

Introductory Econometrics

Lecture 14: Dummy variables

Instructor: Ma, Jun

Renmin University of China

April 26, 2023

Interval, Ordinal, and Categorical Variables

- ▶ Interval variable: is one where the difference between two values is meaningful. Example: “Education” when measured in years. There is a meaning to the difference between 12 and 10 years of education.
- ▶ In some data sets, education is reported as an ordinal variable: only the order between its values matters, but the difference has no meaning. Example: The following two variables are equivalent.

$$Education_i = \begin{cases} 1 & \text{if high-school graduate,} \\ 2 & \text{if college graduate,} \\ 3 & \text{if advanced degree.} \end{cases}$$

$$Education_i = \begin{cases} 1 & \text{if high-school graduate,} \\ 10 & \text{if college graduate,} \\ 234 & \text{if advanced degree.} \end{cases}$$

- ▶ Categorical variable is one that has one or more categories, but there is no natural ordering to the categories
Examples: Gender, race, marital status, geographic location.
- ▶ The following two variables are equivalent:
$$Gender_i = \begin{cases} 1 & \text{if observation } i \text{ corresponds to a } woman, \\ 2 & \text{if observation } i \text{ corresponds to a } man. \end{cases}$$
$$Gender_i = \begin{cases} 1 & \text{if observation } i \text{ corresponds to a } man, \\ 2 & \text{if observation } i \text{ corresponds to a } woman. \end{cases}$$
- ▶ Categorical and ordinal variables are also called qualitative.
- ▶ Qualitative variables cannot be simply included in regression, because the regression technique assumes that all variables are interval.

Dummy variables

- ▶ A dummy variable is a binary zero-one variable which takes on the value one if some condition is satisfied and zero if that condition fails:

- ▶ $Female_i = \begin{cases} 1 & \text{if observation } i \text{ corresponds to a woman,} \\ 0 & \text{if observation } i \text{ corresponds to a man.} \end{cases}$

- ▶ $Male_i = \begin{cases} 1 & \text{if observation } i \text{ corresponds to a man,} \\ 0 & \text{if observation } i \text{ corresponds to a woman.} \end{cases}$

- ▶ Note that $Female_i + Male_i = 1$ for all observations i .

- ▶ $Married_i = \begin{cases} 1 & \text{if married,} \\ 0 & \text{otherwise.} \end{cases}$

Example

TABLE 7.1

A Partial Listing of the Data in WAGE1.RAW

<i>person</i>	<i>wage</i>	<i>educ</i>	<i>exper</i>	<i>female</i>	<i>married</i>
1	3.10	11	2	1	0
2	3.24	12	22	1	1
3	3.00	11	2	0	0
4	6.00	8	44	0	1
5	5.30	12	7	0	1
.
.
.
525	11.56	16	5	0	1
526	3.50	14	5	1	0

A single dummy independent variable

- ▶ Consider the following regression:

$$Wage_i = \beta_0 + \delta_0 Female_i + \beta_1 Educ_i + \beta_3 Exper_i + \beta_4 Tenure_i + U_i,$$

and assume that conditionally on all independent variables, $E[U_i] = 0$.

- ▶ If observation i corresponds to a woman, $Female_i = 1$, and

$$E[Wage_i | Female_i = 1, Educ_i, Exper_i, Tenure_i] = \beta_0 + \delta_0 + \beta_1 Educ_i + \beta_3 Exper_i + \beta_4 Tenure_i.$$

- ▶ If observation i corresponds to a man, $Female_i = 0$, and

$$E[Wage_i | Female_i = 0, Educ_i, Exper_i, Tenure_i] = \beta_0 + \beta_1 Educ_i + \beta_3 Exper_i + \beta_4 Tenure_i.$$

- ▶ Thus,

$$\delta_0 = E[Wage_i | Female_i = 1, Educ_i, Exper_i, Tenure_i] - E[Wage_i | Female_i = 0, Educ_i, Exper_i, Tenure_i].$$

An intercept shift

- ▶ The model:

$$Wage_i = \beta_0 + \delta_0 Female_i + \beta_1 Educ_i + \beta_3 Exper_i + \beta_4 Tenure_i + U_i$$

- ▶ For men ($Female_i = 0$):, we can write the model as

$$Wage_i^M = \beta_0 + \beta_1 Educ_i + \beta_3 Exper_i + \beta_4 Tenure_i + U_i.$$

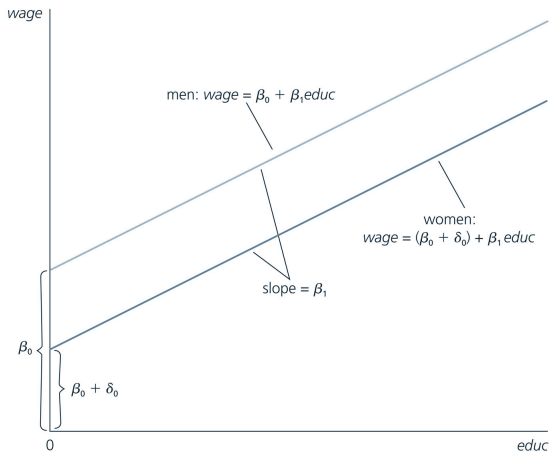
- ▶ For women ($Female_i = 1$):, we can write the model as

$$Wage_i^F = (\beta_0 + \delta_0) + \beta_1 Educ_i + \beta_3 Exper_i + \beta_4 Tenure_i + U_i.$$

- ▶ In this case, men play the role of the base group.
- ▶ δ_0 measures the difference relatively to the base group.

FIGURE 7.1

Graph of $wage = \beta_0 + \delta_0 female + \beta_1 educ$ for $\delta_0 < 0$.



Example

- ▶ Estimated equation:

$$\widehat{Wage}_i = - 1.57 \quad - 1.81 \quad Female_i + 0.572 \quad Educ_i \\ (0.72) \quad (0.26) \quad (0.049) \\ + 0.025 \quad Exper_i + 0.141 \quad Tenure_i. \\ (0.012) \quad (0.021)$$

- ▶ The dependent variable is the wage per hour.
- ▶ $\hat{\delta}_0 = -1.81$ implies that a women earns \$1.81 less per hour than a man with the same level of education, experience, and tenure. (These are 1976 wages.)
- ▶ The difference is also statistically significant.

When the dependent variable is in the logarithmic form

- ▶ The model:

$$\log(\text{Wage}) = \beta_0 + \delta_0 \text{Female} + \beta_1 \text{Educ} + \beta_3 \text{Exper} + \beta_4 \text{Tenure} + U.$$

- ▶ In this case,

$$\begin{aligned}\delta_0 &= \log(\text{Wage}^F) - \log(\text{Wage}^M) \\ &= \log\left(\frac{\text{Wage}^F}{\text{Wage}^M}\right) \\ &= \log\left(\frac{\text{Wage}^M + (\text{Wage}^F - \text{Wage}^M)}{\text{Wage}^M}\right) \\ &= \log\left(1 + \frac{\text{Wage}^F - \text{Wage}^M}{\text{Wage}^M}\right) \\ &\approx \frac{\text{Wage}^F - \text{Wage}^M}{\text{Wage}^M}.\end{aligned}$$

- ▶ When the dependent variable is in the log form, δ_0 has a percentage interpretation.

Example

- ▶ Estimated equation:

$$\begin{aligned}\widehat{\log(Wage_i)} = & 0.417 & - & 0.297 & Female_i & + & 0.080 & Educ_i \\ & (0.099) & & (0.036) & & & (0.007) & \\ & + & 0.029 & Exper_i & - & 0.00058 & Exper_i^2 & \\ & & (0.005) & & & (0.00010) & & \\ & + & 0.032 & Tenure_i & - & 0.00059 & Tenure_i^2 & \\ & & (0.007) & & & (0.00023) & & \end{aligned}$$

- ▶ $\hat{\delta}_0 = -0.297$ implies that a woman earns 29.7% less than a man with the same level of education, experience and tenure.

Changing the base group

- ▶ Instead of

$$\log(\text{Wage}_i) = \beta_0 + \delta_0 \text{Female}_i + \beta_1 \text{Educ}_i + \beta_3 \text{Exper}_i + \beta_4 \text{Tenure}_i + U_i$$

consider:

$$\log(\text{Wage}_i) = \theta_0 + \gamma_0 \text{Male}_i + \theta_1 \text{Educ}_i + \theta_3 \text{Exper}_i + \theta_4 \text{Tenure}_i + U_i.$$

- ▶ Since $\text{Male}_i = 1 - \text{Female}_i$,

$$\begin{aligned} \log(\text{Wage}_i) &= \theta_0 + \gamma_0 \text{Male}_i + \theta_1 \text{Educ}_i + \theta_3 \text{Exper}_i + \theta_4 \text{Tenure}_i + U_i \\ &= \theta_0 + \gamma_0 (1 - \text{Female}_i) + \theta_1 \text{Educ}_i + \theta_3 \text{Exper}_i + \theta_4 \text{Tenure}_i + U_i \\ &= (\theta_0 + \gamma_0) - \gamma_0 \text{Female}_i + \theta_1 \text{Educ}_i + \theta_3 \text{Exper}_i + \theta_4 \text{Tenure}_i + U_i. \end{aligned}$$

- ▶ We conclude that $\delta_0 = -\gamma_0$, $\beta_0 = \theta_0 - \delta_0$, $\beta_1 = \theta_1$, and etc.:

$$\log(\text{Wage}_i) = (\beta_0 + \delta_0) - \delta_0 \text{Female}_i + \beta_1 \text{Educ}_i + \beta_3 \text{Exper}_i + \beta_4 \text{Tenure}_i + U_i.$$

- ▶ Thus, changing the base group has no effect on the conclusions.

The dummy variable trap

- ▶ Consider the equation:

$$\begin{aligned}\log(Wage_i) = & \beta_0 + \delta_0 Female_i + \gamma_0 Male_i \\ & + \beta_1 Educ_i + \beta_3 Exper_i + \beta_4 Tenure_i + U_i.\end{aligned}$$

- ▶ Recall that the intercept is a regressor that takes the value one for all observations.
- ▶ Since $Male_i + Female_i - 1 = 0$ for all observations i , we have the case of perfect multicollinearity, and such an equation cannot be estimated.
- ▶ One cannot include an intercept and dummies for all the groups!

- ▶ One of the dummies has to be omitted and the corresponding group becomes the base group:
 - ▶ Men are the base group: $\log(Wage_i) = \beta_0 + \delta_0 Female_i + \beta_1 Educ_i + \beta_3 Exper_i + \beta_4 Tenure_i + U_i.$
 - ▶ Women are the base group: $\log(Wage_i) = \theta_0 + \gamma_0 Male_i + \beta_1 Educ_i + \beta_3 Exper_i + \beta_4 Tenure_i + U_i.$
- ▶ Alternatively, one can include both dummies without the intercept: $\log(Wage_i) = \pi_0 Female_i + \pi_1 Male_i + \beta_1 Educ_i + \beta_3 Exper_i + \beta_4 Tenure_i + U_i.$
 - ▶ In Stata regression with no intercept can be estimated by using the option "no constant":
`regress Y X, noconstant`
 - ▶ The coefficients on the dummy variables lose the difference interpretation.

A slope shift and interactions

- ▶ We can also allow the returns to education to be different for men and women:

$$\log(Wage_i) = \beta_0 + \delta_0 Female_i + \beta_1 Educ_i + \delta_1 (Female_i \cdot Educ_i) \\ + \beta_3 Exper_i + \beta_4 Tenure_i + U_i.$$

- ▶ The variable $(Female_i \cdot Educ_i)$ is called an interaction.
- ▶ The equation for men ($Female_i = 0$):

$$\log(Wage_i^M) = \beta_0 + \beta_1 Educ_i + \beta_3 Exper_i + \beta_4 Tenure_i + U_i.$$

- ▶ The equation for women ($Female_i = 1$):

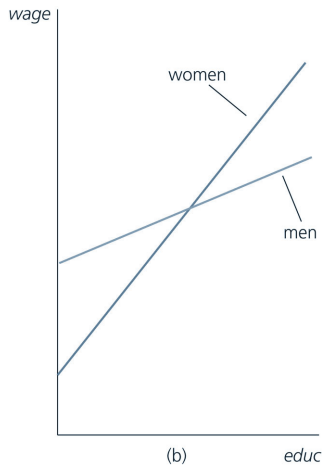
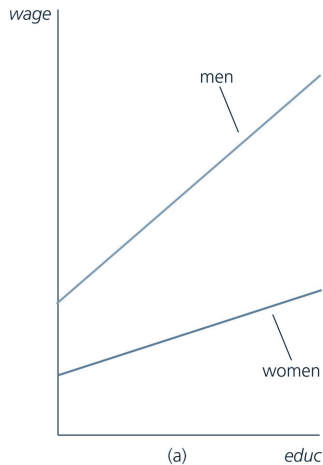
$$\log(Wage_i^F) = (\beta_0 + \delta_0) + (\beta_1 + \delta_1) Educ_i \\ + \beta_3 Exper_i + \beta_4 Tenure_i + U_i.$$

- ▶ δ_1 can be interpreted as the difference in return to education between the women and men (the base group) after controlling for experience and tenure.

A slope shift

FIGURE 7.2

Graphs of equation (7.16): (a) $\delta_0 < 0, \delta_1 < 0$; (b) $\delta_0 < 0, \delta_1 > 0$.



Example

- ▶ Estimated equation:

$$\begin{aligned}\widehat{\log(Wage_i)} = & 0.389 & - & 0.227 & Female_i \\ & (0.119) & & (0.168) & \\ & + & 0.082 & Educ_i & - & 0.0056 & Female_i \cdot Educ_i \\ & & (0.008) & & & (0.0131) & \\ & + & 0.029 & Exper_i & - & 0.00058 & Exper_i^2 \\ & & (0.005) & & & (0.00011) & \\ & + & 0.032 & Tenure_i & - & 0.00059 & Tenure_i^2. \\ & & (0.007) & & & (0.00024) & \end{aligned}$$

- ▶ $\hat{\delta}_1 = -0.0056$ suggesting that the return to education for women is 0.56% less than for men, however it is not statistically significant. Thus, we can conclude that the return to education is the same for men and women.

Multiple categories

- ▶ In the previous examples, *Educ* was a quantitative variable: years of education.
- ▶ Suppose now that instead the education variable is ordinal:

$$Education = \begin{cases} 1 & \text{if high-school dropout,} \\ 2 & \text{if high-school graduate,} \\ 3 & \text{if some college,} \\ 4 & \text{if college graduate,} \\ 5 & \text{if advanced degree.} \end{cases}$$

- ▶ Only the order is important, and there is no meaning to the distance between the values.
- ▶ Adding such a variable to the regression will give a meaningless result.

$$Education_i = \begin{cases} 1 & \text{if high-school dropout,} \\ 2 & \text{if high-school graduate,} \\ 3 & \text{if some college,} \\ 4 & \text{if college graduate,} \\ 5 & \text{if advanced degree.} \end{cases}$$

- ▶ Define 5 new dummy variables:

$$E_{1,i} = \begin{cases} 1 & \text{if high-school dropout,} \\ 0 & \text{otherwise.} \end{cases} \quad E_{2,i} = \begin{cases} 1 & \text{if high-school graduate,} \\ 0 & \text{otherwise.} \end{cases}$$

$$E_{3,i} = \begin{cases} 1 & \text{if some college,} \\ 0 & \text{otherwise.} \end{cases} \quad E_{4,i} = \begin{cases} 1 & \text{if college graduate,} \\ 0 & \text{otherwise.} \end{cases}$$

$$E_{5,i} = \begin{cases} 1 & \text{if advanced degree,} \\ 0 & \text{otherwise.} \end{cases}$$

- ▶ To avoid the dummy variable trap, one of the dummies has to be omitted:

$$Wage_i = \beta_0 + \delta_0 Female_i + \delta_2 E_{2,i} + \delta_3 E_{3,i} + \delta_4 E_{4,i} + \delta_5 E_{5,i} + \text{Other Factors}$$

- ▶ Group 1 (high-school dropout) becomes the base group.
- ▶ δ_2 measures the wage difference between high-school graduates and high-school dropouts.
- ▶ δ_3 measures the wage difference between individuals with some college education and high-school dropouts.

Testing for structural breaks or differences in regression functions across groups

- ▶ Suppose for simplicity we have two groups. For example,
 - ▶ Male and female workers.
 - ▶ Observations before and after a certain date.
- ▶ We want to test if the intercept and all slopes are the same across the two groups.
- ▶ The model:

$$Y_i = \beta_{1,0} + \beta_{1,1}X_{1,i} + \dots + \beta_{1,k}X_{k,i} + U_i \text{ if } i \text{ belongs to Group 1}$$

$$Y_i = \beta_{2,0} + \beta_{2,1}X_{1,i} + \dots + \beta_{2,k}X_{k,i} + U_i \text{ if } i \text{ belongs to Group 2}$$

- ▶ The hypotheses:

$$H_0 : \beta_{1,0} = \beta_{2,0}, \beta_{1,1} = \beta_{2,1}, \dots, \beta_{1,k} = \beta_{2,k}.$$

$$H_1 : \beta_{1,j} \neq \beta_{2,j} \text{ at least for one } j \in \{0, 1, \dots, k\}.$$

$$Y_i = \beta_{1,0} + \beta_{1,1}X_{1,i} + \dots + \beta_{1,k}X_{k,i} + U_i \text{ if } i \text{ belongs to Group 1}$$

$$Y_i = \beta_{2,0} + \beta_{2,1}X_{1,i} + \dots + \beta_{2,k}X_{k,i} + U_i \text{ if } i \text{ belongs to Group 2}$$

- ▶ The Chow F statistic:

$$F^{Chow} = \frac{(SSR_r - SSR_{ur}) / (k + 1)}{SSR_{ur} / (n - 2(k + 1))} = \frac{(SSR_r - (SSR_1 + SSR_2)) / (k + 1)}{(SSR_1 + SSR_2) / (n - 2(k + 1))},$$

where

- ▶ SSR_1 is the SSR obtained by estimating the model using only the observations from Group 1.
- ▶ SSR_2 is the SSR obtained by estimating the model using only the observations from Group 2.
- ▶ SSR_r is the SSR obtained by pooling the groups and estimating a single equation:

$$Y_i = \gamma_0 + \gamma_1 X_{1,i} + \dots + \gamma_k X_{k,i} + U_i \text{ for all } i\text{'s (Groups 1 and 2).}$$

- ▶ H_0 of constancy or no structural break is rejected when

$$F^{Chow} > F_{k+1, n-2(k+1), 1-\alpha}.$$

- ▶ The Chow test can also be performed using the dummy variables, and the two approaches are numerically equivalent.

- ▶ Define

$$D_i = \begin{cases} 1 & \text{observation } i \text{ belongs to Group 1,} \\ 0 & \text{otherwise.} \end{cases}$$

- ▶ Estimate the following single equation using all observations (Groups 1 and 2):

$$Y_i = \beta_0 + \beta_1 X_{1,i} + \dots + \beta_k X_{k,i} \\ + \delta_0 D_i + \delta_1 (D_i \cdot X_{1,i}) + \dots + \delta_k (D_i \cdot X_{k,i}) + U_i.$$

- ▶ Test:

$$H_0 \quad : \quad \delta_0 = \delta_1 = \dots = \delta_k = 0.$$

$$H_1 \quad : \quad \delta_j \neq 0 \text{ for at least one } j \in \{0, 1, \dots, k\}.$$