hat{\beta}_k x_k) \text{ and } \hat{\alpha}_0 = \frac{1}{n} \sum_{i=1}^n \exp(\hat{u}_i).$$

$$E[\exp(u)] \geq 1$$

- Recall that $E[u] = 0$, so

$$E[\exp(u)] \geq \exp(E[u]) = \exp(0) = 1.$$

- In the following figure, suppose u takes only two values u_1 and u_2 with probability $\frac{1}{2}$ and $\frac{1}{2}$, respectively. Since $E[u] = \frac{1}{2}(u_1 + u_2) = 0$, $u_1 = -u_2$.
- Now,

$$\begin{aligned} E[\exp(u)] &= \frac{1}{2}(\exp(u_1) + \exp(u_2)) \\ &\geq \exp\left(\frac{1}{2}(u_1 + u_2)\right) \\ &= \exp(E[u]) = 1, \end{aligned}$$

where the equality is achieved only if $u_1 = u_2 = 0$, i.e., $u = 0$.

- As a result, $\tilde{y} = \hat{m}(\mathbf{x}) = \exp(\widehat{\log(y)})$ under-estimates $E[y|\mathbf{x}]!$

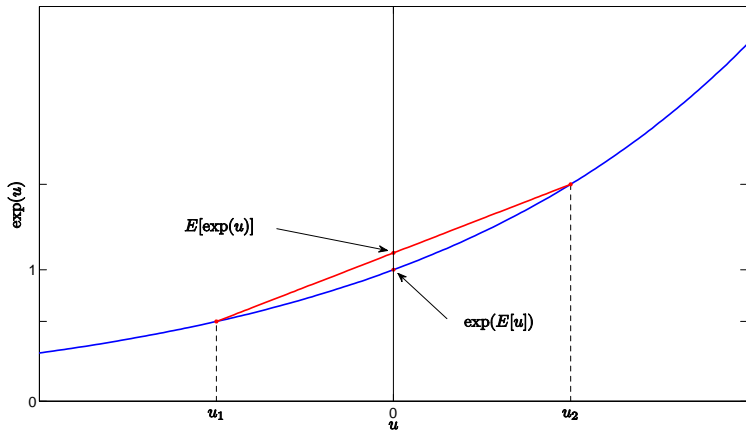


Figure: Illustration of Jensen's Inequality

Comparing R -Squared of a Logged and an Unlogged Specification

- Reconsider the CEO salary problem:

$$\widehat{\text{salary}} = 613.43 + .0190\text{sales} + .0234\text{mktval} + 12.70\text{ceoten}$$

$$(65.23) \quad (.0100) \quad (.0095) \quad (5.61)$$

$$n = 177, R^2 = .201$$

and

$$\widehat{\text{lsalary}} = 4.504 + .163\text{lsales} + .0109\text{mktval} + .0117\text{ceoten}$$

$$(0.257) \quad (.039) \quad (.050) \quad (.0053)$$

$$n = 177, \tilde{R}^2 = .318$$

- R^2 and \tilde{R}^2 are the R -squareds for the predictions of the unlogged salary variable (although the second regression is originally for logged salaries). Both R -squareds can now be directly compared

About \tilde{R}^2

- Recall that

$$R^2 = \widehat{\text{Corr}}(y, \hat{y})^2,$$

where \hat{y} is the predicted value of y .

- When *Salary* is the dependent variable, the predicted value of y is $\hat{m}(\mathbf{x})\hat{\alpha}_0 = \hat{\alpha}_0\tilde{y}$.
- Since $\hat{\alpha}_0 > 0$,

$$\widehat{\text{Corr}}(y, \hat{y}) = \widehat{\text{Corr}}(y, \hat{\alpha}_0\tilde{y}) = \widehat{\text{Corr}}(y, \tilde{y})$$

invariant to $\hat{\alpha}_0$, where recall that for any $a > 0$,

$$\text{Corr}(X, aY) = \text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}.$$

- As a result,

$$\tilde{R}^2 = \widehat{\text{Corr}}(y, \tilde{y})^2.$$