

Multiple Regression Analysis with Qualitative Information



Chapter 6 (Ch. 7 of the textbook)

Wooldridge: Introductory Econometrics:
A Modern Approach, 5e

Multiple Regression Analysis: Qualitative Information

- **Qualitative Information**

- Examples: gender, race, industry, region, rating grade, ...
- A way to incorporate qualitative information is to use dummy variables
- They may appear as the dependent or as independent variables

- **A single dummy independent variable**

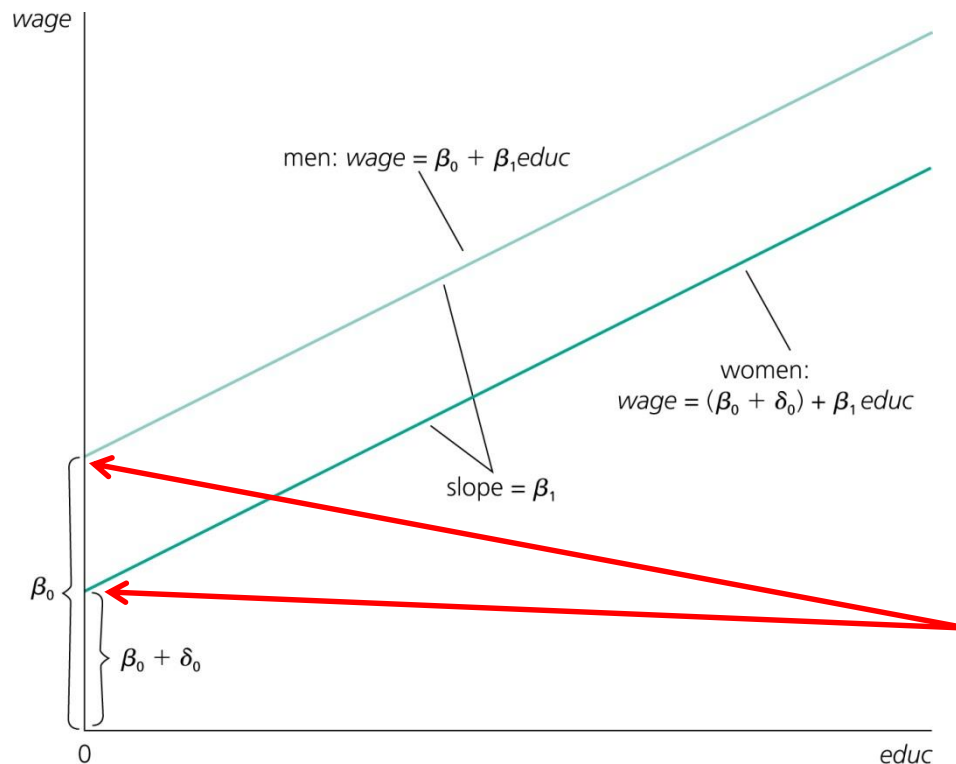
$$wage = \beta_0 + \delta_0 female + \beta_1 educ + u$$

→ = the wage gain/loss if the person is a woman rather than a man (holding other things fixed)

← Dummy variable:
=1 if the person is a woman
=0 if the person is man

Multiple Regression Analysis: Qualitative Information

■ Graphical Illustration



Alternative interpretation of coefficient:

$$\delta_0 = E(\text{wage} | \text{female} = 1, \text{educ}) \\ - E(\text{wage} | \text{female} = 0, \text{educ})$$

i.e. the difference in mean wage between men and women with the same level of education.

Intercept shift

Multiple Regression Analysis: Qualitative Information

- **Dummy variable trap**

This model cannot be estimated (perfect collinearity)

$$wage = \beta_0 + \gamma_0 \text{male} + \delta_0 \text{female} + \beta_1 educ + u$$

When using dummy variables, one category always has to be omitted:

$$wage = \beta_0 + \delta_0 \text{female} + \beta_1 educ + u$$

← The base category are men

$$wage = \beta_0 + \gamma_0 \text{male} + \beta_1 educ + u$$

← The base category are women

Alternatively, one could omit the intercept:

$$wage = \gamma_0 \text{male} + \delta_0 \text{female} + \beta_1 educ + u$$

← Disadvantages:

- 1) More difficult to test for differences between the parameters
- 2) R-squared formula only valid if regression contains intercept

Multiple Regression Analysis: Qualitative Information

- **Estimated wage equation with intercept shift**

$$\widehat{wage} = -1.57_{(.72)} - 1.81_{(.26)} female + .572_{(.049)} educ$$

$$+ .025_{(.012)} exper + .141_{(.021)} tenure$$

Holding education, experience, and tenure fixed, women earn 1.81\$ less per hour than men

$$n = 526, R^2 = .364$$

- **Does that mean that women are discriminated against?**

- Not necessarily. Being female may be correlated with other productivity characteristics that have not been controlled for.

Multiple Regression Analysis: Qualitative Information

- **Comparing means of subpopulations described by dummies**

$$\widehat{wage} = 7.10 - 2.51 \text{ female}$$

(.21) (.26)

Not holding other factors constant, women earn 2.51\$ per hour less than men, i.e. the difference between the mean wage of men and that of women is 2.51\$.

$$n = 526, R^2 = .116$$

- **Discussion**

- It can easily be tested whether difference in means is significant
- The wage difference between men and women is larger if no other things are controlled for; i.e. part of the difference is due to differences in education, experience and tenure between men and women

Multiple Regression Analysis: Qualitative Information

- **Further example: Effects of training grants on hours of training**

Hours training per employee

Dummy indicating whether firm received training grant

$$\widehat{hrsemp} = 46.67 + 26.25 \textit{grant} - 0.98 \log(\textit{sales})$$

(43.41) (5.59) (3.54)

$$- 6.07 \log(\textit{employ}), \quad n = 105, R^2 = .237$$

(3.88)

- **This is an example of program evaluation**

- Treatment group (= grant receivers) vs. control group (= no grant)
- Is the effect of treatment on the outcome of interest causal?

Multiple Regression Analysis: Qualitative Information

- Using dummy explanatory variables in equations for $\log(y)$

$$\widehat{\log(price)} = -1.35 + .168 \log(lotsize) + .707 \log(sqft)$$

(.65) (.038) (.093)

$$+ .027 bdrms + .054 \textit{colonial}$$

(.029) (.045)

← Dummy indicating whether house is of colonial style

$$n = 88, R^2 = .649$$

$$\Rightarrow \frac{\partial \log(price)}{\partial \textit{colonial}} = \frac{\% \Delta price}{\Delta \textit{colonial}} = 5.4\%$$

← As the dummy for colonial style changes from 0 to 1, the house price increases by 5.4%

Multiple Regression Analysis: Qualitative Information

- **Using dummy variables for multiple categories**
 - 1) Define membership in each category by a dummy variable
 - 2) Leave out one category (which becomes the base category)

$$\widehat{\log(wage)} = .321 + .213 \text{ marrmale} - .198 \text{ marrfem} \\ - .110 \text{ singfem} + .079 \text{ educ} + .027 \text{ exper} - .00054 \text{ exper}^2 \\ + .029 \text{ tenure} - .00053 \text{ tenure}^2$$

Holding other things fixed, married women earn 19.8% less than single men (= the base category)

Multiple Regression Analysis: Qualitative Information

- **Using dummy variables for multiple categories**

- **Exact variations:** $(e^{0.213} - 1) \times 100\% = 23.74\%$

$$(e^{-0.198} - 1) \times 100\% = -17.96\% \quad (e^{-0.11} - 1) \times 100\% = -10.42\%$$

$$\widehat{\log}(wage) = .321 + .213 \text{ marrmale} - .198 \text{ marrfem}$$

(.100)
(.055)
(.058)

$$- .110 \text{ singfem} + .079 \text{ educ} + .027 \text{ exper} - .00054 \text{ exper}^2$$

(.056)
(.007)
(.005)
(.00011)

Holding other things fixed, married women earn 19.8% less than single men (= the base category)

$$+ .029 \text{ tenure} - .00053 \text{ tenure}^2$$

(.007)
(.00023)

Multiple Regression Analysis: Qualitative Information

- Using dummy variables for multiple categories

Changing the base category

$$\log(\text{wage}) = 0.123 + 0.411 \text{marrmale} + 0.198 \text{singmale} + 0.088 \text{singfemale} + 0.079 \text{educ} + \dots$$

Note: In the original image, the coefficient 0.198 for singmale is circled in red, and a red arrow points from this coefficient to the explanatory text below. The variables marrmale, singmale, and singfemale are also enclosed in red dashed boxes.

Holding other things fixed, single men earn 19.8% more than married women (= the base category)

Multiple Regression Analysis: Qualitative Information

- **Incorporating ordinal information using dummy variables**
- **Example: City credit ratings and municipal bond interest rates**

Municipal bond rate

Credit rating from 0-4 (0=worst, 4=best)

$$MBR = \beta_0 + \beta_1 CR + \text{other factors}$$

This specification would probably not be appropriate as the credit rating only contains ordinal information. A better way to incorporate this information is to define dummies:

$$MBR = \beta_0 + \delta_1 CR_1 + \delta_2 CR_2 + \delta_3 CR_3 + \delta_4 CR_4 + \text{other factors}$$

Dummies indicating whether the particular rating applies, e.g. $CR_1=1$ if $CR=1$ and $CR_1=0$ otherwise. All effects are measured in comparison to the worst rating (= base category).

Multiple Regression Analysis: Qualitative Information

- **Interactions involving dummy variables**
- **Allowing for different slopes**

$$\log(\text{wage}) = \beta_0 + \delta_0 \text{female} + \beta_1 \text{educ} + \delta_1 \text{female} \cdot \text{educ} + u$$

Interaction term

$$\beta_0 = \text{intercept men} \qquad \beta_1 = \text{slope men}$$

$$\beta_0 + \delta_0 = \text{intercept women} \qquad \beta_1 + \delta_1 = \text{slope women}$$

- **Interesting hypotheses**

$$H_0 : \delta_1 = 0$$

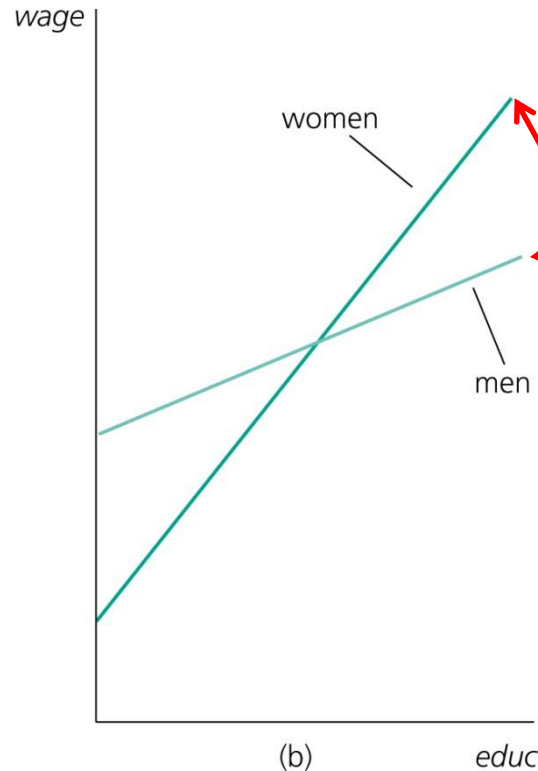
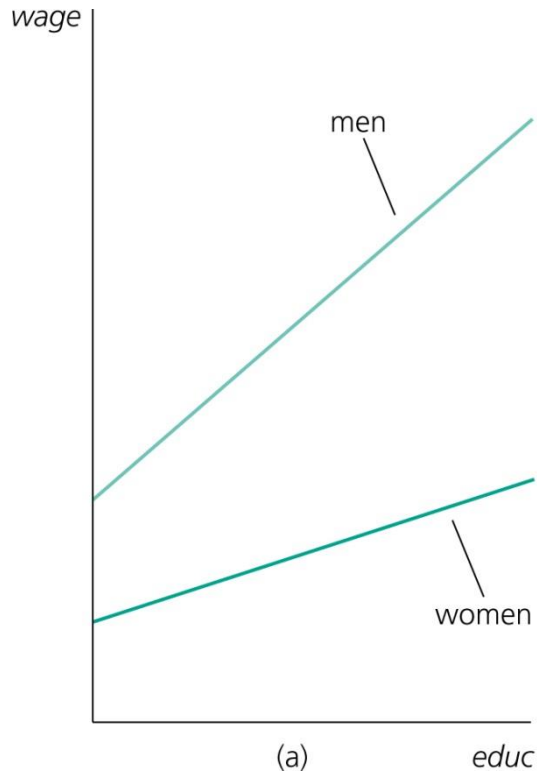
The return to education is the same for men and women

$$H_0 : \delta_0 = 0, \delta_1 = 0$$

The whole wage equation is the same for men and women

Multiple Regression Analysis: Qualitative Information

■ Graphical illustration



Interacting both the intercept and the slope with the female dummy enables one to model completely independent wage equations for men and women

Multiple Regression Analysis: Qualitative Information

■ Estimated wage equation with interaction term

$$\widehat{\log(wage)} = .389 - .227 \text{ female} - .082 \text{ educ} \\ (.119) \quad (.168) \quad (.008) \\ - .0056 \text{ female} \cdot \text{educ} + .029 \text{ exper} - .00058 \text{ exper}^2 \\ (.0131) \quad (.005) \quad (.00011) \\ + .032 \text{ tenure} - .00059 \text{ tenure}^2, n = 526, R^2 = .441 \\ (.007) \quad (.00024)$$

No evidence against hypothesis that the return to education is the same for men and women

Does this mean that there is no significant evidence of lower pay for women at the same levels of educ, exper, and tenure? No: this is only the effect for educ = 0. To answer the question one has to recenter the interaction term, e.g. around educ = 12.5 (= average education).

Multiple Regression Analysis: Qualitative Information

- **Testing for differences in regression functions across groups**
- **Unrestricted model (contains full set of interactions)**

College grade point average Standardized aptitude test score High school rank percentile

$$\begin{aligned} cumgpa = & \beta_0 + \delta_0 female + \beta_1 sat + \delta_1 female sat + \beta_2 hsperc \\ & + \delta_2 female hsperc + \beta_3 tothrs + \delta_3 female tothrs + u \end{aligned}$$

Total hours spent in college courses

- **Restricted model (same regression for both groups)**

$$cumgpa = \beta_0 + \beta_1 sat + \beta_2 hsperc + \beta_3 tothrs + u$$

Multiple Regression Analysis: Qualitative Information

- **Null hypothesis**

$$H_0 : \delta_0 = 0, \delta_1 = 0, \delta_2 = 0, \delta_3 = 0$$

All interaction effects are zero, i.e. the same regression coefficients apply to men and women

- **Estimation of the unrestricted model**

$$\begin{aligned} \widehat{cumgpa} = & 1.48 - .353 \text{ female} + .0011 \text{ sat} + .00075 \text{ female} \cdot \text{sat} \\ & (.21) \quad (.411) \quad (.0002) \quad (.00039) \\ & - .0085 \text{ hisperc} - .00055 \text{ female} \cdot \text{hisperc} \\ & (.0014) \quad (.00316) \\ & + .0023 \text{ tothrs} - .00012 \text{ female} \cdot \text{tothours} \\ & (.0009) \quad (.00163) \end{aligned}$$

Tested individually, the hypothesis that the interaction effects are zero cannot be rejected

Multiple Regression Analysis: Qualitative Information

■ Joint test with F-statistic

Null hypothesis is rejected

$$F = \frac{(SSR_r - SSR_{ur})/q}{SSR_{ur}/(n - k - 1)} = \frac{(85.515 - 78.355)/4}{78.355/(366 - 7 - 1)} \approx 8.18$$

■ Alternative way to compute F-statistic

- Run separate regressions for men and for women; the unrestricted SSR is given by the sum of the SSR of these two regressions: $SSR_1 + SSR_2$
- Run regression for the restricted model and store SSR: SSR_p
- If the test is computed in this way it is called the Chow-Test
- Important: Test assumes a constant error variance accross groups

Multiple Regression Analysis: Qualitative Information

- **Alternative way to compute F-statistic**

CHOW test: $H_0 : \beta_{01} = \beta_{02}, \beta_{11} = \beta_{12}, \beta_{21} = \beta_{22}, \dots, \beta_{k1} = \beta_{k2}$

$$F = \frac{SSR_p - (SSR_1 + SSR_2)}{SSR_1 + SSR_2} \times \frac{n - 2(k + 1)}{k + 1} \sim F(k + 1, n - 2(k + 1))$$

If the restricted regression contains a dummy for an intercept shift the statistic to test changes only in the slope coefficients is:

$$H_0 : \beta_{11} = \beta_{12}, \beta_{21} = \beta_{22}, \dots, \beta_{k1} = \beta_{k2}$$

$$F = \frac{SSR_p^* - (SSR_1 + SSR_2)}{SSR_1 + SSR_2} \times \frac{n - 2(k + 1)}{k} \sim F(k, n - 2(k + 1))$$

With SSR_p^* the sum of squared residuals of the restricted regression

Multiple Regression Analysis: Qualitative Information



■ Example

Dependent Variable: LOG(WAGE)

Method: Least Squares

Included observations: 526

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	0.127998	0.105932	1.208296	0.2275
EDUC	0.090366	0.007468	12.10041	0.0000
EXPER	0.041009	0.005197	7.891606	0.0000
EXPER^2	-0.000714	0.000116	-6.163888	0.0000
Sum squared resid	103.7904			

Multiple Regression Analysis: Qualitative Information



■ **Example** Dependent Variable: LOG(WAGE)
Method: Least Squares
Sample: **IF FEMALE=1**
Included observations: **252**

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	0.266084	0.141535	1.879985	0.0613
EDUC	0.079195	0.010369	7.637838	0.0000
EXPER	0.022372	0.006664	3.356971	0.0009
EXPER^2	-0.000423	0.000148	-2.853701	0.0047

Sum squared resid 38.38393

Multiple Regression Analysis: Qualitative Information

■ **Example** Dependent Variable: LOG(WAGE)
Method: Least Squares
Sample: **IF FEMALE=0**
Included observations: **274**

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	0.157291	0.136403	1.153132	0.2499
EDUC	0.090354	0.009267	9.750477	0.0000
EXPER	0.054017	0.006743	8.011343	0.0000
EXPER^2	-0.000914	0.000150	-6.079015	0.0000

Sum squared resid 47.35130

Multiple Regression Analysis: Qualitative Information



■ Example

Dependent Variable: LOG(WAGE)

Method: Least Squares

Included observations: 526

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	0.390483	0.102210	3.820413	0.0001
FEMALE	-0.337187	0.036321	-9.283424	0.0000
EDUC	0.084136	0.006957	12.09407	0.0000
EXPER	0.038910	0.004824	8.066683	0.0000
EXPER^2	-0.000686	0.000107	-6.388842	0.0000

Sum squared resid	89.05862
-------------------	----------

Multiple Regression Analysis: Qualitative Information

The CHOW test for the example:

$$H_0 : \beta_{01} = \beta_{02}, \beta_{11} = \beta_{12}, \beta_{21} = \beta_{22}, \beta_{31} = \beta_{32}$$

$$F = \frac{103.7904 - (38.3839 + 47.3513)}{38.3839 + 47.3513} \times \frac{526 - 2(3+1)}{3+1} = 27.272 > F_{0.05}(4, 518) = 2.39$$

To test changes only in the slope coefficients in the example:

$$H_0 : \beta_{11} = \beta_{12}, \beta_{21} = \beta_{22}, \beta_{31} = \beta_{32}$$

$$F = \frac{89.0586 - (38.3839 + 47.3513)}{38.3839 + 47.3513} \times \frac{526 - 2(3+1)}{3} = 6.6931 > F_{0.05}(3, 518) = 2.62$$