

# Regression Diagnostics

## (SW 9.2 and 8.2)

Dragos Ailoae  
dailoae@gradcenter.cuny.edu

Advanced Economics and Business Statistics  
ECON-4400w

Brooklyn College

# Regression User's Guide (1 of 2)

What Can Go Wrong?	What Are the Consequences?	How Can It Be Detected?	How Can It Be Corrected?
<b>Omitted Variable</b>			
The omission of a relevant independent variable	Bias in the coefficient estimates (the $\hat{\beta}$ s) of the included Xs.	Theory, significant unexpected signs, or surprisingly poor fits.	Include the omitted variable or a proxy.
<b>Irrelevant Variable</b>			
The inclusion of a variable that does not belong in the equation	Decreased precision in the form of higher standard errors, lower $t$ -scores and wider confidence intervals.	<ol style="list-style-type: none"> <li>1. Theory</li> <li>2. <math>t</math>-test on <math>\hat{\beta}</math></li> <li>3. <math>\bar{R}^2</math></li> <li>4. Impact on other coefficients if X is dropped.</li> </ol>	Delete the variable if its inclusion is not required by the underlying theory.
<b>Incorrect Functional Form</b>			
The functional form is inappropriate	Biased estimates, poor fit, and difficult interpretation.	Examine the theory carefully; think about the relationship between X and Y.	Transform the variable or the equation to a different functional form.

# Regression User's Guide (2 of 2)

What Can Go Wrong?	What Are the Consequences?	How Can It Be Detected?	How Can It Be Corrected?
<b>Multicollinearity</b> Some of the independent variables are (imperfectly) correlated	No biased $\hat{\beta}$ s, but estimates of the separate effects of the Xs are not reliable, i.e., high $SE(\hat{\beta})$ s and low $t$ -scores.	Pairwise correlations or scatterplots	Drop redundant variables, but to drop others might introduce bias. Often doing nothing is best.
<b>Serial Correlation</b> Observations of the error term are correlated, as in: $\epsilon_t = \rho\epsilon_{t-1} + u_t$	No biased $\hat{\beta}$ s, but OLS no longer is minimum variance, and hypothesis testing and confidence intervals are unreliable.	Use residual plots	If impure, fix the specification.
<b>Heteroskedasticity</b> The variance of the error term is not constant for all observations, as in: $VAR(\epsilon_i) = \sigma^2 Z_i$	Same as for serial correlation.	Use residual plots	If impure, fix the specification. Otherwise, use robust std. errors or reformulate the variables.

# Functional form (SW 8.2)

## Logarithms refresher

10x Larger ↑  
10x Smaller ↓

Number	How Many 10s	Base-10 Logarithm	
.. etc..			
1000	$1 \times 10 \times 10 \times 10$	$\log_{10}(1000)$	$= 3$
100	$1 \times 10 \times 10$	$\log_{10}(100)$	$= 2$
10	$1 \times 10$	$\log_{10}(10)$	$= 1$
1	1	$\log_{10}(1)$	$= 0$
0.1	$1 \div 10$	$\log_{10}(0.1)$	$= -1$
0.01	$1 \div 10 \div 10$	$\log_{10}(0.01)$	$= -2$
0.001	$1 \div 10 \div 10 \div 10$	$\log_{10}(0.001)$	$= -3$
.. etc..			

# Functional form (SW 8.2)

## Converting between log bases

$$\begin{aligned}\log_5(12) &= \frac{\log_{10}(12)}{\log_{10}(5)} \\ &= \frac{\log(12)}{\log(5)} \\ &= \frac{1.079181246...}{0.6989700043...}\end{aligned}$$

$$\log_5(12) \approx 1.544$$

# Functional form (SW 8.2)

## Natural logs (ln)

If  $e$  (a constant equal to 2.71828) to the " $b$ th power" produces  $x$ , then  $b$  is the log of  $x$  :

$b$  is the log of  $x$  to the base  $e$  if:  $e^b = x$

Thus, a log (or logarithm) is the exponent to which a given base must be taken in order to produce a specific number. While logs come in more than one variety, we'll use only natural logs (logs to the base  $e$ ) in this text.

The symbol for a natural log is "ln," so  $\ln(x) = b$  means that  $(2.71828)^b = x$  or, more simply,

$\ln(x) = b$  means that  $e^b = x$

For example, since  $e^2 = (2.71828)^2 = 7.389$ , we can state that:

$$\ln(7.389) = 2$$

Thus, the natural log of 7.389 is 2! Two is the power of  $e$  that produces 7.389 . Let's look at some other natural log calculations:

$$\ln(100) = 4.605$$

$$\ln(1000) = 6.908$$

## Functional form (SW 8.2)

### Logarithmic functions of $Y$ and/or $X$

- $\ln(X)$  = the natural logarithm of  $X$
- Logarithmic transforms permit modeling relations in “percentage” terms (like elasticities), rather than linearly.

*Here's why:*  $\ln(x + \Delta x) - \ln(x) = \ln\left(1 + \frac{\Delta x}{x}\right) \cong \frac{\Delta x}{x}$

$$(\text{calculus: } \frac{d \ln(x)}{dx} = \frac{1}{x})$$

*Numerically:*

$$\ln(1.01) = .00995 \cong .01;$$

$$\ln(1.10) = .0953 \cong .10 \text{ (sort of)}$$

# Functional form (SW 8.2)

## Interpreting coefficients

The best way to choose a functional form for a regression model is to select the specification that best matches the underlying theory of the equation. In a majority of cases, the linear form will be adequate, and for most of the rest, common sense will point out a fairly easy choice from the following alternatives:

Functional Form	Equation (one X only)	The Change in Y when X Changes
Linear	$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$	If X increases by one unit, Y will change by $\beta_1$ units.
Double-log	$\ln Y_i = \beta_0 + \beta_1 \ln X_i + \epsilon_i$	If X increases by one percent, Y will change by $\beta_1$ percent. (Thus $\beta_1$ is the elasticity of Y with respect to X.)
Semilog (lnX)	$Y_i = \beta_0 + \beta_1 \ln X_i + \epsilon_i$	If X increases by one percent, Y will change by $\beta_1/100$ units.
Semilog (lnY)	$\ln Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$	If X increases by one unit, Y will change by roughly $100\beta_1$ percent.
Polynomial	$Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \epsilon_i$	If X increases by one unit, Y will change by $(\beta_1 + 2\beta_2 X)$ units.



# Functional form (SW 8.2)

## *Example: $\ln(\text{TestScore})$ vs. $\ln(\text{Income})$*

- First defining a new dependent variable,  $\ln(\text{TestScore})$ , **and** the new regressor,  $\ln(\text{Income})$
- The model is now a linear regression of  $\ln(\text{TestScore})$  against  $\ln(\text{Income})$ , which can be estimated by OLS:

$$\ln(\widehat{\text{TestScore}}) = 6.336 + 0.0554 \times \ln(\text{Income}_i)$$

(0.006)      (0.0021)

An 1% increase in *Income* is associated with an increase of .0554% in *TestScore* (*Income* up by a factor of 1.01, *TestScore* up by a factor of 1.000554)

# Functional form (SW 8.2)

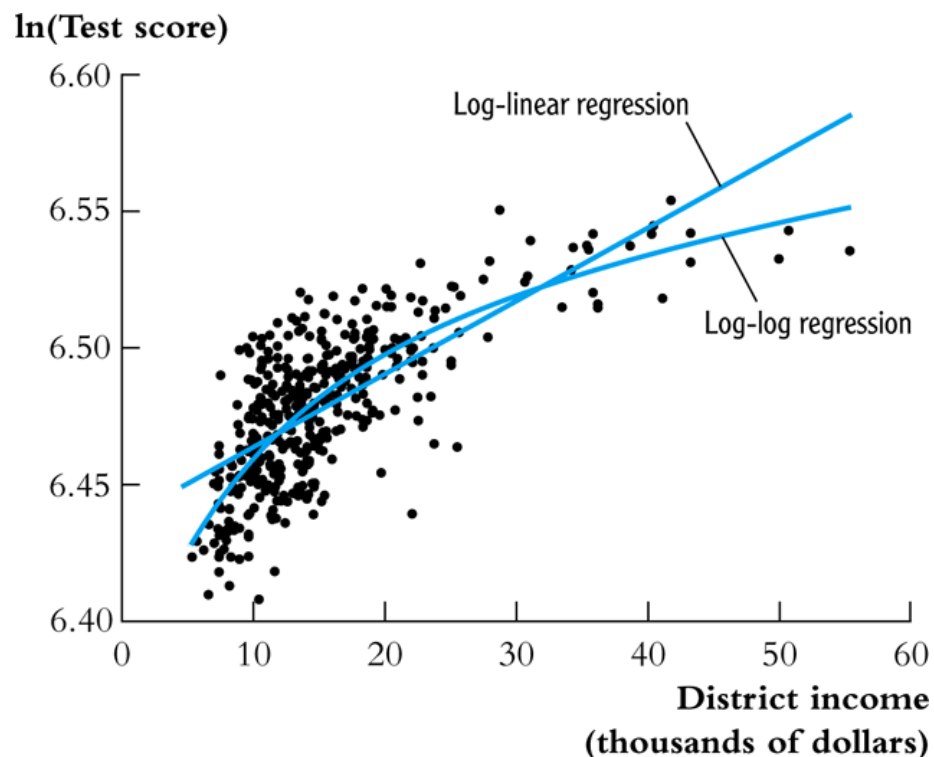
## *Example: $\ln(\text{TestScore})$ vs. $\ln(\text{Income})$*

$$\ln(\widehat{\text{TestScore}}) = 6.336 + 0.0554 \times \ln(\text{Income}_i) \\ (0.006) \quad (0.0021)$$

- For example, suppose income increases from \$10,000 to \$11,000, or by 10%. Then *TestScore* increases by approximately  $.0554 \times 10\% = .554\%$ . If *TestScore* = 650, this corresponds to an increase of  $.00554 \times 650 = 3.6$  points.
- How does this compare to the log-linear model?

# Functional form (SW 8.2)

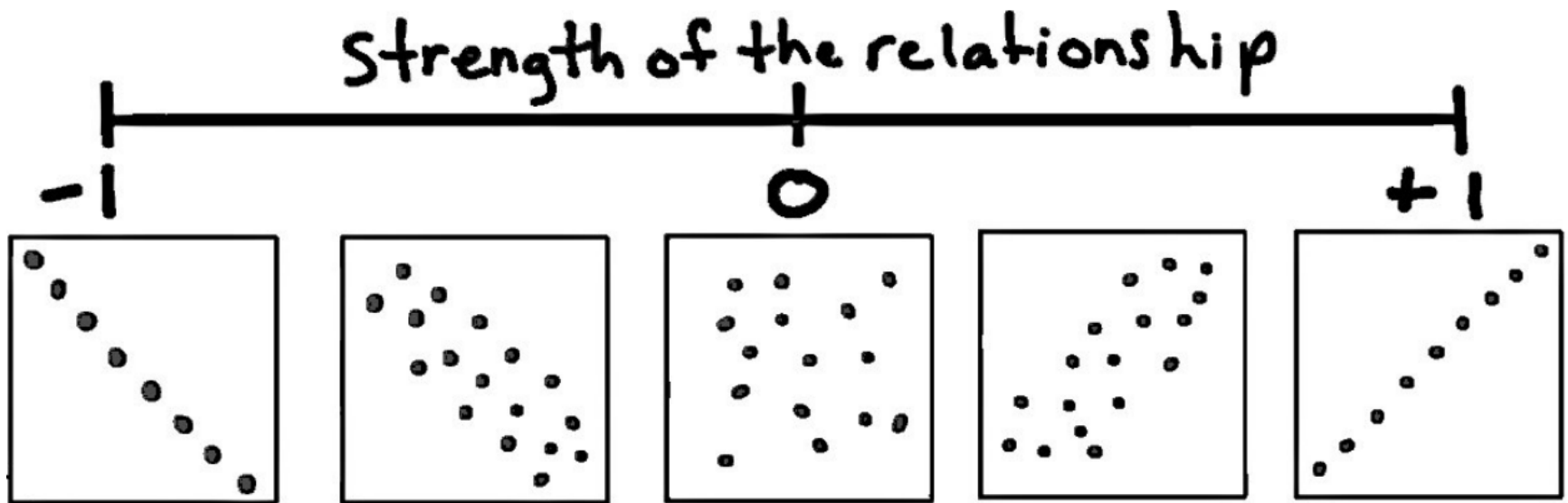
*Example:  $\ln(\text{TestScore})$  vs.  $\ln(\text{Income})$*



- Note vertical axis
- The log-linear model doesn't seem to fit as well as the log-log model, based on visual inspection.

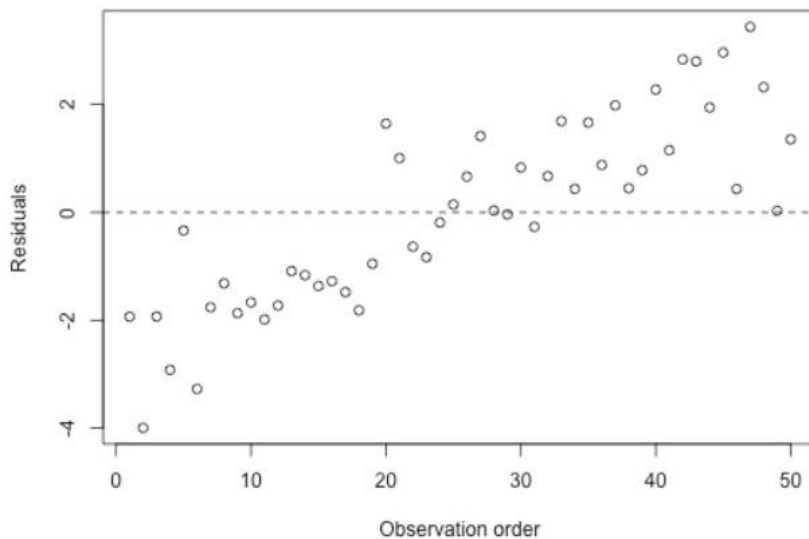
# Multicollinearity

Check pairwise correlations and scatterplots of the suspected independent variables

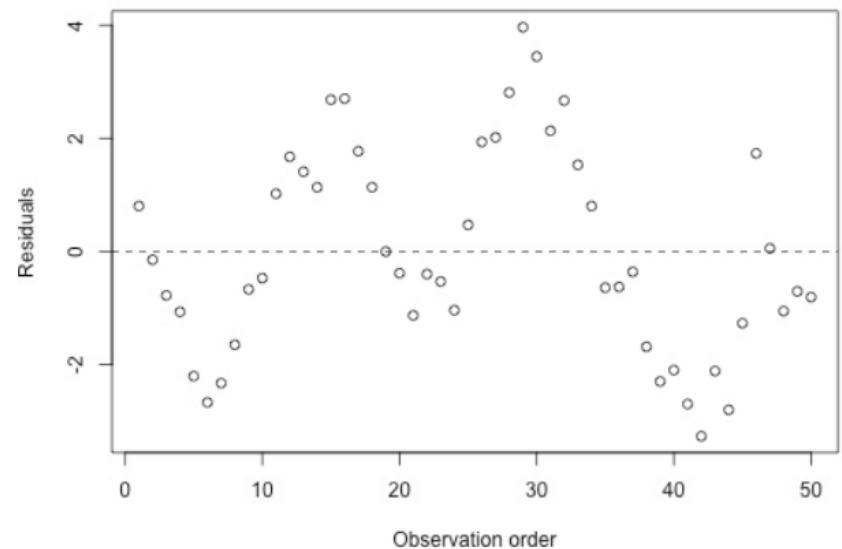


# Serial correlation

A residuals vs. order plot that exhibits (positive) trend suggests that some of the variation in the response is due to time

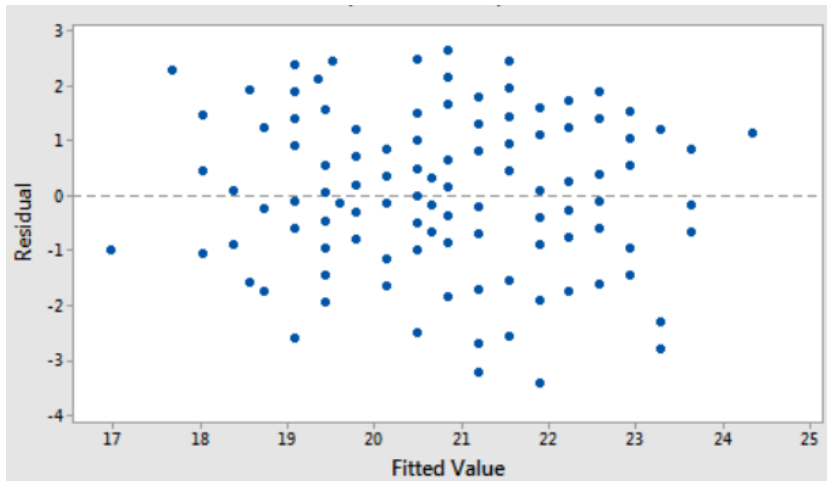


A residuals vs. order plot that suggests that there is "positive serial correlation" among the error terms. The plot suggests that the assumption of independent error terms is violated.



# Heteroskedasticity

A Good Residual Plot



Indications that Assumption of Constant Variance is Not Valid

