

## ARE LOG TRANSFORMED?

### Introduction

In this page, we will discuss how to interpret a regression model when some variables in the model have been log transformed. The example data can be downloaded [here \(https://stats.idre.ucla.edu/wp-content/uploads/2016/02/lqtrans.csv\)](https://stats.idre.ucla.edu/wp-content/uploads/2016/02/lqtrans.csv) (the file is in .csv format). The variables in the data set are writing, reading, and math scores ( **write**, **read** and **math**), the log transformed writing (**lgwrite**) and log transformed math scores (**lgmath**) and **female**. For these examples, we have taken the natural log (ln). All the examples are done in Stata, but they can be easily generated in any statistical package. In the examples below, the variable **write** or its log transformed version will be used as the outcome variable. The examples are used for illustrative purposes and are not intended to make substantive sense. Here is a table of different types of means for variable **write**.

Variable	Type	Obs	Mean	[95% Conf. Interval]	
write	Arithmetic	200	52.775	51.45332	54.09668
	Geometric	200	51.8496	50.46854	53.26845
	Harmonic	200	50.84403	49.40262	52.37208

### Outcome variable is log transformed

Very often, a linear relationship is hypothesized between a log transformed outcome variable and a group of predictor variables. Written mathematically, the relationship follows the equation

$$\log(y_i) = \beta_0 + \beta_1 x_{1i} + \cdots + \beta_k x_{ki} + e_i,$$

where  $y$  is the outcome variable and  $x_1, \dots, x_k$  are the predictor variables. In other words, we assume that  $\log(\mathbf{y}) - \mathbf{X}^T \boldsymbol{\beta}$  is normally distributed, (or  $\mathbf{y}$  is log-normal conditional on all the covariates). Since this is just an ordinary least squares regression, we can easily interpret a regression coefficient, say  $\beta_1$ , as the expected change in log of  $y$  with respect to a one-unit increase in  $x_1$  holding all other variables at any fixed value, assuming that  $x_1$  enters the model only as a main effect. But what if we want to know what happens to the outcome variable  $y$  itself for a one-unit increase in  $x_1$ ? The natural way to do this is to interpret the exponentiated regression coefficients,  $\exp(\beta)$ , since exponentiation is the inverse of logarithm function.

Let's start with the intercept-only model.

lgwrite	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
intercept	3.948347	.0136905	288.40	0.000	3.92135	3.975344

$$\log(\mathbf{write}) = \beta_0 = 3.95.$$

We can say that 3.95 is the unconditional expected mean of log of **write**. Therefore the exponentiated value is  $\exp(3.948347) = 51.85$ . This is the geometric mean of **write**. The emphasis here is that it is the geometric mean instead of the arithmetic mean. OLS regression of the original variable  $y$  is used to estimate the expected arithmetic mean and OLS regression of the log transformed outcome variable is to estimate the expected geometric mean of the original variable.

Now let's move on to a model with a single binary predictor variable.

lgwrite	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
female	.1032614	.0265669	3.89	0.000	.050871	.1556518
intercept	3.89207	.0196128	198.45	0.000	3.853393	3.930747

$$\begin{aligned}\log(\mathbf{write}) &= \beta_0 + \beta_1 \times \mathbf{female} \\ &= 3.89 + .10 \times \mathbf{female}.\end{aligned}$$

Before diving into the interpretation of these parameters, let's get the means of our dependent variable, **write**, by gender.

males						
Variable	Type	Obs	Mean	[95% Conf. Interval]		
write	Arithmetic	91	50.12088	47.97473	52.26703	
	Geometric	91	49.01222	46.8497	51.27457	
	Harmonic	91	47.85388	45.6903	50.23255	
females						
Variable	Type	Obs	Mean	[95% Conf. Interval]		
write	Arithmetic	109	54.99083	53.44658	56.53507	
	Geometric	109	54.34383	52.73513	56.0016	
	Harmonic	109	53.64236	51.96389	55.43289	

Now we can map the parameter estimates to the geometric means for the two groups. The intercept of 3.89 is the log of geometric mean of **write** when **female** = 0, i.e., for males. Therefore, the exponentiated value of it is the geometric mean for the male group:  $\exp(3.892) = 49.01$ . What can we say about the coefficient for **female**? In the log scale, it is the difference in the expected geometric means of the log of **write** between the female students and male students. In the original scale of the variable **write**, it is the ratio of the geometric mean of **write** for female students over the geometric mean of **write** for male students,  $\exp(.1032614) = 54.34383/49.01222 = 1.11$ . In terms of percent change, we can say that switching from male students to female students, we expect to see about 11% increase in the geometric mean of writing scores.

Last, let's look at a model with multiple predictor variables.

lgwrite	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
female	.114718	.0195341	5.87	0.000	.076194	.153242
read	.0066305	.0012689	5.23	0.000	.0041281	.0091329
math	.0076792	.0013873	5.54	0.000	.0049432	.0104152
intercept	3.135243	.0598109	52.42	0.000	3.017287	3.253198

$$\begin{aligned}\log(\text{write}) &= \beta_0 + \beta_1 \times \text{female} + \beta_2 \times \text{read} + \beta_3 \times \text{math} \\ &= 3.135 + .115 \times \text{female} + .0066 \times \text{read} + .0077 \times \text{math}.\end{aligned}$$

The exponentiated coefficient  $\exp(\beta_1)$  for **female** is the ratio of the expected geometric mean for the female students group over the expected geometric mean for the male students group, when **read** and **math** are held at some fixed value. Of course, the expected geometric means for the male and female students group will be different for different values of **read** and **math**. However, their ratio is a constant:  $\exp(\beta_1)$ . In our example,  $\exp(\beta_1) = \exp(.114718) \approx 1.12$ . We can say that writing scores will be 12% higher for the female students than for the male students. For the variable **read**, we can say that for a one-unit increase in **read**, we expect to see about a 0.7% increase in writing score, since  $\exp(.0066305) = 1.006653 \approx 1.007$ . For a ten-unit increase in **read**, we expect to see about a 6.9% increase in writing score, since  $\exp(.0066305 \times 10) = 1.0685526 \approx 1.069$ .

The intercept becomes less interesting when the predictor variables are not centered and are continuous. In this particular model, the intercept is the expected mean for  $\log(\text{write})$  for male (**female** = 0) when **read** and **math** are equal to zero.

In summary, when the outcome variable is log transformed, it is natural to interpret the exponentiated regression coefficients. These values correspond to changes in the ratio of the expected geometric means of the original outcome variable.

## Some (not all) predictor variables are log transformed

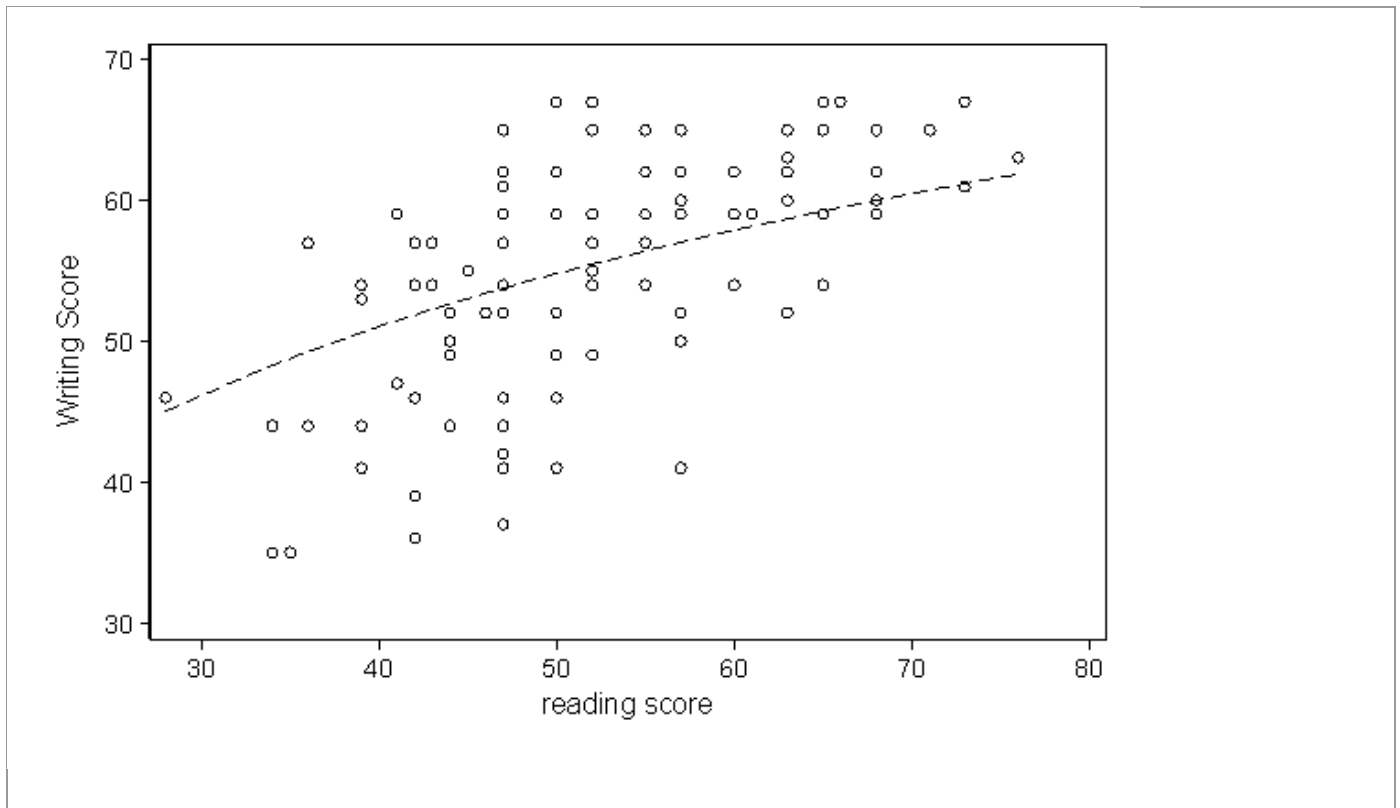
Occasionally, we also have some predictor variables being log transformed. In this section, we will take a look at an example where some predictor variables are log-transformed, but the outcome variable is in its original scale.

write	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
female	5.388777	.9307948	5.79	0.000	3.553118	7.224436
lgmath	20.94097	3.430907	6.10	0.000	14.17473	27.7072
lgread	16.85218	3.063376	5.50	0.000	10.81076	22.89359
intercept	-99.16397	10.80406	-9.18	0.000	-120.4711	-77.85685

Written in equation, we have

$$\begin{aligned}\text{write} &= \beta_0 + \beta_1 \times \text{female} + \beta_2 \times \log(\text{math}) + \beta_3 \times \log(\text{read}) \\ &= -99.164 + 5.389 \times \text{female} + 20.941 \times \log(\text{math}) + 16.852 \times \log(\text{read}).\end{aligned}$$

Since this is an OLS regression, the interpretation of the regression coefficients for the non-transformed variables are unchanged from an OLS regression without any transformed variables. For example, the expected mean difference in writing scores between the female and male students is about 5.4 points, holding the other predictor variables constant. On the other hand, due to the log transformation, the estimated effects of **math** and **read** are no longer linear, even though the effect of  $\log(\mathbf{math})$  and  $\log(\mathbf{read})$  are linear. The plot below shows the curve of predicted values against the reading scores for the female students group holding math score constant.



How do we interpret the coefficient of 16.852 for the variable of log of reading score? Let's take two values of reading score,  $r_1$  and  $r_2$ . The expected mean difference in writing score at  $r_1$  and  $r_2$ , holding the other predictor variables constant, is  $\mathbf{write}(r_2) - \mathbf{write}(r_1) = \beta_3 \times [\log(r_2) - \log(r_1)] = \beta_3 \times [\log(r_2/r_1)]$ . This means that as long as the percent increase in **read** (the predictor variable) is fixed, we will see the same difference in writing score, regardless where the baseline reading score is. For example, we can say that for a 10% increase in reading score, the difference in the expected mean writing scores will be always  $\beta_3 \times \log(1.10) = 16.85218 \times \log(1.1) \approx 1.61$ .

### Note:

Recalling the Taylor expansion of the function  $f(x) = \log(1+x)$  around  $x_0 = 0$ , we have  $\log(1+x) = x + \mathcal{O}(x^2)$ . Therefore, for a small change in the predictor variable we can approximate the difference in the expected mean of the dependent variable by multiplying the coefficient by the change in the

predictor variable. In our example we can say that for a 1% increase in reading score, the difference in the expected mean writing scores will be approximately  $\beta_3 \times 0.01 = 16.85218 \times 0.01 = .1685218$ . If we use the log, the exact value will be  $\beta_3 \times \log(1.01) = 16.85218 \times \log(1.01) = .1676848$ .

## Both the outcome variable and some predictor variables are log transformed

What happens when both the outcome variable and predictor variables are log transformed? We can combine the two previously described situations into one. Here is an example of such a model.

lgwrite	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
female	.1142399	.0194712	5.87	0.000	.07584	.1526399
lgmath	.4085369	.0720791	5.67	0.000	.2663866	.5506872
read	.0066086	.0012561	5.26	0.000	.0041313	.0090859
intercept	1.928101	.2469391	7.81	0.000	1.441102	2.415099

Written as an equation, we can describe the model:

$$\begin{aligned}\log(\mathbf{write}) &= \beta_0 + \beta_1 \times \mathbf{female} + \beta_2 \times \log(\mathbf{math}) + \beta_3 \times \mathbf{read} \\ &= 1.928101 + .1142399 \times \mathbf{female} + .4085369 \times \log(\mathbf{math}) + .0066086 \times \mathbf{read}.\end{aligned}$$

For variables that are not transformed, such as **female**, its exponentiated coefficient is the ratio of the geometric mean for the female to the geometric mean for the male students group. For example, in our example, we can say that the expected percent increase in geometric mean from male student group to female student group is about 12% holding other variables constant, since  $\exp(.1142399) \approx 1.12$ . For reading score, we can say that for a one-unit increase in reading score, we expected to see about 0.7% of increase in the geometric mean of writing score, since  $\exp(.0066086) = 1.007$ .

Now, let's focus on the effect of **math**. Take two values of **math**,  $m_1$  and  $m_2$ , and hold the other predictor variables at any fixed value. The equation above yields

$$\log(\mathbf{write}(m_2)) - \log(\mathbf{write}(m_1)) = \beta_2 \times [\log(m_2) - \log(m_1)]$$

It can be simplified to  $\log[\mathbf{write}(m_2)/\mathbf{write}(m_1)] = \beta_2 \times [\log(m_2/m_1)]$ , leading to

$$\frac{\mathbf{write}(m_2)}{\mathbf{write}(m_1)} = \left(\frac{m_2}{m_1}\right)^{\beta_2}$$

This tells us that as long as the ratio of the two math scores,  $m_2/m_1$  stays the same, the expected ratio of the outcome variable, **write**, stays the same. For example, we can say that for any 10% increase in **math** score, the expected ratio of the writing score will be  $(1.10)^{\beta_2} = (1.10)^{.4085369} = 1.0397057$ . In other words, we expect about 4% increase in writing score when math score increases by 10%.

### Note:

Here also we can use an approximation method. Since,  $(1 + x)^a \approx 1 + ax$  for a small value of  $|a|x$ , therefore for a small change in the predictor variable we can approximate the expected ratio of the of the dependent variable by multiplying the coefficient by the ratio of the change in the predictor variable. For example, we can say that for any 1% increase in **math** score, the expected ratio of the writing score is approximately

$1 + .01 \times \beta_2 = 1 + .01 \times .4085369 = 1.004085$ . The exact value will be  $(1.01)^{\beta_2} = (1.01)^{.4085369} = 1.004073$ .

[Click here to report an error on this page or leave a comment](#)

[How to cite this page \(https://stats.idre.ucla.edu/other/mult-pkg/faq/general/faq-how-do-i-cite-web-pages-and-programs-from-the-ucla-statistical-consulting-group/\)](https://stats.idre.ucla.edu/other/mult-pkg/faq/general/faq-how-do-i-cite-web-pages-and-programs-from-the-ucla-statistical-consulting-group/)