

# Ch07. Multiple Regression Analysis with Qualitative Information: Binary (or Dummy) Variables

Ping Yu

HKU Business School  
The University of Hong Kong

# Describing Qualitative Information

# Quantitative and Qualitative Information

- **Quantitative** Variables: hourly wage, years of education, college GPA, amount of air pollution, firm sales, number of arrests, etc., where the magnitude of variable conveys useful information.
- **Qualitative** Variable: gender, race, industry (manufacturing, retail, finance, etc.), region (South, North, West, etc.), rating grade (A, B, C, D, F, etc), etc.
- A way to incorporate qualitative information is to use **dummy variables**.
- A dummy variable is also called a **binary variable** or a **zero-one variable**.
- Dummy variables may appear as the dependent or as independent variables.
- We consider only **independent dummy** variables.

# A Single Dummy Independent Variable

## Example: A Simple Wage Equation

- Suppose

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

where

$$female = \begin{cases} 1, & \text{if the person is a woman,} \\ 0, & \text{if the person is a man,} \end{cases} \quad \text{is a dummy variable.}$$

- $\delta_0$  is the wage gain/loss if the person is a woman rather than a man (holding other things fixed).
- Alternative interpretation of  $\delta_0$ :

$$\begin{aligned} \delta_0 &= E[wage | female = 1, educ] - E[wage | female = 0, educ] \\ &= \beta_0 + \delta_0 + \beta_1 educ - (\beta_0 + \beta_1 educ), \end{aligned}$$

i.e. the difference in mean wage between men and women with the same level of education. [\[figure here\]](#)

- Note that the mean wage difference is the same at **all** levels of education, i.e., the mean wage equations for men and women are parallel.

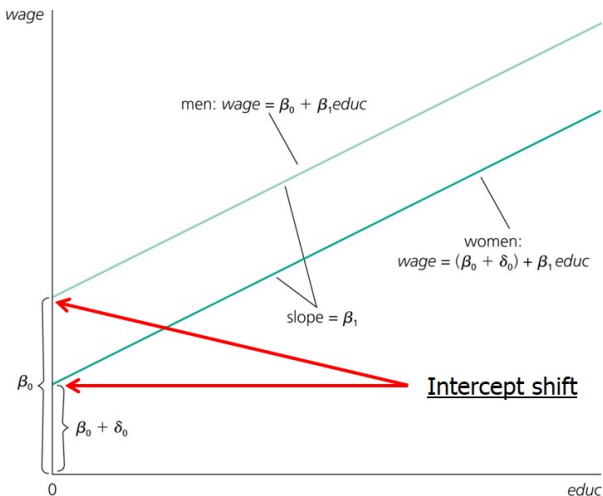


Figure: Graph of  $\text{wage} = \beta_0 + \delta_0 \text{female} + \beta_1 \text{educ}$  for  $\delta_0 < 0$

## Dummy Variable Trap

- The model

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

cannot be estimated due to perfect collinearity.

- Why?** There is an **exact** relationship among the independent variables:  
 $1 = male + female$ .
- When using dummy variables, one category always has to be omitted:

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

where men is the **base group** or **benchmark group**, i.e., the group with the dummy equal to zero/used for comparison, or

$$wage = \beta_0 + \gamma_0 male + \beta_1 educ + u,$$

where women is the base group (or category).

- Alternatively, one could omit the intercept,

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

## Disadvantage Without Intercept

- More difficult to test for differences between the parameters:

$$H_0 : \gamma_0 = \delta_0,$$

and the  $t$  statistic is

$$t = \frac{\hat{\gamma} - \hat{\delta}}{\text{se}(\hat{\gamma} - \hat{\delta})},$$

but  $\text{se}(\hat{\gamma} - \hat{\delta})$  is not available from the output of standard econometric softwares (as discussed in chapter 4).

- The  $R$ -squared formula is valid only if the regression contains an intercept: Recall that

$$R^2 = 1 - \frac{SSR}{SST} \text{ with } SST = \sum_{i=1}^n (y_i - \bar{y})^2,$$

where  $\bar{y}$  is  $\hat{\beta}_0$  in the regression

$$y = \beta_0 + u,$$

so  $SST$  can be treated as the  $SSR$  for the restricted regression of

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + u \text{ with } \beta_1 = \dots = \beta_k = 0.$$



## Example: Hourly Wage Equation with Intercept Shift

- The fitted wage equation is

$$\widehat{wage} = -1.57 - 1.81 \text{female} + .572 \text{educ} + .025 \text{exper} + .141 \text{tenure}$$

$$\begin{array}{cccccc}
 & (.72) & (.26) & (.049) & (.012) & (.021)
 \end{array}$$

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

- Holding education, experience, and tenure fixed, women earn  $\hat{\delta}_0 = \$1.81$  less per hour than men.
- Does that mean that women are discriminated against?
- Not necessarily. Being female may be correlated with other productivity characteristics (e.g., baby birth) that have not been controlled for.

## continue

- Let's compare means of subpopulations described by dummies:

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

(.21) (.30)

$$n = 526, R^2 = .116 (< .364 \text{ as expected})$$

- 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,
 
$$|t| = \left| \frac{-2.51}{.30} \right| = |-8.37| > 1.96.$$
  - 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.
  - When more factors (such as baby birth) are controlled for, then we expect  $|\widehat{\delta}_0|$  would be even smaller (until insignificance?).

## Example: Effects of Training Grants on Hours of Training

- The fitted regression line is

$$\widehat{hrsemp} = 46.67 + 26.25 grant - .98 \log(sales) - 6.07 \log(employ)$$

$$(43.41) (5.59) \quad (3.54) \quad (3.88)$$

$$n = 105, R^2 = .237$$

where

$hrsemp$  = hours training per employee, at the firm level

$grant$  = dummy indicating whether firm received training grant

$employ$  = number of employees

- This is an example of program evaluation:
  - treatment group (= grant receivers) vs. control group (= no grant).
  - $t_{grant} = 4.70 > 1.96$ , but is the effect of treatment on the outcome of interest causal? The answer depends on whether  $E[u|grant] = 0$ . It might be that to get grants, some firms give more training to their employees.

a: Using Dummy Explanatory Variables in Equations for  $\log(y)$ 

- **Example** (Housing Price Regression): The fitted regression line is

$$\begin{aligned} \widehat{\log(\text{price})} &= -1.35 + .168\log(\text{lotsize}) + .707\log(\text{sqft}) \\ &\quad (.65) \quad (.038) \quad \quad (.093) \\ &\quad + .027\text{bdrms} + .054\text{colonial} \\ &\quad (.029) \quad \quad (.045) \\ n &= 88, R^2 = .649 \end{aligned}$$

where

*colonial* = dummy for the colonial style [[figure here](#)]

- Now,

$$\frac{\partial \log(\text{price})}{\partial \text{colonial}} = \frac{\partial \text{price} / \text{price}}{\partial \text{colonial}} = 5.4\%$$

- As the dummy for colonial style changes from 0 to 1, the house price increases by 5.4 percentage points.

## American Colonial Architecture

- American colonial architecture includes several building design styles associated with the colonial period of the United States, including First Period English (late-medieval), French Colonial, Spanish Colonial, Dutch Colonial and Georgian. These styles are associated with the houses, churches and government buildings of the period from about 1600 through the 19th century.

- From Wiki



Figure: Corwin House, Salem, Massachusetts, built about 1660, First Period English

# Using 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).
- Example** (Log Hourly Wage Equation): The fitted regression line is

$$\begin{aligned} \log(\widehat{\text{wage}}) = & - .321 + .213 \text{marrmale} - .198 \text{marrfem} + \\ & (.100) (.055) \qquad (.0058) \\ & - .110 \text{singfem} + .079 \text{educ} + .027 \text{exper} - .00054 \text{exper}^2 \\ & (.056) \qquad (.007) \qquad (.005) \qquad (.00011) \\ & + .029 \text{tenure} - .00053 \text{tenure}^2 \\ & (.007) \qquad (.00023) \\ n = & 526, R^2 = .461 \end{aligned}$$

- Holding other things fixed, married women earn 19.8% less than single men (= the base category); similarly, married men earn 21.3% more and single women earn 11.0% (< 19.8%) less than single men. [[economic intuition here](#)]

## a: Incorporating Ordinal Information by Using Dummy Variables

- **Example** (City Credit Ratings and Municipal Bond Interest Rates): We can consider two specifications of the regression line.
- The first specification is

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

where

$MBR$  = municipal bond interest rate

$CR$  = credit rating from 0 – 4 (0 = worst, 4 = best)

- 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:

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

where  $CR_1, \dots, CR_4$  are 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).



## Difference Between These Two Specifications

- Specification 1:

$$CR = 0 \implies MBR = \beta_0,$$

$$CR = 1 \implies MBR = \beta_0 + \beta_1,$$

$$CR = 2 \implies MBR = \beta_0 + 2\beta_1,$$

$$CR = 3 \implies MBR = \beta_0 + 3\beta_1,$$

$$CR = 4 \implies MBR = \beta_0 + 4\beta_1,$$

where the increase in  $MBR$  for each rating improvement is the same -  $\beta_1$ .

- Specification 2:

$$CR = 0 \implies MBR = \beta_0,$$

$$CR = 1 \implies MBR = \beta_0 + \delta_1,$$

$$CR = 2 \implies MBR = \beta_0 + \delta_2,$$

$$CR = 3 \implies MBR = \beta_0 + \delta_3,$$

$$CR = 4 \implies MBR = \beta_0 + \delta_4,$$

where the increase in  $MBR$  for each rating improvement can be different due to the arbitrariness of  $\delta_1, \dots, \delta_4$ .

# Interactions Involving Dummy Variables

## a: Interactions among Dummy Variables

- Reconsider the female and marital status effect on  $\log(\text{wage})$  by adding the *female* · *married* interaction term:

$$\log(\widehat{\text{wage}}) = - .321 - .110\text{female} + .213\text{married} - .301\text{female} \cdot \text{married} + \dots$$

$$(.100) \quad (.056) \quad (.055) \quad (.072)$$

- These two specifications are equivalent: four categories are generated. [\[2 × 2 table in the next slide\]](#)
- marmale*: setting *married* = 1 and *female* = 0, we get  $\widehat{\delta}_2 = .213$  as before.
- marrfem*: setting *married* = 1 and *female* = 1, we get  $\widehat{\delta}_1 + \widehat{\delta}_2 + \widehat{\delta}_3 = -.110 + .213 - .301 = -.198$  as before.
- singfem*: setting *married* = 0 and *female* = 1, we get  $\widehat{\delta}_1 = -.110$  as before.

(\*) What is the Meaning of the Coefficient of *female* · *married*,  $\hat{\delta}_3$ ?

- Four Categories:

	Female	Male
Single	$\delta_1$	0
Married	$\delta_1 + \delta_2 + \delta_3$	$\delta_2$

- DID:

$$\begin{aligned}
 & (E[\log(\text{wage}) | \text{female} = 1, \text{married} = 1] - E[\log(\text{wage}) | \text{female} = 1, \text{married} = 0]) \\
 & - (E[\log(\text{wage}) | \text{female} = 0, \text{married} = 1] - E[\log(\text{wage}) | \text{female} = 0, \text{married} = 0]) \\
 = & [(\delta_1 + \delta_2 + \delta_3) - \delta_1] - [\delta_2 - 0] \\
 = & \text{difference (in gender) in difference (in marriage)} \\
 = & [(\delta_1 + \delta_2 + \delta_3) - \delta_2] - [\delta_1 - 0] \\
 = & \text{difference (in marriage) in difference (in gender)}
 \end{aligned}$$

- Analog:

$$\frac{\partial^2 y}{\partial x_1 \partial x_2}$$

## b: Allowing for Different Slopes

- Consider the model

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

where

$$\begin{aligned} \beta_0 &= \text{intercept of men, } \beta_1 = \text{slope of men,} \\ \beta_0 + \delta_0 &= \text{intercept of women, } \beta_1 + \delta_1 = \text{slope of women.} \end{aligned}$$

- Interacting both the intercept and the slope with the female dummy enables one to model completely independent wage equations for men and women. [[figure here](#)]
- Interested Hypotheses:**

$$H_0 : \delta_1 = 0,$$

i.e., the return to education is the same for men and women, and

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

i.e., the whole wage equation is the same for men and women.

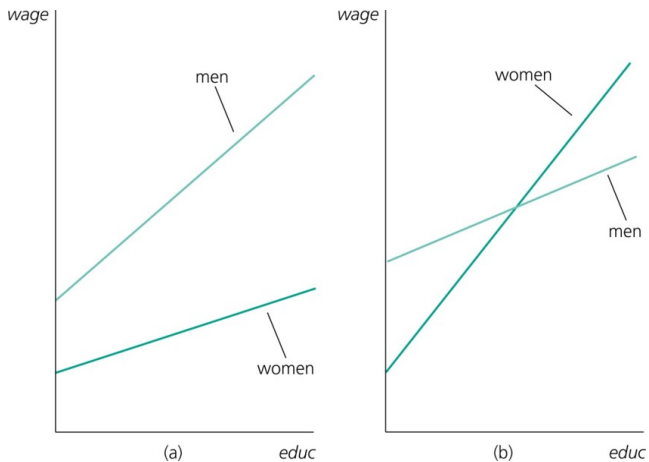


Figure: (a)  $\delta_0 < 0, \delta_1 < 0$ ; (b)  $\delta_0 < 0, \delta_1 > 0$

## Example: Log Hourly Wage Equation

- The fitted regression line is

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

- $|t_{\text{female-educ}}| = \left| \frac{-.0056}{.0131} \right| = |-0.43| < 1.96$ : No evidence against hypothesis that the return to education is the same for men and women.
- $|t_{\text{female}}| = \left| \frac{-.227}{.168} \right| = |-1.35| < 1.96$ : 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  $\text{educ} = 0$  since

$$\frac{\partial \log(\text{wage})}{\partial \text{female}} = -.227 - .0056 \text{educ}.$$

- To answer the question one has to recenter the interaction term, e.g., around  $educ = 12.5$  (= average education) to have  $female \cdot (educ - 12.5)$ :

$$\frac{\partial \log(\text{wage})}{\partial \text{female}} = \hat{\delta}_0 + \hat{\delta}_1 (\text{educ} - 12.5) \text{ with new } \hat{\delta}_0 = -.297$$

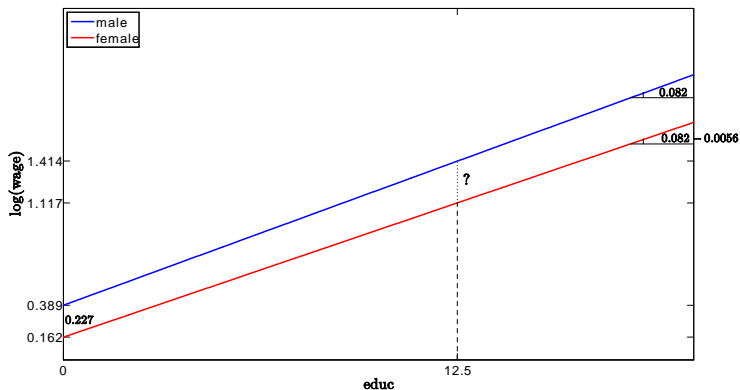


Figure: The New new  $\hat{\delta}_0 = -.227 + 12.5 \times (-.0056) = -.297 < -.227$



## c: Testing for Differences in Regression Functions across Groups

- This is a **special**  $F$  test with the unrestricted model containing full set of interactions,

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

and the restricted model with same regression for both groups,

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

where

$cumgpa$  = college GPA

$sat$  = standardized aptitude test score

$hsperc$  = high school rank percentile

$tothrs$  = total hours spent in college courses

- The null hypothesis is

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

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

## Estimation of the Unrestricted Model

- The estimated unrestricted model is

$$\widehat{cumgpa} = 1.48 - .353female + .0011sat + .00075female \cdot sat$$

$$(.21)(.411) \quad (.0002) \quad (.00039)$$

$$- .0085hsperc - .00055female \cdot hsperc + .0023tothrs$$

$$(.0014) \quad (.00316) \quad (.0009)$$

$$- .00012female \cdot tothrs$$

$$(.00163)$$

$$n = 366, R^2 = .406, \bar{R}^2 = .394$$

- It can be shown that [proof not required]

$$SSR_{ur} = SSR_{male} + SSR_{female},$$

where  $SSR_{male}$  is the SSR in the regression

$$cumgpa = \beta_0 + \beta_1sat + \beta_2hsperc + \beta_3tothrs + u,$$

using only the data of male, and  $SSR_{female}$  is the SSR using only the data of female.

## Testing Results

- Tested individually, the hypothesis that the interaction effects are zero cannot be rejected.
- Tested jointly, the  $F$  statistic is

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

and the null is rejected.

- $SSR_{ur} = SSR_{male} + SSR_{female} = 58.752 + 19.603 = 78.355$ ,  $n_{male} = 276$ ,  $n_{female} = 90$  and  $n = 366$ .
- This relationship is true only if **all** interaction terms are included in the unrestricted model.

## Chow Test

- In general, if the test is computed in this way, i.e.,

$$F = \frac{[SSR_P - (SSR_1 + SSR_2)] / (k + 1)}{(SSR_1 + SSR_2) / [n - 2(k + 1)]},$$

it is called the **Chow-Test**. [\[photo here\]](#)

- $SSR_P$  is the SSR for the pooled (restricted) regression; [see below]
- $SSR_1$  and  $SSR_2$  are the SSRs for the two separate regressions; [see below]
- The number of restrictions is  $k + 1$  with  $k$  being the number of nonconstant regressors in the **pooled** regression; [see below]
- The total number of parameters in the unrestricted model is  $2(k + 1)$ . [see below]
- $SSR_P$ :  $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i, i = 1, \dots, n.$
- $SSR_{ur}$ :  $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \delta_0 D_i + \delta_1 D_i x_{i1} + \dots + \delta_k D_i x_{ik} + u_i, i = 1, \dots, n.$
- $SSR_1$ :  $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i, i = 1, \dots, n_1.$
- $SSR_2$ :  $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + u_i, i = n_1 + 1, \dots, n.$
- $H_0: \delta_0 = \delta_1 = \dots = \delta_k = 0.$



Gregory C. Chow (1929-),  
Princeton, 1955ChicagoPhD

- Chow, G.C., 1960, Tests of Equality Between Sets of Coefficients in Two Linear Regressions, *Econometrica*, 28, 591-605.
- **Caution:** Chow-Test assumes a constant error variance across groups as assumed in the  $F$  test.
  - $E[u_i|D_i] = \sigma^2$  for  $D_i = 0$  and 1.

# Applications of the Chow-Test

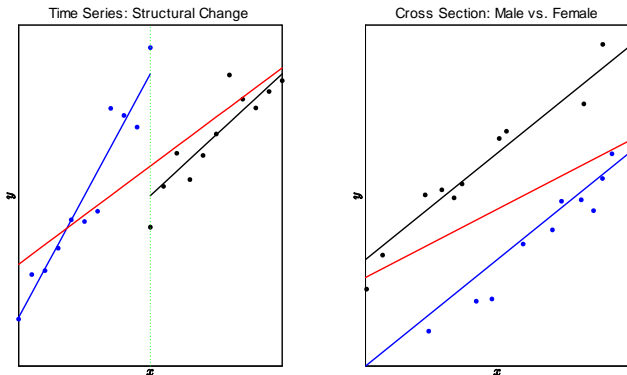


Figure: Restricted and Unrestricted Models in the Chow Test

- You must pass the Chow test to pass this course!

## (\*\*) A Binary Dependent Variable: The Linear Probability Model

## Linear Regression When the Dependent Variable is Binary

- Suppose

$$y = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k + u, \quad (1)$$

with  $E[u|\mathbf{x}] = 0$ ; then

$$E[y|\mathbf{x}] = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k.$$

- If the dependent variable  $y$  only takes on the values 1 and 0, then

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

i.e.,

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

where  $P(y = 1|\mathbf{x})$  is called the **response probability**.

- Since  $P(y = 1|\mathbf{x})$  is linear in  $\mathbf{x}$ , the model (1) is called the **linear probability model (LPM)**.
- In the LPM,

$$\beta_j = \partial P(y = 1|\mathbf{x}) / \partial x_j$$

describes the effect of the explanatory variable  $x_j$  on the probability that  $y = 1$ .



## Example: Labor Force Participation of Married Women

- The fitted regression line is

$$\widehat{infl} = .596 - .0034nwifeinc + .038educ + .039exper$$

$$\begin{array}{cccc}
 (.154) & (.0014) & (.007) & (.006) \\
 -.00060exper^2 - .016age - .262kidslt6 + .0130kidsge6 \\
 (.00018) & (.002) & (.0034) & (.0132)
 \end{array}$$

$$n = 753, R^2 = .264$$

where

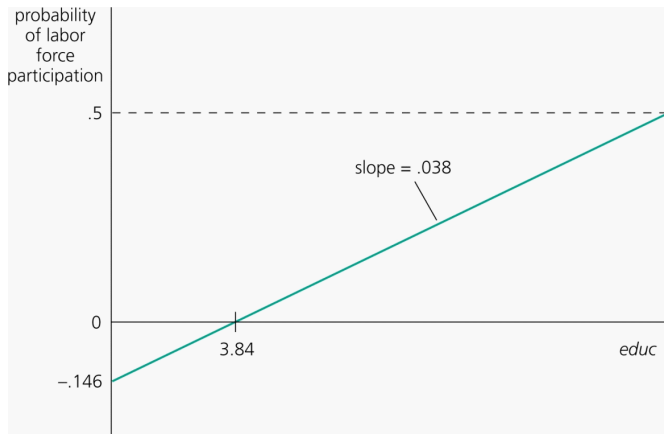
*infl* = dummy for "in the labor force" of a married women

*nwifeinc* = husband's earnings (in thousands of dollars)

*kidslt6* = number of children less than six years old

*kidsge6* = number of kids between 6 and 18 years of age

- All variables except *kidsge6* are statistically significant, and all of the significant variables have correct signs.
- If the number of kids under six years increases by one, the probability that the woman works falls by 26.2%.

Graph for  $nwifeinc=50$ ,  $exper=5$ ,  $age=30$ ,  $kindslt6=1$ ,  $kidsge6=0$ 

- The maximum level of education in the sample is  $educ = 17$ . For the given case, this leads to a predicted probability to be in the labor force of about 50%.
- Negative predicted probability but no problem because no woman in the sample has  $educ < 5$ .

## Disadvantages and Advantages of the LPM

### Disadvantages:

- Predicted probabilities may be larger than one or smaller than zero.
- Marginal probability effects sometimes logically impossible, e.g., the first small child would reduce the probability by a large amount, but subsequent children would have a smaller marginal effect; going from zero to four young children reduces the probability of working by  $\Delta \widehat{infl} = .262 \times 4 = 1.048!$
- The LPM is necessarily heteroskedastic:

$$\text{Var}(y|\mathbf{x}) = P(y = 1|\mathbf{x})(1 - P(y = 1|\mathbf{x}))$$

by the variance formula of a Bernoulli random variable.

- Heterosceasticity consistent standard errors need to be computed.

### Advantages:

- Easy estimation and interpretation.
- Estimated effects and predictions often reasonably good in practice.