

## Linearity in Parameters

The OLS estimation technique requires that our model be *linear in parameters*. What this means is, each  $\beta$  term must appear essentially as a constant: we cannot have  $\beta_1^2$  or  $\log(\beta_1)$  or  $\beta_1\beta_2$ , for instance. This is because the OLS technique is only able to solve explicitly for each  $\beta$  if they appear in a linear fashion.

However, this does not necessitate that the model be linear in *variables*. There is no reason why can't specify a model of the form

$$y = \beta_1 + \beta_2 \log(x) + u$$

if we think it is useful to do so. And it certainly might be useful to do so. Consider the relationship between healthcare expenditure and life expectancy. We would expect more healthcare expenditure to be correlated with higher life expectancy, but at a diminishing rate (since there is a natural limit to life expectancy that medical treatment cannot overcome). So it wouldn't make sense to impose a linear relationship between healthcare expenditure and life expectancy; we would want to use a log to capture diminishing returns to healthcare expenditure.

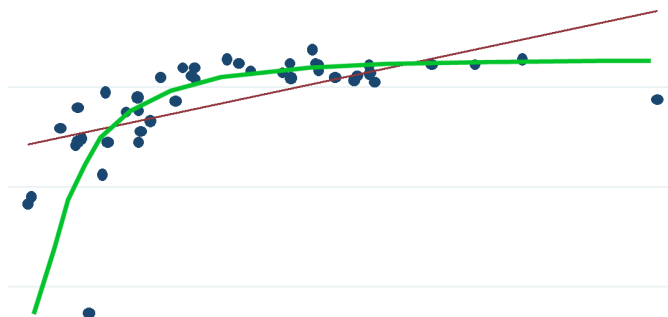


FIGURE 1: Specifying a logarithmic relationship (green) generates a better fit of the data compared to a linear relationship (red).

## Functional Forms

There are an infinite number of different ways we could specify such a model. I will focus on three due to their salient economic interpretations.

## Linear-Log Regression

A **linear-log** regression is of the form

$$y = \beta_1 + \beta_2 \log(x) + u. \quad (1)$$

It is named as such because the dependent variable is linear (it's simply  $y$ ) and the regressor variable is a logarithm.

We are not really interested in finding out how a change in  $\log(x)$  will affect  $y$ , however; we are interested in how a change in  $x$  will affect  $y$ . We'll have to do a little work to squeeze that information out of the regression, but it's not so bad. First, we can take the derivative of both sides with respect to  $x$ , which yields

$$\frac{dy}{dx} = \frac{\beta_2}{x}.$$

Now multiply both sides by  $dx$ , and then multiply the right-hand side by  $100/100$ . We end up with

$$dy = \frac{\beta_2}{100} \left( \frac{dx}{x} \times 100 \right).$$

Notice that the term in parentheses is the percentage change in  $x$ . So the interpretation in words is: the change in the level of  $y$  equals  $\beta_2/100$  times the percentage change in  $x$ ,

$$\Delta y = \frac{\beta_2}{100} \times \% \Delta x. \quad (2)$$

## Log-Linear Regression (Semi-Elasticity)

A **log-linear** regression is of the form

$$\log(y) = \beta_1 + \beta_2 x + u. \quad (3)$$

To interpret, start by taking the derivative of both sides with respect to  $x$ , which yields

$$\frac{d \log(y)}{dx} = \beta_2.$$

In order to squeeze  $y$  out of the left-hand side, we will appeal to the chain rule of calculus. In particular, we can write

$$\frac{d \log(y)}{dy} \frac{dy}{dx} = \beta_2.$$

We know that  $d\log(y)/dy = 1/y$ , so let's make that substitution. Let's also multiply both sides by  $100 \times dx$ , which yields

$$\frac{dy}{y} \times 100 = (100 \times \beta_2)dx.$$

In words: the percentage change in  $y$  equals  $100 \times \beta_2$  times the change in the level of  $x$ ,

$$\% \Delta y = 100 \beta_2 \times \Delta x. \quad (4)$$

In this form, coefficient  $\beta_2$  is referred to as the **semi-elasticity** of  $y$  with respect to  $x$ .

## Log-Log Regression (Elasticity)

A **log-log** regression is of the form

$$\log(y) = \beta_1 + \beta_2 \log(x) + u. \quad (5)$$

To interpret, take the derivative of both sides with respect to  $x$ , which gives

$$\frac{d\log(y)}{dx} = \frac{\beta_2}{x}.$$

Use the chain rule again on the right-hand side so that

$$\frac{d\log(y)}{dy} \frac{dy}{dx} = \frac{\beta_2}{x}.$$

We know that  $d\log(y)/dy = 1/y$ , so let's make that substitution. Also multiply both sides by  $dx$  and both sides by 100. Doing so yields

$$\frac{dy}{y} \times 100 = \beta_2 \left( \frac{dx}{x} \times 100 \right).$$

In words: the percentage change in  $y$  is equal to  $\beta_2$  times the percentage change in  $x$ ,

$$\% \Delta y = \beta_2 \times \% \Delta x. \quad (6)$$

In this form, coefficient  $\beta_2$  is referred to as the **elasticity** of  $y$  with respect to  $x$ , which you hopefully remember from a microeconomics course.

## Summary

Again, there are a multitude of other functional forms we could consider, e.g. quadratic forms, that are useful in certain contexts. Those will be discussed later as they become pertinent. But for now, here is a table of the four functional forms we have seen thus far.

Model	Dependent Variable	Regressor	Interpretation of $\beta_2$
linear	$y$	$x$	$\Delta y = \beta_2 \times \Delta x$
linear-log	$y$	$\log(x)$	$\Delta y = \frac{\beta_2}{100} \times \% \Delta x$
log-linear (semi-elasticity)	$\log(y)$	$x$	$\% \Delta y = 100 \beta_2 \times \Delta x$
log-log (elasticity)	$\log(y)$	$\log(x)$	$\% \Delta y = \beta_2 \times \% \Delta x$

## Example: Life Expectancy and Healthcare Expenditure

Download `hcle.csv` from my website and import the data into R. The dataset has three variables: country, life expectancy at birth, and healthcare spending per-capita in 2015. If we estimate a linear regression of the form

$$\text{lifeexpect} = \beta_1 + \beta_2 \text{hcspending} + u,$$

then we find goodness-of-fit measure  $R^2 = 0.363$ . On the other hand, if we do a linear-log regression like suggested earlier,

$$\text{lifeexpect} = \beta_1 + \beta_2 \log(\text{hcspending}) + u,$$

then we find goodness-of-fit measure  $R^2 = 0.542$ , implying a better fit.

Note that there is an **important caveat** in comparing the  $R^2$  of different models. In particular, it is *not* meaningful to compare the  $R^2$  of models that have different dependent variables! If the regressors are different but the dependent variables are the same, as is the case here, then it is fine to compare  $R^2$ .

Here is the R code I used to compare the two regressions. Notice that I had to generate `log(hcspending)` on line 11 because the data originally comes in levels only.

```
1 library("rio")
2 library("stargazer")
3
4 hcle <- import("hcle.csv")
5
6 ### linear model
7 linreg <- lm(hcle$lifeexpect ~ hcle$hcspending)
8 stargazer(linreg, type = "text")
9
10 ### linear-log model
11 hcle$loghcsp = log(hcle$hcspending)
12 linlogreg <- lm(hcle$lifeexpect ~ hcle$loghcsp)
13 stargazer(linlogreg, type = "text")
```