## ECON 351* -- Introduction to NOTE 21

# Introduction to Using Dummy Variable Regressors in Regression Models

- Consider the regression model for the average weekly earnings of individual workers given by the following **population regression equation (PRE)**:

$$\text{earn}_i = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{fem}_i + \beta_3 \text{fem}_i \text{ed}_i + u_i \tag{1}$$

where $u_i$ is an iid (independently and identically distributed) random error term, and the observable variables are defined as follows:

*earn$_i$* = person i's **average weekly earnings** during a calendar year;

*ed$_i$*   = person i's years of **completed formal education**;

*fem$_i$*  = a **female *indicator* variable**, or ***dummy* variable**, defined such that
            fem$_i$ = 1 if person i is female, = 0 if person i is male.

- The earnings regression equation (1) contains two explanatory variables: the continuous variable *ed$_i$* and the indicator (dummy) variable regressor *fem$_i$*.

- The dummy variable regressor *fem$_i$* enters earnings regression equation (1) in two distinct ways:

  1. it enters *additively* on its own as *fem$_i$*;

  2. it enters *multiplicatively* interacted with *ed$_i$* in the regressor *fem$_i$ed$_i$*.

## 1. Interpreting the Simple Earnings Regression with a Female Dummy Variable Regressor

$$earn_i = \beta_0 + \beta_1 ed_i + \beta_2 fem_i + \beta_3 fem_i ed_i + u_i \tag{1}$$

♦ *Question 1:* **What is the slope coefficient $\beta_2$ of the female indicator variable** *fem$_i$***?**

**Answer:** The slope coefficient $\beta_2$ of the dummy variable regressor *fem$_i$* is the **female-male** *difference* **in intercept coefficients**, i.e.,

$\beta_2$ = the *female* **intercept coefficient − the** *male* **intercept coefficient**

♦ *Question 2:* **What is the slope coefficient $\beta_3$ of the female indicator interaction term** *fem$_i$ed$_i$***?**

**Answer:** The slope coefficient $\beta_3$ of the dummy variable interaction regressor *fem$_i$ed$_i$* is the **female-male** *difference* **in the slope coefficients of the regressor** *ed$_i$*.

$\beta_3$ = the *female* **slope coefficient of** *ed$_i$* − the *male* **slope coefficient of** *ed$_i$*

♦ *Question 3:* **What is the equation intercept coefficient $\beta_0$ in equation (1)?**

**Answer:** The equation intercept coefficient $\beta_0$ is the intercept of the earnings equation for males, for whom the female dummy variable *fem$_i$* equals 0, i.e., for whom *fem$_i$* = 0.

♦ *Question 4:* **What is the slope coefficient $\beta_1$ of the regressor** *ed$_i$* **in equation (1)?**

**Answer:** The slope coefficient $\beta_1$ of the regressor *ed$_i$* in equation (1) is the slope coefficient of *ed$_i$* **for** *males* for whom *fem$_i$* = 0 by definition.

## 2. Demonstrating the Interpretation of the Regression Coefficients in Equation (1)

$$\text{earn}_i = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{fem}_i + \beta_3 \text{fem}_i \text{ed}_i + u_i \tag{1}$$

♦ The **population regression function (PRF)** corresponding to regression equation (1) gives the **conditional mean weekly earnings of female and male workers** with different levels of formal education (i.e., different values of the continuous explanatory variable $ed_i$).

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i\right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{fem}_i + \beta_3 \text{fem}_i \text{ed}_i \tag{2}$$

The population regression function (2) contains two separate regression functions: (1) a female regression function; and (2) a male regression function.

♦ The *female* **population regression function** is obtained by setting the **female indicator variable** $fem_i$ equal to 1 everywhere it appears in regression function (2). It gives the conditional mean earnings of female workers as a function of $ed_i$, years of formal education.

Setting $fem_i = \mathbf{1}$ in regression function (2) gives:

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 1\right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 + \beta_3 \text{ed}_i = \beta_0 + \beta_2 + (\beta_1 + \beta_3) \text{ed}_i \tag{2f}$$

♦ The *male* **population regression function** is obtained by setting the **female indicator variable** $fem_i$ equal to 0 everywhere it appears in regression function (2). It gives the conditional mean earnings of male workers as a function of $ed_i$, years of formal education.

Setting $fem_i = \mathbf{0}$ in regression function (2) gives:

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 0\right) = \beta_0 + \beta_1 \text{ed}_i \tag{2m}$$

♦ The *female-male* **difference in mean weekly earnings** is obtained by subtracting the male regression function (2m) from the female regression function (2f).

• The *female* **mean earnings function** is the female regression function (2f):

$$E\left( \text{earn}_i \mid \text{ed}_i, \text{fem}_i = 1 \right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 + \beta_3 \text{ed}_i = \beta_0 + \beta_2 + (\beta_1 + \beta_3)\text{ed}_i \qquad (2f)$$

Note: $\beta_0 + \beta_2$ is the *intercept* of the mean earnings function **for** *females*; $\beta_1 + \beta_3$ is the **slope coefficient** of the regressor $\text{ed}_i$ **for** *females***.**

• The *male* **mean earnings function** is the male regression function (2m):

$$E\left( \text{earn}_i \mid \text{ed}_i, \text{fem}_i = 0 \right) = \beta_0 + \beta_1 \text{ed}_i \qquad (2m)$$

Note: $\beta_0$ is the *intercept* of the mean earnings function **for** *males*; $\beta_1$ is the **slope coefficient** of the regressor $\text{ed}_i$ **for** *males***.**

• Subtracting the male regression function (2m) from the female regression function (2f) gives the expression for the *female-male difference* **in mean weekly earnings**:

$$E\left( \text{earn}_i \mid \text{ed}_i, \text{fem}_i = 1 \right) - E\left( \text{earn}_i \mid \text{ed}_i, \text{fem}_i = 0 \right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 + \beta_3 \text{ed}_i - \beta_0 - \beta_1 \text{ed}_i = \beta_2 + \beta_3 \text{ed}_i$$

i.e.,

$$E\left( \text{earn}_i \mid \text{ed}_i, \text{fem}_i = 1 \right) - E\left( \text{earn}_i \mid \text{ed}_i, \text{fem}_i = 0 \right) = \beta_2 + \beta_3 \text{ed}_i \qquad (3)$$

**_Result:_** Regression equation (2) implies that the *female-male difference* **in mean weekly earnings** is the linear function (3) of the explanatory variable $\text{ed}_i$:

**Interpretation of the Regression Coefficients in Regression Equation (1)**

$$\text{earn}_i = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{fem}_i + \beta_3 \text{fem}_i \text{ed}_i + u_i \tag{1}$$

♦ The **population regression function (PRF)** corresponding to regression equation (1) is:

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i\right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{fem}_i + \beta_3 \text{fem}_i \text{ed}_i \tag{2}$$

• The *female* **mean earnings function** is the female regression function (2f):

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 1\right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 + \beta_3 \text{ed}_i = \beta_0 + \beta_2 + (\beta_1 + \beta_3)\text{ed}_i \tag{2f}$$

• The *male* **mean earnings function** is the male regression function (2m):

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 0\right) = \beta_0 + \beta_1 \text{ed}_i \tag{2m}$$

The *female* **intercept** coefficient   $= \boldsymbol{\beta_0 + \beta_2}$

The *male* **intercept** coefficient   $= \boldsymbol{\beta_0}$

Therefore  $\boldsymbol{\beta_2}$ **= the** *female* **intercept coefficient** $(\beta_0 + \beta_2)$ **minus the** *male* **intercept coefficient** $(\beta_0)$

The *female* **slope coefficient of** *ed$_i$*   $= \boldsymbol{\beta_1 + \beta_3}$

The *male* **slope coefficient of** *ed$_i$*   $= \boldsymbol{\beta_1}$

Therefore  $\boldsymbol{\beta_3}$ **= the** *female* **slope coefficient of** *ed$_i$* $(\beta_1 + \beta_3)$ **minus the** *male* **slope coefficient of** *ed$_i$* $(\beta_1)$

## 4. A More General Earnings Regression Model with a Female Dummy Variable Regressor

Consider now an expanded earnings regression model for female and male workers that allows for non-constant marginal earnings effects of $ed_i$. The population regression equation for this model is:

$$\text{earn}_i = \beta_0 + \beta_1 ed_i + \beta_2 ed_i^2 + \beta_3 \text{fem}_i + \beta_4 \text{fem}_i ed_i + \beta_5 \text{fem}_i ed_i^2 + u_i \tag{3}$$

Note that regression equation (3) allows for the possibly that the **marginal earnings effect of $ed$** may be *increasing* or *decreasing* in $ed_i$.

♦ The **population regression function (PRF)** corresponding to regression equation (3) gives the **conditional mean weekly earnings of female and male workers** with different levels of formal education (i.e., different values of the continuous explanatory variable $ed_i$).

$$\text{E}\left(\text{earn}_i \mid ed_i, \text{fem}_i\right) = \beta_0 + \beta_1 ed_i + \beta_2 ed_i^2 + \beta_3 \text{fem}_i + \beta_4 \text{fem}_i ed_i + \beta_5 \text{fem}_i ed_i^2 \tag{3p}$$

The population regression function (3p) contains two separate regression functions: (1) a female regression function; and (2) a male regression function.

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i\right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{ed}_i^2 + \beta_3 \text{fem}_i + \beta_4 \text{fem}_i \text{ed}_i + \beta_5 \text{fem}_i \text{ed}_i^2 \qquad (3p)$$

♦ The *female* **population regression function** is obtained by setting the **female indicator variable** *fem_i* equal to 1 everywhere it appears in regression function (3p). It gives the conditional mean earnings of female workers as a function of *ed_i*, years of formal education.

Setting *fem_i* = **1** in regression function (3p) gives the *female* **mean earnings function**:

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 1\right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{ed}_i^2 + \beta_3 + \beta_4 \text{ed}_i + \beta_5 \text{ed}_i^2 = \beta_0 + \beta_3 + (\beta_1 + \beta_4)\text{ed}_i + (\beta_2 + \beta_5)\text{ed}_i^2 \quad (3f)$$

♦ The *male* **population regression function** is obtained by setting the **female indicator variable** *fem_i* equal to 0 everywhere it appears in regression function (3p). It gives the conditional mean earnings of male workers as a function of *ed_i*, years of formal education.

Setting *fem_i* = **0** in regression function (3p) gives the *male* **mean earnings function**:

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 0\right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{ed}_i^2 \qquad (3m)$$

♦ The *female-male* **difference in mean weekly earnings** is obtained by subtracting the male regression function (3m) from the female regression function (3f):

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 1\right) - E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 0\right)$$
$$= \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{ed}_i^2 + \beta_3 + \beta_4 \text{ed}_i + \beta_5 \text{ed}_i^2 - \beta_0 - \beta_1 \text{ed}_i - \beta_2 \text{ed}_i^2$$
$$= \beta_3 + \beta_4 \text{ed}_i + \beta_5 \text{ed}_i^2 \qquad (3d)$$

**_Result:_** Regression equation (3) implies that the *female-male difference* **in mean weekly earnings** is the quadratic function (3d) of the continuous explanatory variable *ed_i*.

## Summary Interpretation of the Regression Coefficients in Regression Equation (3)

♦ The **population regression function (PRF)** corresponding to regression equation (3) is:

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i\right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{ed}_i^2 + \beta_3 \text{fem}_i + \beta_4 \text{fem}_i \text{ed}_i + \beta_5 \text{fem}_i \text{ed}_i^2 \qquad (3p)$$

• The *female* **mean earnings function** is the female regression function (3f):

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 1\right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{ed}_i^2 + \beta_3 + \beta_4 \text{ed}_i + \beta_5 \text{ed}_i^2 = \beta_0 + \beta_3 + (\beta_1 + \beta_4)\text{ed}_i + (\beta_2 + \beta_5)\text{ed}_i^2 \qquad (3f)$$

• The *male* **mean earnings function** is the male regression function (3m):

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 0\right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{ed}_i^2 \qquad (3m)$$

The *female* **intercept** coefficient   $= \beta_0 + \beta_3$
The *male* **intercept** coefficient    $= \beta_0$

Therefore  $\beta_3$ = **the** *female* **intercept coefficient** ($\beta_0 + \beta_3$) **minus** **the** *male* **intercept coefficient** ($\beta_0$)

The *female* **slope coefficient of** $ed_i$  $= \beta_1 + \beta_4$
The *male* **slope coefficient of** $ed_i$   $= \beta_1$

Therefore  $\beta_4$ = **the** *female* **slope coefficient of** $ed_i$ ($\beta_1 + \beta_4$) **minus** **the** *male* **slope coefficient of** $ed_i$ ($\beta_1$)

The *female* **slope coefficient of** $ed_i^2$  $= \beta_2 + \beta_5$
The *male* **slope coefficient of** $ed_i^2$   $= \beta_2$

Therefore  $\beta_5$ = **the** *female* **slope coefficient of** $ed_i^2$ ($\beta_2 + \beta_5$) **minus** **the** *male* **slope coefficient of** $ed_i^2$ ($\beta_2$)

## The Marginal Earnings Effect of ed in Regression Equation (3)

♦ The **marginal earnings effect of $ed_i$** is obtained by partially differentiating the **population regression function (PRF)** for regression equation (3) with respect to $ed_i$:

$$E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i\right) = \beta_0 + \beta_1 \text{ed}_i + \beta_2 \text{ed}_i^2 + \beta_3 \text{fem}_i + \beta_4 \text{fem}_i \text{ed}_i + \beta_5 \text{fem}_i \text{ed}_i^2 \qquad (3p)$$

$$\frac{\partial E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i\right)}{\partial \text{ed}_i} = \beta_1 + 2\beta_2 \text{ed}_i + \beta_4 \text{fem}_i + 2\beta_5 \text{fem}_i \text{ed}_i \qquad (4)$$

♦ The **marginal earnings effect of $ed_i$ for *females*** is obtained by setting the female indicator variable $fem_i = 1$ in the above expression (4) for the marginal earnings effect of $ed_i$:

$$\frac{\partial E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 1\right)}{\partial \text{ed}_i} = \beta_1 + 2\beta_2 \text{ed}_i + \beta_4 + 2\beta_5 \text{ed}_i = \left(\beta_1 + \beta_4\right) + 2\left(\beta_2 + \beta_5\right)\text{ed}_i \qquad (4f)$$

♦ The **marginal earnings effect of $ed_i$ for *males*** is obtained by setting the female indicator variable $fem_i = 0$ in expression (4) for the marginal earnings effect of $ed_i$:

$$\frac{\partial E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 0\right)}{\partial \text{ed}_i} = \beta_1 + 2\beta_2 \text{ed}_i \qquad (4m)$$

♦ The *female-male difference* in the marginal earnings effect of $ed_i$ is obtained by subtracting (4m) from (4f):

$$\frac{\partial E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 1\right)}{\partial \text{ed}_i} - \frac{\partial E\left(\text{earn}_i \mid \text{ed}_i, \text{fem}_i = 0\right)}{\partial \text{ed}_i} = \beta_4 + 2\beta_5 \text{ed}_i \qquad (4d)$$