

1)

4 Testing Hypotheses about a Single Linear Combination of the Parameter  $\beta_1, \beta_2$  (non-linear)

Consider

$\log(wage) = \beta_0 + \beta_1 jc + \beta_2 univ + \beta_3 exper + u$

where  $jc$  = number of years attending a two-year college  
 $univ$  = number of years at a four-year college  
 $exper$  = months in the workforce.

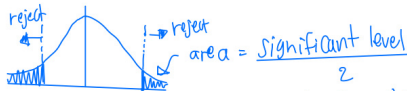
We want to test whether  $\beta_1 = \beta_2$ .

$H_0 : \beta_1 = \beta_2 \Rightarrow H_0 : \beta_1 - \beta_2 = 0$

Against

$H_a : \beta_1 \neq \beta_2 \Rightarrow H_a : \beta_1 - \beta_2 \neq 0$

2-tailed test



$t = \frac{(\hat{\beta}_1 - \hat{\beta}_2) - 0}{S.E.(\hat{\beta}_1 - \hat{\beta}_2)}$

We compute this t-statistic and compare with the critical value

Where  $S.E.(\hat{\beta}_1 - \hat{\beta}_2) = \sqrt{Var(\hat{\beta}_1 - \hat{\beta}_2)}$   
 $= \sqrt{Var(\hat{\beta}_1) + Var(\hat{\beta}_2) - 2Cov(\hat{\beta}_1, \hat{\beta}_2)}$

not very straight forward to calculate  
 we use a variable transformation  
 trick see notes!!!

2)

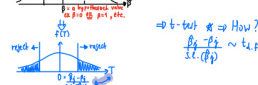
Inference  $\rightarrow$  Hypothesis testing about  $\beta$  the true parameter.

$Wage = \beta_0 + \beta_1 educ + \beta_2 experience + \dots + u$

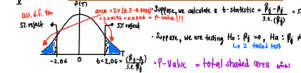
We want to test hypothesis about the true impact of each X variables (educ, experience) on the independent variable (Y)

BUT, We don't know what the true  $\beta$  are, so we use  $\hat{\beta}$  (estimator) and S.E. ( $\hat{\beta}$ ) to test the hypothesis

Test if  $\beta =$  some number  
 eg.  $\beta_1 = 0 \rightarrow X_1$  has no impact on Y.  
 $\beta_1 = 1 \rightarrow 1$  month in X<sub>1</sub> correspond to 1 unit in Y.



\* (Significant level = total area in the rejection region)  $\alpha$



\*  $\alpha$  = Significant level which we will reject the  $H_0$  and risk that we will reject  $H_0$ .  
 If  $\alpha$  value < significant level  $\rightarrow$  always reject  $H_0$ !!!

3)

another possible hypothesis test (one-tailed alternative)

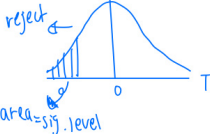
$H_0 : \beta_1 = \beta_2 \Rightarrow H_0 : \beta_1 - \beta_2 = 0$

$H_a : \beta_1 < \beta_2 \Rightarrow H_a : \beta_1 - \beta_2 < 0$

It is assumed that  $\beta_1$  would not be more than  $\beta_2$   
 (returns to a 2-year college would never be more than returns to University education)

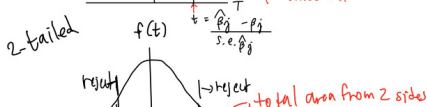
$t = \frac{(\hat{\beta}_1 - \hat{\beta}_2) - 0}{S.E.(\hat{\beta}_1 - \hat{\beta}_2)}$

\* Then, go to the extra note



5 Computing p-Values for t-Tests

What is the significance level given the computed t-statistics?



p-value:  $P(|T| > |t|)$   
 $T = t$ -distributed random variable with d. f. =  $n - k - 1$   
 $t =$  computed t-statistic.

$\rightarrow$  P-value = probability that a random T value will be greater (in the 1 term) than our T in the H<sub>0</sub> test

4)

In-class exercise

Consider the multiple regression model, assume MLR 1-6 are satisfied.

$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$

You would like to test  $H_0 : \beta_1 - 3\beta_2 = 1$   
 $H_a$ : otherwise is true

write the t-statistic for testing  $H_0$

$t = \frac{(\hat{\beta}_1 - 3\hat{\beta}_2) - 1}{S.E.(\hat{\beta}_1 - 3\hat{\beta}_2)}$

Define  $\theta_1 = \hat{\beta}_1 - 3\hat{\beta}_2 \Rightarrow H_0 : \theta_1 = 1, H_a : \theta_1 \neq 1$   
 $t = \frac{\hat{\theta}_1 - 1}{S.E.(\hat{\theta}_1)}$   
 we need our regression to have  $\theta_1$  in it. So OLS estimation will automatically give  $\hat{\theta}_1$  & S.E.  $\hat{\theta}_1$

Now,  $\hat{\beta}_1 = \hat{\theta}_1 + 3\hat{\beta}_2$   
 $\beta_1 = \theta_1 + 3\beta_2$

Substitute in the main regression & get

$Y = \beta_0 + (\theta_1 + 3\beta_2)X_1 + \beta_2 X_2 + \beta_3 X_3 + u$   
 $= \beta_0 + \theta_1 X_1 + 3\beta_2 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$   
 $= \beta_0 + \theta_1 X_1 + \beta_2 (X_2 + 3X_1) + \beta_3 X_3 + u$

Now, the explanatory variables are going to be  $X_1, X_2 + 3X_1$  &  $X_3$

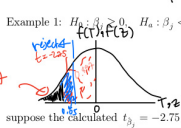
We can calculate  $t = \frac{\hat{\theta}_1 - 1}{S.E.(\hat{\theta}_1)}$

5)

for z-table

Example 1:  $H_0: \beta_j \geq 0, H_a: \beta_j < 0, d.f. = 140$ .  $\rightarrow$  z-table  
 $\rightarrow$  p-value = what should be the significant level, given the critical value of -2.75?  $\rightarrow$  find the shaded area

$0.5 - 0.997$   
 $\Rightarrow 0.003$



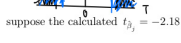
suppose the calculated  $t_{\beta_j} = -2.75 \rightarrow t_{\beta_j} = \frac{\hat{\beta}_j - \beta_j}{s.e.(\hat{\beta}_j)}$

From the z-table, the value -2.75 corresponds to area = 0.003

Thus, p-value = 0.003

Would we reject  $H_0$  if we use the significance level = 5%? Yes.  
~~X rule!~~ We reject  $H_0$  if p-value < sig-level

Example 2:  $H_0: \beta_j = a_j, H_a: \beta_j \neq a_j, d.f. = 18$ .  $\rightarrow$  use t-table



suppose the calculated  $t_{\beta_j} = -2.18$

From the t-table, the value -2.18 corresponds to area = 0.02 to 0.05

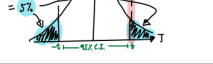
Thus, p-value = is between 0.02 - 0.05

Would we reject  $H_0$  if we use the significance level = 5%? Yes, reject  $H_0$  bco the area is less than 0.05 or p-value < 0.05

6 Confidence Intervals (CI)

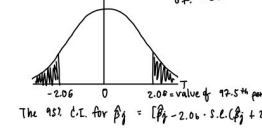
Confidence Intervals for the POPULATION PARAMETER ( $\beta_j$ )  
 The range of values that would capture the true  $\beta_j$  at a 5% chance

A 95% CI of  $\beta_j$  is given by  $\hat{\beta}_j \pm 2 \cdot s.e.(\hat{\beta}_j)$   
 CI  $\Rightarrow \hat{\beta}_j \pm 2 \cdot s.e.(\hat{\beta}_j)$   
 CI is the 95% percentile in the t-distribution with  $n-k-1$  d.f.

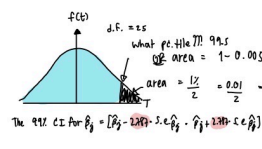


6)

Example 1: 95% CI  $f(t)$   $d.f. = 25$



Example 2: 99% CI



for power

F-test motivation

We want to test the significance of the group of hypothesis (Multiple Hypothesis)

Grade 325 =  $\beta_0 + \beta_1 \cdot \text{times\_front} + \beta_2 \cdot \text{times\_back} + \beta_3 \cdot \text{times\_study} + \beta_4 \cdot \text{past\_GPA} + \beta_5 \cdot \text{gender} + u$

$H_0: \text{seat position doesn't have impact on GPA}$   
 $\beta_1 = 0 \ \& \ \beta_2 = 0 \Rightarrow \beta_1 = \beta_2 = 0$

$H_a: \text{seat position matters}$   
 $\beta_1 \neq 0 \ \& \ \beta_2 \neq 0$

OR  $\beta_1 \neq 0 \ \& \ \beta_2 = 0$  at least one of the  $\beta_1, \beta_2 \neq 0$   
 OR  $\beta_1 = 0 \ \& \ \beta_2 \neq 0$

7)

7 Testing Multiple Linear Restrictions: The F-test

Suppose the model is specified by

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + u$$

$H_0: \beta_1 = 0 \ \& \ \beta_2 = 0$   $\rightarrow$  want to test if  $x_1$  &  $x_2$  BOTH have no impact on  $y$   
 $H_a, H_1: H_0$  is not true

We can use the F-test to test this type of "multiple hypotheses".

1. Our full model is called the "unrestricted" model (ur). Suppose it can be expressed as:  
 $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + u$  is true  $\Rightarrow$  reject  $H_0$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_k x_k + u$$

2. The model which takes out  $x$  (which we think its associated  $\beta = 0$ ) is called the restricted model (r).  $\rightarrow$  small model  
 $Y = \beta_0 + \beta_1 x_1 + u$  is true  $\Rightarrow$  do not reject  $H_0$

\* Suppose there are "q" no. of  $\beta$  that we would like to perform a joint-test of  $= 0$

e.g. in this model  $q = 2$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

$H_0: \beta_{k-q+1} = \beta_{k-q+2} = \dots = \beta_k = 0$

(the last q  $\beta_j = 0$ )  
 $H_a: H_0$  is not true.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{k-q} x_{k-q} + \beta_{k-q+1} x_{k-q+1} + \beta_{k-q+2} x_{k-q+2} + \dots + \beta_k x_k + u$$

(r) (ur) unrestricted model

$$F = \frac{(SSR_r - SSR_{ur}) / q}{SSR_{ur} / (n-k-1)}$$

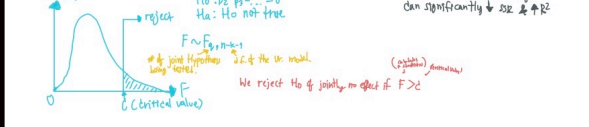
This is always  $> 0$  bco  $SSR_r < SSR_{ur}$ . Every time you add 1 more  $x_j$ , the model will be better explained.  
 d.f. of the "ur" model.

8)

So, if every time you add 1 more  $x$  variable, the SSR  $\downarrow$  and  $R^2 \uparrow$ , why don't we just keep the additional  $x$  in the model??

Because everytime we add 1 more  $x$ ,  $Var(\hat{\beta}_j)$  will increase, making the prediction of  $\beta$  less precise. So, we only keep the additional  $x$ , if it / they can improve the model enough.

can't SSR ( $\uparrow R^2$ ) enough can significantly  $\downarrow$  SSR  $\& \ \uparrow R^2$



3. Some useful facts

①  $R^2_{ur} > R^2_r$  because any additional X will increase  $R^2$  (improve fit)  
 $\Rightarrow SSR_{ur} < SSR_r$

② By including more  $X_j$  the model is certainly better explained. However, we would like to reject  $H_0$  if the inclusion of extra variables does not improve the model enough.

4. Other ways to calculate the F-statistics:

$\Rightarrow$  from  $R^2 = 1 - \frac{SSR}{SST}$

(Now) we have  $F = \frac{(R^2_{ur} - R^2_r) / (k_2 - k_1)}{(1 - R^2_{ur}) / (n - k_2)}$

will we need to test the overall significance of the model?  
 $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$   
 $H_a: \text{otherwise}$   
 $F = \frac{R^2 / k}{(1 - R^2) / (n - k)}$

Example: Suppose we are interested in understanding the determinant of a baseball player's salary.

- salary = season salary
- years = years in major leagues
- gamesyr = games per year in the league
- bavg = career batting average
- hrunsyr = homers per year
- rbsyr = runs batted in per year

If we want to test whether performance has any impact on salary.  
 $H_0: \beta_{\text{gamesyr}} = \beta_{\text{bavg}} = \beta_{\text{hrunsyr}} = \beta_{\text{rbsyr}} = 0$   
 $H_a: \text{otherwise}$

the unrestricted model (ur) is defined by

UR Model

```
regress log_salary years gamesyr bavg hrunsyr rbsyr
```

Source	SS	df	MS	Number of obs = 353
Model	308.989208	5	61.7978416	F( 5, 347) = 117.06
Residual	183.186329	347	.52784487	Prob > F = 0.0000
Total	492.175535	352	1.39822595	R-squared = 0.6278
				Adj R-squared = 0.6224
				Root MSE = .72658

log_salary	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
years	.0688626	.0121145	5.68	0.000	.0450355 .0926898
gamesyr	-.0125521	.0026468	-4.74	0.000	-.0078464 -.0172578
bavg	-.0009786	.0011035	-0.89	0.376	-.0031818 .0012046
hrunsyr	.0144295	.016057	0.90	0.369	-.0171518 .046107
rbsyr	-.0107637	.007175	-1.50	0.134	-.0239462 .0024176
_cons	11.12942	.288229	38.75	0.000	10.62433 11.76048

the restricted model (r) is defined by

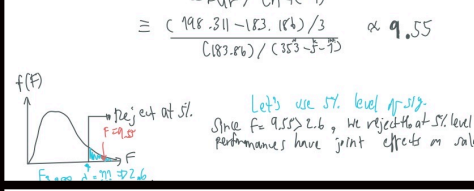
```
regress log_salary years gamesyr
```

Source	SS	df	MS	Number of obs = 353
Model	293.864058	2	146.932029	F( 2, 350) = 259.32
Residual	198.311477	350	.566604221	Prob > F = 0.0000
Total	492.175535	352	1.39822595	R-squared = 0.5971
				Adj R-squared = 0.5949
				Root MSE = .75273

Now, our  $H_0$  and  $H_a$  becomes

$$F = \frac{(SSR_r - SSR_{ur}) / (k_2 - k_1)}{SSR_{ur} / (n - k_2)}$$

$$= \frac{(198.311477 - 183.186329) / (350 - 5)}{(183.186329) / (353 - 5 - 7)} \approx 9.55$$



8 How the Hypothesis Testing is done in Practice

1. Check the values of  $t$ -statistic reported by the statistical software (i.e. STATA, SPSS, SAS)

- $\Rightarrow$  These  $t$ -statistics are to test  $H_0: \beta_1 = 0$
- $\Rightarrow$  If the d.f. > 30, then when  $t > 1.96$ , we can reject  $H_0$
- $\Rightarrow$  When  $t > 1.96$ , we can say that  $\beta_1$  is statistically significant at 5% level. (value of  $\beta_1 \neq 0$ )
- $\Rightarrow$  When  $t < 1.96$  we can say that  $\beta_1$  is not statistically significant at 5% level.
- $\Rightarrow$  If  $t < 1.96$  we can drop  $x_i$  from the model
- $\Rightarrow$  After we drop  $x_i$ , we estimate the new regression function and obtain a new set of  $\beta$ .

2. We can also perform other hypothesis testings of interest.

- e.g.  $H_0: \beta_1 = \beta_2$
- or  $H_0: \beta_1 = 5$  etc.
- or perform an F-test for testing multiple linear restrictions

3. Usually, in economics, the estimation results are reported using this form

Dependent Variable: log(salary)			
Independent Variables	(1)	(2)	(3)
log(sales)	.224 (.027)	.158 (.040)	.188 (.040)
log(mktval)	—	.112 (.050)	.100 (.049)
profmarg	—	-.0023 (.0022)	-.0022 (.0021)
ceoten	—	—	.0171 (.0055)
comten	—	—	-.0092 (.0033)
intercept	4.94 (0.20)	4.62 (0.25)	4.57 (0.25)
Observations	177	177	177
R-squared	.281	.304	.353

Other company performance  
 CEO characteristics

like a simple regression with 2X

Multiple Regression Analysis : Further Issues

1 Data scaling on OLS statistics

When we change the unit of measurement of a variable, the value of estimators would change accordingly. For example

$$bweight = \beta_0 + \beta_1 cigs + \beta_2 faminc$$

where  $bweight$  = child birth weight, in grams.  
 $cigs$  = number of cigarettes smoked by the mother while pregnant, per day.  
 $faminc$  = annual family income, in thousands of dollars.

What if we use  $bweight$  in kilograms?  
 $bweight_{kg} = \frac{\beta_0}{1000} + \beta_1 cigs + \frac{\beta_2}{1000} faminc$   
 $\Rightarrow \alpha_0 = \frac{\beta_0}{1000}, \alpha_1 = \beta_1, \alpha_2 = \frac{\beta_2}{1000}$

What if we use  $faminc$  in USD (instead of 1000 USD)?  
 $bweight = \beta_0 + \beta_1 cigs + \beta_2 faminc_{USD}$   
 $\Rightarrow \beta_1 = \frac{\beta_2}{1000}$

in other words  $\beta_2$  = Impact of 1 USD in income.  
 $\beta_2 = \frac{1}{1000} USD$

What if we use  $bweight$  in kg & income in THB?  
 $bweight_{kg} = \frac{\beta_0}{1000} + \beta_1 cigs + \left(\frac{\beta_2}{1000}\right) faminc_{THB}$   
 This value is going to be 3000 times greater than before.

13)

2.4 More on Contour

2 More on functional forms

- Logarithmic Functional Form

usually means natural log

$$\log(Y) = \beta_0 + \beta_1 \log(X_1) + \beta_2 X_2 + u$$

$$\beta_1 = \frac{d \log(Y)}{d \log(X_1)} = \frac{\frac{1}{Y} dY}{\frac{1}{X_1} dX_1} = \frac{dY}{Y} \cdot \frac{X_1}{dX_1} = 100 \times \frac{\Delta Y}{Y} \cdot \frac{X_1}{\Delta X_1} = \frac{\% \Delta Y}{\% \Delta X_1}$$

with the log Y & log X formula, the coefficient is going to be the elasticity! (X always)

$$\beta_2 = \frac{d \log(Y)}{d X_2} = \frac{\frac{1}{Y} dY}{d X_2} = \frac{1}{Y} \frac{dY}{d X_2}$$

if we want the upper term to be % change, then

$$100 \beta_2 = \frac{100 \frac{dY}{d X_2}}{Y}$$

$$100 \beta_2 = \frac{\% \Delta Y}{\Delta X_2}$$

∴ 100 β<sub>2</sub> = % Δ in Y given that X<sub>2</sub> increases by 1 Unit.

- Models with Quadratics (Squares)

⇒ Capture increasing/decreasing marginal effects (slope of the relationship btw X & Y is not constant)

Cost-It example  
Y (in \$) vs X (in \$)

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + u$$

$$\frac{dY}{dX} = \beta_1 + 2\beta_2 X$$

decreasing returns  
⇒ Y = β<sub>0</sub> + β<sub>1</sub>X + β<sub>2</sub>X<sup>2</sup> + u  
 $\frac{dY}{dX} = \beta_1 + 2\beta_2 X$   
 F.O.C.  $\frac{\partial \pi}{\partial X} = 0 = 90 - 2Q$  ⇒ Q = 45  
 P<sub>1</sub> is positive

Example: Effects of Pollution on Housing Prices

$$\log(\text{price}) = \beta_0 + \beta_1 \log(\text{nox}) + \beta_2 \log(\text{dist}) + \beta_3 \text{rooms} + \beta_4 \text{room}^2 + \beta_5 \text{stratio} + u$$

14)

price = housing price  
nox = level of pollution  
dist = distance from downtown  
rooms = number of rooms  
stratio = average student per teacher ratio

The estimation result is given by

Source	SS	df	MS	F	Prob > F
Model	51.4933152	5	10.298663	1.172429	0.3000
Residual	33.0889098	500	.06617782		
Total	84.582225	505	.16748954		

Log(price)	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
log(nox)	-.9787545	-.0959398	-9.81	0.000	-1.172429 - .780886
dist	-.021972	-.0094013	-3.42	0.001	-.030668 - .0132764
rooms	-.528052	.1612965	-3.27	0.001	-.8697056 - .230007
room <sup>2</sup>	.0624697	-.0124867	5.00	0.000	-.0179386 .0800025
stratio	-.0486667	.0058131	-8.37	0.000	-.0600679 -.0372455
_cons	13.39154	.5692901	24.95	0.000	12.4813 14.7018

Log(price)

all variables are significant

Consider the effect of "room"

$$\frac{d \log(\text{price})}{d \text{rooms}} = \beta_3 + 2\beta_4 \text{rooms} = -0.528 + 2(0.062) \cdot \text{rooms}$$

at how many rooms does 1 additional room have a positive impact on log(price)?

$$0 = -0.528 + 2(0.062) \cdot \text{rooms}$$

rooms = 4.2

Answer ⇒ at 4.2 rooms or more  
at 5 rooms or more

What would be the % change in price when the number of room increases from 5 to 6?

$$\frac{d \log(\text{price})}{d \text{rooms}} = -0.528 + 2(0.062) \cdot \text{rooms}$$

$$100 \cdot \frac{d \text{price}}{\text{price}} = 100(-0.528 + 2(0.062) \cdot 5)$$

$$= 100 \times 0.064 = 6.4\% \text{ increase}$$

what about % increase when room increases from 5 to 7?

$$\% \Delta \text{Price} = 100(-0.528 + 2(0.062) \cdot 6) = 19.1\%$$

15)

3 Models with Interaction Terms

Consider

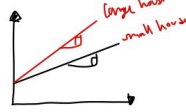
$$\text{price} = \beta_0 + \beta_1 \text{sqft} + \beta_2 \text{bdrms} + \beta_3 \text{sqft} \times \text{bdrms} + \beta_4 \text{bthrms} + u$$

where

price = housing price  
sqft = house size (square feet)  
bdrms = number of bedrooms  
bthrms = number of bathrooms

$\frac{\partial \text{price}}{\partial \text{bdrms}} = \beta_2 + \beta_3 \text{sqft}$

⇒ if β<sub>2</sub> > 0 then, an additional bedroom would increase price more for a larger house!



16)

4 More on the Goodness-of-Fit and Selection of Regressors

- Adding more regressors ALWAYS improve fit → R<sup>2</sup> always ↑
- But we lose the "degree of freedom" (d.f. = free data points used to estimate the parameter) ⇒ 1 data point is sacrificed every time we estimate a parameter.
- Using R<sup>2</sup> would not punish "having too many regressors".
- We use adjusted-R<sup>2</sup> or R<sup>2</sup> when we want to punish adding too many regressors.

$$R^2 = 1 - \frac{SSR}{SST} = 1 - \frac{SSR/k}{SST/k}$$

$$\text{adj. } R^2 = 1 - \frac{SSR/(n-k-1)}{SST/(n-1)}$$

If we have more k, d.f. = n - k - 1 ↓, SSR/(n-k-1) ↑, adj. -R<sup>2</sup> ↓

(∴ adj. -R<sup>2</sup> ↓ ⇒ we analyse additional no. of k) ⇒ depend on

Using adjusted R-squared to choose between non-nested models (one model is not a subset of another).

Consider Model 1

$$\widehat{\text{salary}} = 830.63 + 0.0163 \text{sales} + 19.63 \text{roe}$$

(223.90) (0.0089) (11.08)

n = 209, R<sup>2</sup> = 0.029, R<sup>2</sup> = 0.020

Consider Model 2

$$\log(\widehat{\text{salary}}) = 4.36 + 0.2751 \log(\text{sales}) + 0.0179 \text{roe}$$

(0.29) (0.033) (0.004)

n = 209, R<sup>2</sup> = 0.282, R<sup>2</sup> = 0.276

∴ 19.5% of variation in Y is explained. So, this model is better!!

17)

Multiple Regression Analysis with Qualitative Information:

- Outline
  - Describing qualitative information
  - Using a single dummy independent variable
  - Using dummy variables for multiple categories
  - Interactions involving dummy variables
  - A binary dependent variable (Y variable): The linear probability model
- Describing Qualitative Information
  - "Female" and "Married" are qualitative variables.
  - We arbitrarily assign a dummy variable to describe them.

$$female = \begin{cases} 1 & \text{if female} \\ 0 & \text{otherwise (or if male)} \end{cases}$$

$$married = \begin{cases} 1 & \text{if married} \\ 0 & \text{otherwise (of if single)} \end{cases}$$

TABLE 7.1 A Partial Listing of the Data in WAGE1.RAW

person	wage	educ	exper	female	married
1	3.10	11	2	0	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

This is page 1  
Printer: Opac

10)

8. Multiple Regression Analysis with Qualitative Information:

3 Models with a single dummy independent variable

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

where  $female = \begin{cases} 1 & \text{if female} \\ 0 & \text{otherwise (or if male)} \end{cases}$

In this case, the  $\delta_0$  notation is used to highlight the interpretation of the parameters multiplying dummy variables. In other cases, we can use any notation that is the most convenient.

$$E(wage | female, educ) = E(\beta_0 + \delta_0 female + \beta_1 educ + u | female, educ)$$

$$= \beta_0 + \delta_0 female + \beta_1 educ + E(u | female, educ)$$

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

Thus

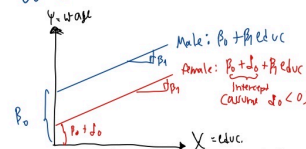
$$E(wage | female=1, educ) = \beta_0 + \delta_0 (1) + \beta_1 educ = \beta_0 + \delta_0 + \beta_1 educ$$

$$E(wage | female=0, educ) = \beta_0 + \delta_0 (0) + \beta_1 educ = \beta_0 + \beta_1 educ$$

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

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

\* given the same value of educ (same education level),  $\delta_0$  is the difference in the expected wage of females & males.



By the way, we model the regression for "female" is going to give a constant impact on wage, regardless of the level of educ.

19)

8. Multiple Regression Analysis with Qualitative Information: 83

4 It is not possible to include all of the dummy alternatives in the same model (as long as there is an interest in the model)

If we include all alternatives of a dummy variable in the same model, we will face the "perfect collinearity" problem.

$$wage = \beta_0 + \delta_0 female + \beta_1 educ + \delta_1 male + u$$

For example:

$$1 = female + male$$

$$X_0 = X_1 + X_2$$

$$female = male + 1$$

	female	male	X <sub>0</sub>
1	1	0	1
2	1	0	1
3	0	1	1
4	0	1	1
5	1	0	1
6	1	0	1

OR If there are "n" categories, we omit "n-1" category to avoid multi collinearity.

$$winter = 1 - spring - summer - fall$$

$$winter = \begin{cases} 1 & \text{if winter} \\ 0 & \text{otherwise} \end{cases}$$

$$spring = \begin{cases} 1 & \text{if spring} \\ 0 & \text{otherwise} \end{cases}$$

	winter	spring	summer	fall	X <sub>0</sub>
1	1	0	0	0	1
2	1	0	0	0	1
3	0	1	0	0	1
4	0	1	0	0	1
5	0	0	1	0	1
6	0	0	1	0	1
7	0	0	0	1	1
8	0	0	0	1	1

At least one alternative has to be dropped. We treat the dropped alternative as the "BASE GROUP" or "BASELINE" or "BENCHMARK GROUP".

```
. regress lwage female male married educ exper
note: male omitted because of collinearity
```

Source	SS	df	MS	Number of obs =
Model	54.3265253	4	13.5816313	526
Residual	94.0032262	521	.180428457	F(4, 521) = 75.27
Total	148.329751	525	.28253286	Prob > F = 0.0000

lwage	Coeff.	Std. Err.	t	P> t	[95% Conf. Interval]
female	-.3251146	.0377061	-8.62	0.000	-.3991892 - .25104
male	0 (omitted)				
married	.1380145	.0411197	3.36	0.001	.0572338 .2187953
educ	.0872644	.0071554	12.20	0.000	.0732075 .1013213
exper	.0076213	.0015314	4.98	0.000	.0046129 .0106297
_cons	.4693918	.1040575	4.51	0.000	.264668 .6735156

Being female workers are expected to have less wage compared to male workers.

20)

8. Multiple Regression Analysis with Qualitative Information:

5 Using dummy variables for multiple categories

Case 1 We can use many dummy variables in the same model

Consider a model which includes 2 dummy variables - female and married.

$$\log(wage) = \beta_0 + \delta_0 female + \beta_1 married + \beta_2 educ + \beta_3 exper + \beta_4 exper^2 + \beta_5 tenure + \beta_6 tenure^2 + u$$

```
regress lwage female married educ exper experq tenure tenuraq
```

Source	SS	df	MS	Number of obs =
Model	65.6482326	7	9.37831895	526
Residual	82.6815188	518	.159616832	F(7, 518) = 58.76
Total	148.329751	525	.28253286	Prob > F = 0.0000

lwage	Coeff.	Std. Err.	t	P> t	[95% Conf. Interval]
female	-.2901838	.0361121	-8.04	0.000	-.3611279 - .2192396
married	-.0529219	.0407561	1.30	0.195	-.071456 .1329994
educ	.0791547	.0068003	11.64	0.000	.0673952 .0905143
exper	.0269535	.0053258	5.06	0.000	.0164907 .0374163
experq	-.0003399	.0001122	-4.81	0.000	-.0007603 - .0003196
tenure	.0313962	.0068492	4.57	0.000	.0178426 .0447499
tenuraq	-.0009744	.0002347	-4.25	0.015	-.0013055 - .0001134
_cons	.4177837	.0988662	4.23	0.000	.2235557 .6120116

1) measures the impact of being married (marriage premium) BUT since |t| < 1.96 OR P > 0.05, we don't reject H0 of no impact

2) measures the expected difference between female & male workers given the same marital status & other factors.

$$\frac{\partial \log(wage)}{\partial female} = \frac{\Delta wage}{wage} = -0.29$$

female workers are paid 29% less than male workers by 29.02%, holding other factors the same.

