

Functional Form in the Linear Model

1	Introduction	2
2	Examples	3
2.1	Polynomial	3
2.2	Inverse	4
2.3	Explanatory Variable in Logs (linear-log)	5
2.4	Dependent Variable in Logs (log-linear)	6
2.5	Dependent and Explanatory Variables in Logs (log-log)	7
2.6	Dummy Variables	8
2.7	Interaction Terms	9
2.8	Interactions with Dummy Variables	10
2.9	Spline Functions	11
3	Implementation in Stata 17	12
4	Implementation in R 4.3.1	13

1 Introduction

Despite its name, the classical *linear* regression model, is not limited to a linear relationship between the dependent and the explanatory variables.

Consider a vector $x'_i = (x_{i1} \ x_{i2} \ \dots \ x_{iK})$ of K variables for each observation i . The L functions $f_1(x_i), f_2(x_i), \dots, f_L(x_i)$ map the K -dimensional vector x_i into L scalars $z_{i1}, z_{i2}, \dots, z_{iL}$. The function $g(y_i)$ is a univariate function of the dependent variable. The non-linear econometric model

$$g(y_i) = \beta_0 + \beta_1 f_1(x_i) + \beta_2 f_2(x_i) + \dots + \beta_L f_L(x_i) + u_i$$

can therefore be written as

$$\begin{aligned} g(y_i) &= \beta_0 + \beta_1 z_{i1} + \beta_2 z_{i2} + \dots + \beta_L z_{iL} + u_i \\ &= z'_i \beta + u_i . \end{aligned}$$

The latter is the usual multiple linear regression model with $L + 1$ regressors as long as all necessary assumptions about the error term and the L *transformed* explanatory variables $z_{i1}, z_{i2}, \dots, z_{iL}$ are satisfied. All properties of OLS are therefore preserved.

Note: While the original model is potentially *non-linear in the variables* x_{ik} , it is *linear in the parameters* β . Also note that the error term u_i is additive.

2 Examples

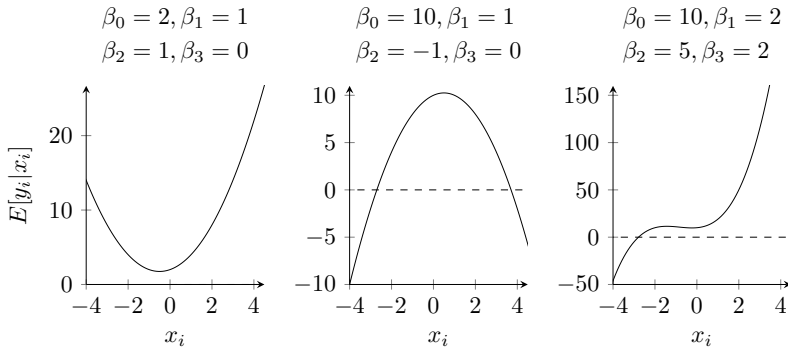
2.1 Polynomial

Functional form (3rd order):

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + u_i$$

Expected value under *OLS3c*:

$$E[y_i|x_i] = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3$$



Marginal effect on the expected dependent variable:

$$\frac{\partial E[y_i|x_i]}{\partial x_i} = \beta_1 + 2\beta_2 x_i + 3\beta_3 x_i^2.$$

Marginal effect on an individual given its error term u_i :

$$\frac{\partial y_i}{\partial x_i} = \beta_1 + 2\beta_2 x_i + 3\beta_3 x_i^2.$$

Note that the marginal effect depends on the value of the explanatory variable x_i . The individual parameters β_1 , β_2 and β_3 often have no direct interpretation.

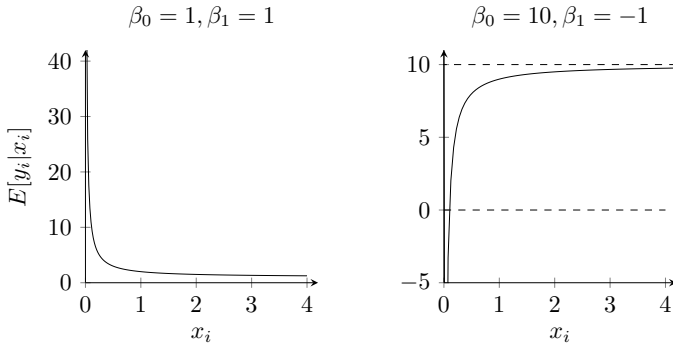
2.2 Inverse

Functional form:

$$y_i = \beta_0 + \beta_1 \frac{1}{x_i} + u_i.$$

Expected value under *OLS*:

$$E[y_i|x_i] = \beta_0 + \beta_1 \frac{1}{x_i}$$



Marginal effect:

$$\frac{\partial E[y_i|x_i]}{\partial x_i} = -\frac{\beta_1}{x_i^2}$$

and

$$\frac{\partial y_i}{\partial x_i} = -\frac{\beta_1}{x_i^2}.$$

Note that a positive sign of β_1 means a negative relationship and vice-versa.

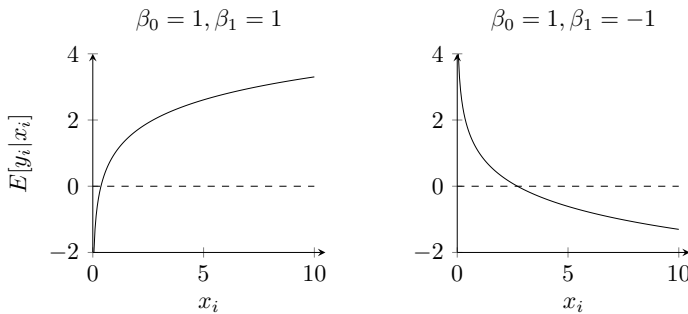
2.3 Explanatory Variable in Logs (linear-log)

Functional form:

$$y_i = \beta_0 + \beta_1 \ln(x_i) + u_i.$$

Expected value under *OLS*³:

$$E[y_i|x_i] = \beta_0 + \beta_1 \ln(x_i)$$



Marginal effect:

$$\frac{\partial E[y_i|x_i]}{\partial x_i} = \frac{\beta_1}{x_i}, \quad \frac{\partial y_i}{\partial x_i} = \frac{\beta_1}{x_i}.$$

The coefficient β_1 is approximately the effect of a 100 % change in the explanatory variable x_i on the dependent variable y_i

$$\beta_1 \approx \frac{\Delta E[y_i|x_i]}{\Delta x_i/x_i}, \quad \beta_1 \approx \frac{\Delta y_i}{\Delta x_i/x_i}$$

where Δ is a small discrete change. Hence β_1 divided by 100 is approximately the effect of a 1 % change in x_i . β_1 is also called a semi-elasticity. For example, $\beta_1 = 56$ means that an increase in x_i by 1 % leads to an increase in the dependent variable y_i by (approximately) 0.56 units.¹

¹The exact effect of a discrete change of the explanatory variable by Δx_i is

$$\Delta E[y_i|x_i] = \beta_1 \cdot \ln \left(1 + \frac{\Delta x_i}{x_i} \right) \quad \text{and} \quad \Delta y_i = \beta_1 \cdot \ln \left(1 + \frac{\Delta x_i}{x_i} \right).$$

For example, $\beta_1 = 56$ means that an increase in x_i by $0.01 = 1\%$ leads to an increase in the dependent variable y_i by exactly $56 \cdot \ln(1 + 0.01) = 56 \cdot 0.00995 = 0.557$ units.

2.4 Dependent Variable in Logs (log-linear)

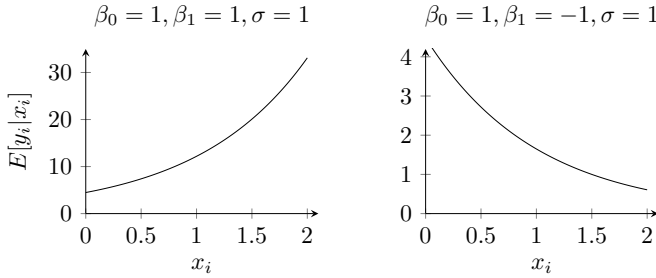
Functional form:

$$\ln(y_i) = \beta_0 + \beta_1 x_i + u_i.$$

Expected value:

$$E[y_i|x_i] = e^{\beta_0 + \beta_1 x_i} \cdot \alpha_i.$$

where $\alpha_i = E[e^{u_i}|x_i] > 1$. If the error is independent of the explanatory variables (*OLSS3b*), $\alpha_i = \alpha$ is constant. If the error is normally distributed with homoscedastic error $V[u_i|x_i] = \sigma^2$ (*OLSS3a*, *OLSS4a*), $\alpha_i = \alpha = e^{\sigma^2/2}$.



Marginal effect:

$$\frac{\partial E[y_i|x_i]}{\partial x_i} = E[y_i|x_i] \cdot \beta_1, \quad \frac{\partial y_i}{\partial x_i} = y_i \cdot \beta_1.$$

The coefficient $\beta_1 \cdot 100\%$ is approximately the percentage effect on the dependent variable y_i of a change in the variable x_i by one unit

$$\beta_1 \approx \frac{\frac{\Delta E[y_i|x_i]}{E[y_i|x_i]}}{\frac{\Delta x_i}{x_i}}, \quad \beta_1 \approx \frac{\frac{\Delta y_i}{y_i}}{\frac{\Delta x_i}{x_i}}$$

where Δ is a small discrete change. β_1 is also called a semi-elasticity. For example, $\beta_1 = 0.06$ means that an increase in x_i by one unit leads to an increase in the dependent variable y_i by (approximately) 6%.²

²The exact effect of a discrete change of the explanatory variable by Δx_i units is

$$\frac{\Delta E[y_i|x_i]}{E[y_i|x_i]} = e^{\beta_1 \Delta x_i} - 1 \quad \text{and} \quad \frac{\Delta y_i}{y_i} = e^{\beta_1 \Delta x_i} - 1.$$

For example, $\beta_1 = 0.06$ means that an increase in x_i by one unit leads to an increase in the dependent variable y_i by exactly $\exp(0.06) - 1 = 0.0618 = 6.18\%$.

2.5 Dependent and Explanatory Variables in Logs (log-log)

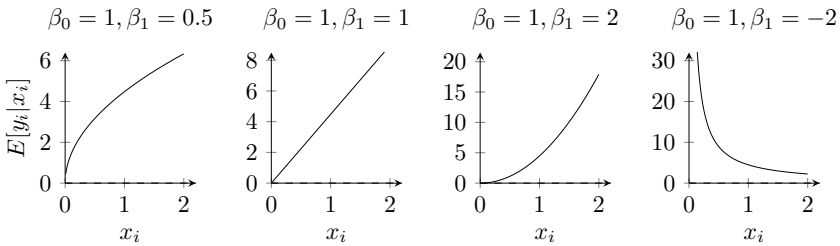
Functional form:

$$\ln(y_i) = \beta_0 + \beta_1 \ln(x_i) + u_i.$$

Expected value:

$$E[y_i|x_i] = e^{\beta_0 + \beta_1 \ln(x_i)} \cdot \alpha_i.$$

where $\alpha_i = E[e^{u_i}|x_i] > 1$. Assuming that the error is independent of the explanatory variables (*OLS3b*), $\alpha_i = \alpha$ is a constant. Assuming that the error is normally distributed with homoscedastic error $V[u_i|x_i] = \sigma^2$ (*OLS3a* and *OLS4a*), $\alpha_i = \alpha = e^{\sigma^2/2}$. For $\sigma = 1$, hence $\alpha = 1.65$:



Marginal effect:

$$\frac{\partial E[y_i|x_i]}{\partial x_i} = E[y_i|x_i] \cdot \frac{\beta_1}{x_i}, \quad \frac{\partial y_i}{\partial x_i} = y_i \cdot \frac{\beta_1}{x_i}.$$

The coefficient β_1 is approximately the percentage effect on the dependent variable y_i of a 1 % change in the explanatory variable x_i

$$\beta_1 \approx \frac{\frac{\Delta E[y_i|x_i]}{E[y_i|x_i]}}{\frac{\Delta x_i}{x_i}}, \quad \beta_1 \approx \frac{\frac{\Delta y_i}{y_i}}{\frac{\Delta x_i}{x_i}},$$

where Δ is a small discrete change. β_1 is also called an elasticity. For example, $\beta_1 = 0.8$ means that an increase in x_i by 1 % leads to an increase in the dependent variable y_i by (approximately) 0.8 %.³

³The exact effect of a discrete change of the explanatory variable by Δx_i is

$$\frac{\Delta E[y_i|x_i]}{E[y_i|x_i]} = e^{\beta_1 \ln(1 + \frac{\Delta x_i}{x_i})} - 1 \quad \text{and} \quad \frac{\Delta y_i}{y_i} = e^{\beta_1 \ln(1 + \frac{\Delta x_i}{x_i})} - 1.$$

For example, $\beta_1 = 0.8$ means that an increase in x_i by 0.01 = 1 % leads to an increase in the dependent variable y_i by exactly $\exp(0.8 \cdot \ln(1 + 0.01)) - 1 = 0.00799 = 0.799\%$.

2.6 Dummy Variables

Functional form:

$$y_i = \beta_0 + \beta_1 d_i + u_i$$

where $d_i \in \{0, 1\}$ is a dummy variable that either takes value 0 or 1.

Expected value under *OLS3c*:

$$E[y_i | d_i] = \begin{cases} \beta_0 & \text{if } d_i = 0 \\ \beta_0 + \beta_1 & \text{if } d_i = 1 \end{cases}$$

Marginal effect:

$$\frac{\partial E[y_i | d_i]}{\partial d_i} = \beta_1 \quad \text{and} \quad \frac{\partial y_i}{\partial d_i} = \beta_1.$$

Note that the notion of a marginal, i.e. infinitesimally small change, is not useful for a dummy variable d_i which can only increase by exactly one unit. The coefficient β_1 is better interpreted as the difference in the means of the (treatment) group with $d_i = 1$ and the (control) group with $d_i = 0$

$$\beta_1 = E[y_i | d_i = 1] - E[y_i | d_i = 0]$$

which constitutes the average treatment effect (ATE).

Note that the model

$$y_i = \beta_0 + \beta_1 \text{treat}_i + \beta_2 \text{control}_i + u_i$$

where the dummy variable treat_i takes value 1 for the treatment group and 0 for the control group while the dummy variable control_i takes the opposite values violates *OLS5* and the parameters β_0, β_1 and β_2 cannot be separately identified. This is called the dummy variable trap.

2.7 Interaction Terms

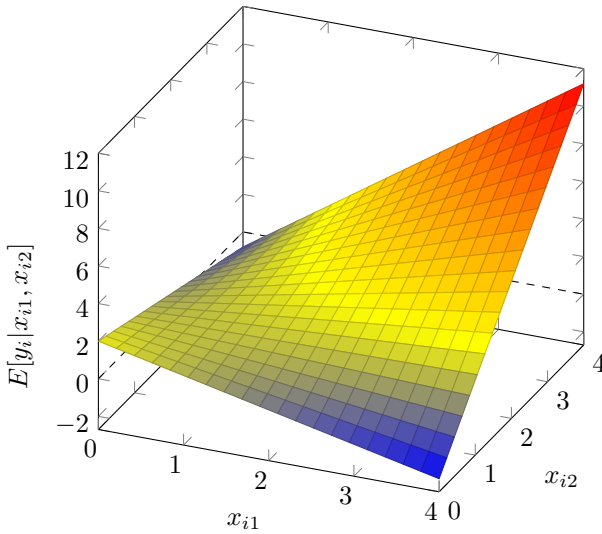
Functional form:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 (x_{i1} \cdot x_{i2}) + u_i$$

Expected value under *OLS3c*:

$$E[y_i | x_{i1}, x_{i2}] = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 (x_{i1} \cdot x_{i2})$$

$$\beta_0 = 2, \beta_1 = -1, \beta_2 = -0.7, \beta_3 = 1$$



Marginal effects:

$$\frac{\partial E[y_i | x_{i1}, x_{i2}]}{\partial x_{i1}} = \frac{\partial y_i}{\partial x_{i1}} = \beta_1 + \beta_3 x_{i2}$$

$$\frac{\partial E[y_i | x_{i1}, x_{i2}]}{\partial x_{i2}} = \frac{\partial y_i}{\partial x_{i2}} = \beta_2 + \beta_3 x_{i1}$$

2.8 Interactions with Dummy Variables

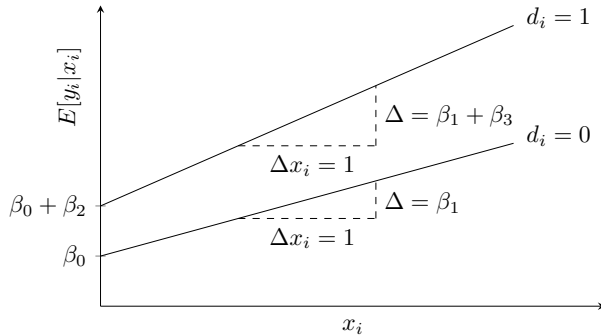
Functional form:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 d_i + \beta_3 (x_i \cdot d_i) + u_i$$

where $d_i \in \{0, 1\}$ is a dummy variable that either takes value 0 or 1 and x_i is a continuous explanatory variable.

Expected value under *OLS3c*:

$$E[y_i | x_i, d_i] = \begin{cases} \beta_0 + \beta_1 x_i & \text{if } d_i = 0 \\ (\beta_0 + \beta_2) + (\beta_1 + \beta_3) x_i & \text{if } d_i = 1 \end{cases}$$



Marginal effect:

$$\frac{\partial E[y_i | x_i, d_i]}{\partial x_{i1}} = \frac{\partial y_i}{\partial x_{i1}} = \begin{cases} \beta_1 & \text{if } d_i = 0 \\ \beta_1 + \beta_3 & \text{if } d_i = 1 \end{cases}$$

Note that the interaction of all variables (here constant and one explanatory variable x_i) with the dummy variable d_i estimates separate linear relationships for the two groups defined by d_i . This yields identical estimates as two separate regressions for both groups. The standard errors may differ between the interacted joint estimation and the separate regressions unless robust standard errors are computed.

2.9 Spline Functions

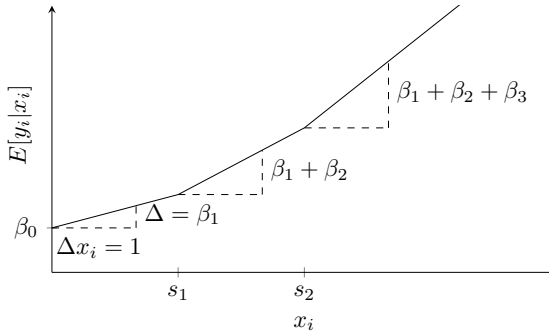
Functional form:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 d_{i1}(x_i - s_1) + \beta_3 d_{i2}(x_i - s_2) + u_i$$

where $d_{i1} = 1$ if $x_i \geq s_1$ and $d_{i2} = 1$ if $x_i \geq s_2$. s_1 and s_2 are known thresholds.

Expected value under *OLS3c*:

$$E[y_i|x_i] = \beta_0 + \beta_1 x_i + \beta_2 d_{i1}(x_i - s_1) + \beta_3 d_{i2}(x_i - s_2)$$



Marginal effect:

$$\frac{\partial E[y_i|x_i]}{\partial x_i} = \frac{\partial y_i}{\partial x_i} = \begin{cases} \beta_1 & \text{if } x_i < s_1 \\ \beta_1 + \beta_2 & \text{if } s_1 \leq x_i < s_2 \\ \beta_1 + \beta_2 + \beta_3 & \text{if } x_i \geq s_2 \end{cases}$$

3 Implementation in Stata 17

Non-linear Functional forms can be estimated with OLS by generating the transformed variables. For example,

```
webuse auto7.dta
generate mpg2 = mpg^2
reg price mpg mpg2
```

estimates a second order polynomial.

Dummy variables are easily created from categorical variables with the `xi` command. For example,

```
xi i.manufacturer
reg price mpg _Imanufactu_*
```

creates 22 dummy variables for the 23 categories in the variable `manufacturer_grp` (excluding the first one for use as reference category) and regresses `price` on `mpg` and all 22 dummy variables plus a constant. This can also be done in one step,

```
reg price mpg i.manufacturer_grp
```

Interactions with categorical variables can be directly can be formed as “factor variables”. For example,

```
reg price i.foreign i.foreign#c.mpg
```

estimates separate intercepts and separate slopes of `mpg` for domestic and for foreign cars. Alternatively,

```
reg price i.foreign##c.mpg
```

reports the difference of intercept and slopes for foreign cars compared to the reference group, i.e. domestic cars.

The variables used for spline functions are conveniently created with the `mkspline` command. For example,

```
mkspline mpg_1 20 mpg_2 25 mpg_3 = mpg
reg price mpg_*
```

regresses `price` on `mpg` using a piecewise linear function. Also consider the option `marginal`.

4 Implementation in R 4.3.1

Non-linear Functional forms can be estimated with OLS by specifying the Functional form in the estimated model. For example,

```
library(haven)
auto <- read_dta("http://www.stata-press.com/data/r17/auto7.dta")
lm(price~mpg+I(mpg^2), data=auto)
```

regresses `price` on a second order polynomial of `mpg`. Note that most mathematical functions need to be wrapped within the `I()` function.

Categorical variables are automatically included as a set of dummy variables if they are defined as factor variables. For example,

```
auto$manufacturer_grpf <- factor(auto$manufacturer_grp)
lm(price~mpg+manufacturer_grpf, data=auto)
```

regresses `price` on `mpg` and 22 dummy variables in `manufacturer_grp` (excluding the first category for use as reference group) plus a constant.

Interactions with categorical variables can be directly estimated when the categorical variable is defined as factor variable. For example,

```
auto$foreignf <- factor(auto$foreign, labels=c("domestic", "foreign"))
lm(price~foreignf+mpg:foreignf, data=auto)
```

estimates separate intercepts and separate slopes of `mpg` for domestic and for foreign cars. Alternatively,

```
lm(price~foreignf+mpg/foreignf, data=auto)
```

reports the difference of intercept and slopes for foreign cars compared to the reference group, i.e. domestic cars. In the case of interaction with a dummy variable, this is equivalent to either

```
lm(price~foreign+mpg+mpg:foreign, data=auto)
lm(price~foreign+mpg+I(mpg*foreign), data=auto)
```

Linear (and polynomial) spline functions are implemented in the `splines` package. See the help for details,

```
library("splines")
?splines
```

References

- Stock, James H. and Mark W. Watson (2020), Introduction to Econometrics, 4th Global ed., Pearson. Chapter 8.
- Wooldridge, Jeffrey M. (2009), Introductory Econometrics: A Modern Approach, 4th ed. South-Western. Section 6.2 and 6.4.
- Jaccard James and Robert Turrissi (2003), Interaction Effects in Multiple Regression, 2nd ed., Quantitative Applications in the Social Sciences 07-72, Sage.
- Kennedy, Peter (2003), A Guide to Econometrics, 5th ed., Blackwell Publishing, chapter 7.