Linear Regression With Special Variables

Junhui Qian

December 21, 2014

Outline

- Standardized Scores
- Quadratic Terms
- Interaction Terms
- Binary Explanatory Variables
- Binary Choice Models

Standardized Scores: An Example

 When we study how salaries depend on IQ score, we may run the following simple regression,

$$\log(salary_i) = \beta_0 + \beta_1 I Q_i + u.$$

- What is the economic interpretation of β_1 above?
- To make the coefficient more meaningful, we can define a new variable

$$z_i = \frac{IQ_i - IQ}{\hat{\sigma}_{IQ}}$$

and run

$$\log(salary_i) = \beta_0 + \beta_1 z_i + u.$$

• Now, what is the economic interpretation of β_1 ?

Standardized Scores

Given a variable x_i, the standardized score is defined as

$$z_{x,i}=\frac{x_i-\bar{x}}{\hat{\sigma}_x},$$

where \bar{x} and $\hat{\sigma}_{x}$ are the sample mean and standard deviation, respectively.

The coefficient on a standardized score is interpreted as the increase in y when x is one standard deviation higher. Standardized Score as Dependent Variable

The standardized score can also be the dependent variable in a regression. For example, let

$$z_{y,i}=\frac{y_i-\bar{y}}{\hat{\sigma}_y},$$

and run

$$z_{\mathbf{y},i} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u.$$

β₁ in this model is interpreted as the increase in y, in terms of standard deviations, with each unit of increase in x₁, holding other factors fixed.

Standardized Coefficients

 We may need a model where all variables are standardized, for example,

$$z_{y,i} = \beta_1 z_{x1,i} + \beta_2 z_{x2,i} + u.$$

- The β₁ is interpreted as the increase in y, in terms of standard deviations, when x₁ is one standard deviation higher, holding other factors fixed.
- Note that we need no constant term in the above model.
- The β 's in such models are called "standardized coefficients".

Models with Quadratic Terms

 We have seen models with quadratic terms. For example, in the income determination model,

$$\log(income) = \beta_0 + \beta_1 edu + \beta_2 edu^2 + \beta_3 expr + u.$$

• How do you interpret the parameters β_1 and β_2 ?

The Level-Dependent Partial Effect

In a model with quadratic terms,

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + u.$$

The partial effect of x on y is dependent on the level of x.
To see why,

$$\frac{\partial y}{\partial x} = \beta_1 + 2\beta_2 x.$$

We estimate the model and obtain,

 $log(income) = 7.55 + 0.0715edu + 0.00560edu^2 - 0.00299expr.$

We may ask, for a person with 10 years' education, how much income increase would be expected if he receives one more year of education?

Interaction Terms

Recall that linear regression is the first-order approximation of a usually nonlinear relationship:

$$y = f(x_1, x_2) + u \approx \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u.$$

▶ To get better approximations, we may add quadratic terms,

$$y = f(x_1, x_2) + u \approx \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_2^2 + \beta_5 x_1 x_2 + u.$$

- β₃ and β₄ measures nonlinearities in x₁ and x₂, and β₅ measures the "interaction" effect between x₁ and x₂.
- The term x_1x_2 is hence called the "interaction term".

Interaction Effect

- There is "interaction effect" when a change in one explanatory variable may affect the slope of another.
- ► For example, in a model of housing price,

 $price = \beta_0 + \beta_1 sqft + \beta_2 rooms + \beta_3 sqft \cdot rooms + \beta_4 bath + u.$

- The partial effect of *rooms* on *price* is $\beta_2 + \beta_3 sqft$.
- If β₃ > 0, an additional room in a large house is more valuable than that in a small house.

Average Partial Effect

In a model with interaction terms,

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

- β_1 is the partial effect of x_1 on y when $x_2 = 0$.
- ► This is not very interesting. Instead, we may be interested in the partial effect of x₁ on y when x₂ = x
 ₂.
- ▶ We may obtain such an "average partial effect" by plug \bar{x}_2 in $\hat{\beta}_1 + \hat{\beta}_3 x_2$.
- To obtain average partial effect directly, we can run

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 (x_1 - \bar{x}_1) (x_2 - \bar{x}_2) + u.$$

• Now β_1 is the desired average partial effect.

Binary (Dummy) Variables

- A binary variable describes qualitative information, in contrast to quantitative information.
- For example, male or female, urban or rural residence, employed or unemployed, a person buys a car or not, etc.
- A typical binary variable is defined as

female = 0, 1,

where 0 corresponds to "not female" and 1 corresponds to "female".

Binary Variable on the Intercept

 Binary variables may influence the intercept only. For example, in the following model,

 $\log(income) = \beta_0 + \beta_1 edu + \beta_2 expr + \beta_3 female + u,$

where *female* is a binary variable.

For females, the model is

 $log(income) = (\beta_0 + \beta_3) + \beta_1 edu + \beta_2 expr + u.$

For males, the model is

$$\log(income) = \beta_0 + \beta_1 edu + \beta_2 expr + u.$$

 The group of males in this model can be called "baseline group".

Binary Variable on the Slope

 Binary variables may also influence the slope. For example, in the following model,

 $log(income) = \beta_0 + \beta_1 edu + \beta_2 expr + \beta_3 female \cdot edu + u.$

For females, the model is

$$log(income) = \beta_0 + (\beta_1 + \beta_3)edu + \beta_2expr + u.$$

For males, the model is

$$log(income) = \beta_0 + \beta_1 edu + \beta_2 expr + u.$$

When There Are Several Categories

- Binary variables are easily defined when there are two categories. For example, the gender of a person can only be male or female.
- Other qualitative information may involve more than two categories. For example, to describe the region of the country where people live and work, we may have three categories (coast, middle, west). For another example, to describe seasonality, we have four categories.

When There Are Several Categories

Obviously, one binary variable is not enough for more than two categories. The solution is to have more than one binary variables. For example, if there are 3 regions (coast, middle, west) in total, we may consider

 $log(income) = \beta_0 + \beta_1 edu + \beta_2 expr + \beta_3 west + \beta_4 middle + u,$

where west is defined as 1 if the individual is in the west and 0 otherwise.

Why not add one more binary variable coast?

Testing the Existence of Qualitative Effects

- We often want to test whether being in some category (such as gender, hukou, race, region, etc.) has any effect on dependent variables, controlling for other factors.
- With multiple linear regression with binary variables, we may easily achieve this.
- For example, to test whether or not females are discriminated in work places, we may run

$$log(income) = \beta_0 + \beta_1 edu + \beta_2 expr + \beta_3 female + u,$$

and test

$$H_0: \beta_3 = 0 \quad H_1: \beta_3 < 0.$$

Modeling Binary Choices

- Decision to migrate or not
- Buy a car or not

► E

- Rent or buy an apartment
- To vote for or against a bill
- ► For college/MBA graduates, find a job within 3 months
- For females, to be a housewife or not

Binary Distribution

Let y be the binary choice variable taking values of 0 and 1. We can describe y by binary (or Bernoulli) distribution with parameter p,

$$f(y) = \begin{cases} p & \text{if } y = 1\\ 1-p & \text{if } y = 0 \end{cases}$$

Equivalently, we may write

$$f(y) = p^{y}(1-p)^{1-y}, y = 0, 1.$$

We have

$$\mathbb{E}y = p \operatorname{var}(y) = p(1-p).$$

Let x be a vector of variables that influences the outcome of y, which takes value of 0 or 1. We may write a **linear probability model** as follows,

$$y=x'\beta+u.$$

- Let $p = \mathbb{E}(y|x) = x'\beta$. It is the probability of y = 1 given x.
- Conditional on x, y is distributed as Binary(p).
- $\operatorname{var}(y|x)$ is then $p(1-p) = x'\beta(1-x'\beta)$.

An Example

Suppose y denotes the decision of an individual on whether or not migrates to city. And let x be the income the individual receives from previous job in the countryside. We may construct a simple migration model,

$$y=\beta_0+\beta_1x+u.$$

Problem of Linear Probability Model

- Heteroscedasticity.
- Nonsense Probability
 - ▶ To see this, observe that after estimation,

$$\hat{y} = \hat{p} = x'\hat{\beta}.$$

- ► To make predictions based on the estimated model, it is not guaranteed that p̂ ∈ [0, 1].
- The nonsense probability comes from linear marginal probability,

$$\frac{\partial p}{\partial x} = \beta.$$

Probit and Logit Model

► One method to limit p within [0, 1] is to use the cumulative distribution function (cdf) of a distribution, say F(·),

$$p = F(x'\beta) = \mathbb{P}(Z \le x'\beta).$$

- ► If F is the cdf of the standard normal distribution (N(0,1)), then the model is called a "probit model".
- If F(s) = 1/(1+e^{-s}), then the model is called a "logit model".
 Note that in this case, F is the cdf of a logistic distribution.

Marginal Effects on Probability

 Probit and logit models are nonlinear models with nonlinear marginal effects on probability,

$$\frac{\partial p}{\partial x} = f(x'\beta)\beta,\tag{1}$$

where f is the pdf of the corresponding distribution.

- When |x'β| is large, the marginal probability is small. When x'β = 0, f reaches maximum and the marginal effects on probability would be the largest.
- If x contains a binary variable d, the marginal probability is given by

$$\mathbb{P}(y = 1 | x_{(d)}, d = 1) - \mathbb{P}(y = 1 | x_{(d)}, d = 0).$$

where $x_{(d)}$ contains all other variables.

 However, the formula in (1) is often accurate enough for binary variables too.

Estimation of Probit and Logit Models

- ▶ We have to estimate the probit and logit models using MLE.
- The likelihood function is given by

$$p(\beta|Y,X) = \prod_{i=1}^{n} p_i(\beta|x_i)^{y_i}(1-p_i(\beta|x_i))^{1-y_i},$$

where

$$p_i(\beta|x_i) = F(x'_i\beta).$$

The log likelihood function is given by

$$\begin{split} \ell(\beta|Y,X) &= \sum_{i=1}^n \left\{ y_i \log F(x_i'\beta) \right. \\ &+ \left(1-y_i\right) \log(1-F(x_i'\beta)) \right\}. \end{split}$$

 The MLE estimator solves the following maximization problem,

$$\max_{\beta} \ell(\beta|Y,X).$$

Generalized Linear Model

- The GLM generalizes a large class of statistical models, including the classical linear regression and the logit/probit model.
- In a GLM, the dependent variables y is assumed to be generated from a particular distribution in the exponential family, which includes the normal, binomial, Poisson and gamma distributions, among others. The mean of the distribution depends on the independent variables x through:

$$\mathbb{E}(y|x) = \mu = g^{-1}(x'\beta),$$

where $g(\cdot)$ is a link function that links the linear predictor $\eta = x'\beta$ with μ , $\eta = g(\mu)$.

- Linear regression: y is normal, $\mu = x'\beta$, g is identity.
- Logit/probit model: y is binary, $\mu = F(x'\beta)$, $g = F^{-1}$.
- Poisson regression: y is Poisson, $\mu = \exp(x'\beta)$, $g = \log$.
- ...