



Multiple Linear Regression (Dummy Variable Treatment)

CIVL 7012/8012





In Today's Class

- Recap
- Single dummy variable
- Multiple dummy variables
- Ordinal dummy variables
- Dummy-dummy interaction
- Dummy-continuous/discrete interaction
- Binary dependent variables



THE UNIVERSITY OF

Introducing Dummy Independent Variable

- <u>Qualitative Information</u>
 - 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





Illustrative Example

Graphical Illustration

THE UNIVERSITY OF

MEMPHIS



Alternative interpretation of coefficient:

$$\delta_0 = E(wage | female = 1, educ)$$

-E(wage|female = 0, educ)

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



Specification of Dummy Variables

• Dummy variable trap This model cannot be estimated (perfect collinearity)

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

When using dummy variables, one category always has to be omitted: $wage = \beta_0 + \delta_0 female + \beta_1 educ + u \leftarrow \text{The base category are men}$ $wage = \beta_0 + \gamma_0 male + \beta_1 educ + u \leftarrow \text{The base category are women}$

Alternatively, one could omit the intercept:

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

Disadvantages: 1) More difficult to test for differences between the parameters

2) R-squared formula only valid if regression contains intercept



Interpretation of Dummy Variables

• Estimated wage equation with intercept shift



(Standard errors in parenthesis)

- 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.

Model with only dummy variables-(Example-1)

• Comparing means of subpopulations described by dummies

 $\widehat{wage} = 7.10 \leftarrow 2.51 female$ $(.21) \quad (.26)$

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

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\$.

• 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

Model with only dummy variables-(Example-2)

• Further example: Effects of training grants on hours of training



- 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?



Dependent log(y) and Dummy Independent

• Using dummy explanatory variables in equations for log(y)

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

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

$$+ .027 bdrms + .054 colonial$$

$$\xrightarrow{\text{Dummy indicating whether house is of colonial style}}$$

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

$$\Rightarrow \frac{\partial \log(price)}{\partial colonial} = \frac{\% \partial price}{\partial colonial} = 5.4\%$$
As the dummy for colonial style changes from 0 to 1, the house price increases by 5.4 percentage points

THE UNIVERSITY OF

Dummy variables for multiple categories

- 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)

$$\widehat{\log}(wage) = .321 + .213 \ marrmale \ .198 \ marrfem \ (.100) \ (.055) \ (.058) \ (.058) \ (.056) \ (.007) \ (.007) \ (.005) \ (.005) \ (.00011) \ (.00011) \ (.00011) \ (.00011) \ (.00011) \ (.00023) \ (.000$$



THE UNIVERSITY OF **MEMPHIS**.

Ordinal Dummy Variables

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



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:

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

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).

Interaction term

THE UNIVERSITY OF

Interactions among dummy variables

- Interactions involving dummy variables
- Allowing for different slopes

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

 β_0 = intercept men β_1 = slope men

 $\beta_0 + \delta_0 =$ intercept women $\beta_1 + \delta_1 =$ slope women

• Interesting hypotheses



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

The <u>whole wage equation</u> is the same for men and women









THE UNIVERSITY OF

Dummy-Continuous /Discrete Interaction (2)

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



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

All interaction effects are zero, i.e. the same regression coefficients

apply to men and women

Dummy-Continuous /Discrete Interaction (3)

• Null hypothesis

THE UNIVERSITY OF

1EMPHIS.

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

• Estimation of the unrestricted model





Craaners. Thinkers. Doers.

Models (with Dummy Variables)

• Joint test with F-statistic

Null hypothesis is rejected

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

- SSRr is the sum of squared residuals from the restricted regression, i.e., the regression where we impose the restriction.
- SSRur is the sum of squared residuals from the full model,
- q is the number of restrictions under the null and
- k is the number of regressors in the unrestricted regression.





- <u>A Binary dependent variable: the linear probability model</u>
- Linear regression when the dependent variable is binary

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + u$$
If the dependent variable only takes on the values 1 and 0
$$\Rightarrow E(y|\mathbf{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

$$E(y|\mathbf{x}) = 1 \cdot P(y = 1|\mathbf{x}) + 0 \cdot P(y = 0|\mathbf{x})$$

$$\Rightarrow P(y = 1|\mathbf{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

$$\Rightarrow \beta_j = \partial P(y = 1|\mathbf{x}) / \partial x_j$$
In the linear probability model, the coefficients describe the effect of the explanatory variables on the probability that $y=1$

THE UNIVERSITY OF **MEMPHIS**

Binary dependent variable:Example-1

• Example: Labor force participation of married women





Binary dependent variable: Example-2

• Example: Female labor participation of married women (cont.)

