

1)

4 Testing Hypotheses about a Single Linear Combination of the Parameter β_1, β_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.

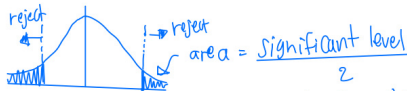
We want to test whether $\beta_1 = \beta_2$.

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

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
 \Rightarrow We use a variable transformation
 trick \Rightarrow see notes!!!

2)

Inference \rightarrow Hypothesis testing about β 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)

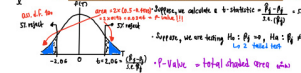
BUT, We don't know what the true β are, so we use $\hat{\beta}$ (estimator) and S.E. ($\hat{\beta}$) to test the hypothesis



Test if $\beta =$ some number
 eg. $\beta_1 = 0 \rightarrow X_1$ has no impact on Y.
 $\beta_1 = 1 \rightarrow 1$ month in X₁ correspond to 1 unit in Y.

\Rightarrow t-test of H_0 ?
 $\frac{\hat{\beta}_1 - \beta_1}{S.E.(\hat{\beta}_1)} \sim t_{df}$

* (Significance level = total area in the rejection region) \Rightarrow α



* t-value < significance level \Rightarrow always reject H_0 !!!

3)

another possible hypothesis test (one-tailed alternative)

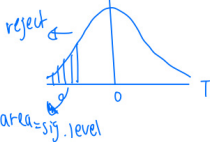
$H_0 : \beta_1 = \beta_2 \Rightarrow H_0 : \beta_1 - \beta_2 = 0$

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

It is assumed that β_1 would not be more than β_2
 (returns to a 2-year college would never be more than returns to University education)

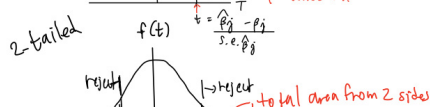
$t = \frac{(\hat{\beta}_1 - \hat{\beta}_2) - 0}{S.E.(\hat{\beta}_1 - \hat{\beta}_2)}$

* 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 H₀ 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

1st) write the t-statistic for testing H_0

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

2nd) 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 θ_1 in it. So, SPSS, STATA or 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

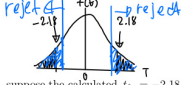
Example 1: $H_0: \beta_j \geq 0, H_a: \beta_j < 0, d.f. = 140$. \rightarrow z-table
 \rightarrow p-value = what should be the significant level, given the critical value of -2.75? \rightarrow find the shaded area

$0.5 - 0.997$
 $\Rightarrow 0.003$

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: $H_0: \beta_j = a_j, H_a: \beta_j \neq a_j, d.f. = 18$. \leftarrow 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

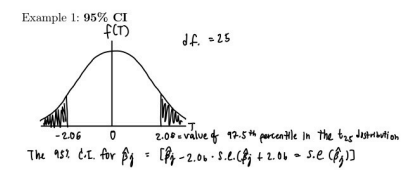
6 Confidence Intervals (CI)

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



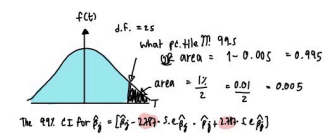
A 95% CI of β_j is given by $\hat{\beta}_j \pm 2 \cdot s.e.(\hat{\beta}_j)$
 CI $\Rightarrow \hat{\beta}_j \pm 2 \cdot s.e.(\hat{\beta}_j)$
 CI is the 95% percentile in the t-distribution with $n-k-1$ d.f.

6)



Example 1: 95% CI
 $d.f. = 25$
 The 95% CI for $\beta_j = [\hat{\beta}_j - 2.06 \cdot s.e.(\hat{\beta}_j) + 2.06 \cdot s.e.(\hat{\beta}_j)]$
 The 95% value of 75.5th percentile in the t_{25} distribution

Example 2: 99% CI



$d.f. = 25$
 what percentile? 99.5
 OF AREA = $1 - 0.005 = 0.995$
 area = $\frac{1-\alpha}{2} = \frac{0.01}{2} = 0.005$
 The 99% CI for $\beta_j = [\hat{\beta}_j - 2.81 \cdot s.e.(\hat{\beta}_j), \hat{\beta}_j + 2.81 \cdot s.e.(\hat{\beta}_j)]$

19.1 Nov

F-test motivation

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

Grade 325 = $\beta_0 + \beta_1 \cdot \# \text{ times front} + \beta_2 \cdot \# \text{ times back} + \beta_3 \cdot \# \text{ study} + \beta_4 \cdot \text{past GPA} + \beta_5 \cdot \text{gender} + u$
 $H_0: \text{seat position doesn't have impact on GPA}$
 $\beta_1 = 0 \ \& \ \beta_2 = 0 \Rightarrow \beta_1 = \beta_2 = 0$
 $H_a: \text{seat position matters}$
 $\beta_1 \neq 0 \ \& \ \beta_2 \neq 0$
 OR $\beta_1 \neq 0 \ \& \ \beta_2 = 0$
 OR $\beta_1 = 0 \ \& \ \beta_2 \neq 0$
 at least one of the $\beta_1, \beta_2 \neq 0$

7)

7 Testing Multiple Linear Restrictions: The F-test

Suppose the model is specified by

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

$H_0: \beta_2 = 0 \ \& \ \beta_3 = 0 \rightarrow$ 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). \leftarrow small model
 $Y = \beta_0 + \beta_1 x_1 + u$ is true \Rightarrow do not reject H_0

* Suppose there are "q" no. of β 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.

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{k-q} X_{k-q} + \beta_{k-q+1} X_{k-q+1} + \beta_{k-q+2} X_{k-q+2} + \dots + \beta_k X_k + u$$

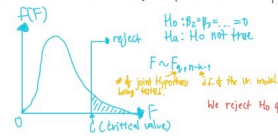
$$F = \frac{\text{ur} - \text{r}}{\text{SSR}_{ur} - \text{SSR}_{r}} \cdot \frac{q}{\text{SSR}_{ur} / (n-k-1)}$$

This is always > 0 bco $\text{SSR}_{ur} < \text{SSR}_{r}$. 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, $\text{Var}(\hat{\beta}_j)$ will increase, making the prediction of β 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 $\& \uparrow R^2$



$H_0: \beta_2 = \beta_3 = \dots = 0$
 $H_a: H_0$ not true
 $F \sim F_{q, n-k-1}$
 We reject H_0 if jointly no effect if $F > c$

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 $F = \frac{(R^2_{ur} - R^2_r) / (k - r)}{(1 - R^2_{ur}) / (n - k - 1)}$

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

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.

$H_0: \beta_{\text{gamesyr}} = \beta_{\text{bavg}} = \beta_{\text{hrunsyr}} = \beta_{\text{rbsyr}} = 0$

$H_a: \text{otherwise to } H_0$

the unrestricted model (ur) is defined by

UR Model

Source	SS	df	MS	Number of obs = 353
Model	308.989208	5	61.7978416	F(5, 347) = 117.06
Residual	183.186329	347	.527914487	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	-.0203462 .0088187
_cons	11.12942	.288229	38.75	0.000	10.62433 11.76048

the restricted model (r) is defined by

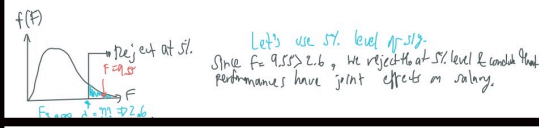
R Model

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 = 0.5948
				Root MSE = .75273

log_salary	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
years	.071318	.012505	5.70	0.000	.0467236 .0959124
gamesyr	-.0201745	.0013429	15.02	0.000	-.0175334 -.0228156
_cons	11.2238	.109312	103.62	0.000	11.01078 11.43689

Now, our H_0 and H_a becomes

$F = \frac{(SSR_r - SSR_{ur}) / (k - r)}{SSR_{ur} / (n - k - 1)}$
 $= \frac{(198.311477 - 183.186329) / (5 - 2)}{(183.186329) / (353 - 5 - 2)} \approx 9.55$



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 β_1 is statistically significant at 5% level. (value of $\beta_1 \neq 0$)

\Rightarrow When $t < 1.96$ we can say that β_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 β .

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

$bweight = \beta_0 + \beta_1 \text{cigs} + \beta_2 \text{faminc}$

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??

$bweight_{kg} = \frac{\beta_0}{1000} + \beta_1 \text{cigs} + \frac{\beta_2}{1000} \text{faminc}$
 $\Rightarrow \alpha_0 = \frac{\beta_0}{1000}, \alpha_1 = \beta_1, \alpha_2 = \frac{\beta_2}{1000}$

What if we use $faminc$ in USD (instead of 1000 USD)

$bweight = \beta_0 + \beta_1 \text{cigs} + \beta_2 \text{faminc}$
 $\Rightarrow \beta_1 = \frac{\beta_2}{1000}$

in other words β_2 = Impact of 1 USD in income.
 $\beta_2 = \frac{1}{1000}$

What if we use $bweight$ in kg & income in THB
 $bweight_{kg} = \frac{\beta_0}{1000} + \frac{\beta_1}{1000} \text{cigs} + \left(\frac{\beta_2}{1000}\right) \text{faminc}_{THB}$
 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

$$100 \beta_2 = \frac{100 \cdot \frac{1}{Y} dY}{d X_2} = \frac{\% \Delta Y}{\Delta X_2}$$

∴ 100 β₂ = % Δ in Y given that X₂ increases by 1 Unit.

- Models with Quadratics (Squares)

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

Cost-Y 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
∴ $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + u$
∴ $\frac{dY}{dX} = \beta_1 + 2\beta_2 X$
∴ $\frac{d^2Y}{dX^2} = 2\beta_2$
∴ $\beta_2 < 0$ ⇒ decreasing returns

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)	price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
lnox	price	-.9787545	-.0959398	-9.81	0.000	-1.172429 - .7810896
dist	price	-.0214972	-.0094013	-3.42	0.001	-.030668 - .0137264
rooms	price	-.5528052	.1612965	-3.43	0.001	-.8697056 - .2359007
room ²	price	.0624657	-.0124867	-5.00	0.000	-.0179386 - .0070025
stratio	price	-.0486667	.0058131	-8.37	0.000	-.0600679 - .0372455
_cons	price	13.39154	.5692901	24.95	0.000	12.48113 14.70198

Log(price)

In the US or many other countries, students don't really go to school in the area with bad air to make any left. So, the lower stratio, the better the school

all variables are significant

Consider the effect of "room"

$$\frac{d \log(\text{price})}{d \text{rooms}} = \beta_3 + 2\beta_4 \text{rooms} = -0.553 + 2(0.062) \cdot \text{rooms}$$

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

$$0 = -0.553 + 2(0.062) \cdot \text{rooms}$$

rooms = 4.4

Answer ⇒ at 4.4 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?

$$\frac{d \log(\text{price})}{d \text{rooms}} = -0.553 + 2(0.062) \cdot \text{rooms}$$

at 5 rooms: $-0.553 + 2(0.062) \cdot 5 = 0.257$

at 6 rooms: $-0.553 + 2(0.062) \cdot 6 = 0.374$

∴ % Δ price = 100(-0.553 + 2(0.062) · 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 β₂ > 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 → R² 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 R² would not punish "having too many regressors".
- We use adjusted-R² or R² when we want to punish adding too many regressors.

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

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

If we have more k, d.f. = n - k - 1 ↓, SSR / (n - k - 1) ↑, adj. -R² ↓

(∴ adj. -R² ↓ ⇒ 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, R² = 0.029, R² = 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, R² = 0.282, R² = 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 δ_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

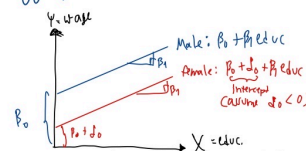
$$E(wage | female=1, educ) = \beta_0 + \delta_0 (1) + \beta_1 educ = \beta_0 + \delta_0 + \beta_1 educ$$

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

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

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

* given the same value of educ (same education level), δ_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	X ₀
1	1	0	1
2	0	1	1
3	1	0	1
4	0	1	1
5	1	0	1
6	0	1	1
7	1	0	1
8	0	1	1

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

$$1 = winter + spring + summer + fall$$

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

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

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 experq tenure tenuraq

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	Coef.	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.

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

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	Coef.	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	.0679792 .0903303
exper	.0269535	.0053258	5.06	0.000	.0164907 .0374163
experq	-.0003399	.0001122	-4.81	0.000	-.0007603 - .0001196
tenure	.0313962	.0068482	4.57	0.000	.0178426 .0447499
tenuraq	-.0009744	.0002347	-4.20	0.000	-.0013355 - .0001134
_cons	.4177837	.0988662	4.23	0.000	.2235557 .6120116

1) measures the impact of being married (marriage premium)

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

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

$$\text{female workers are paid } = 100 \times \frac{\Delta wage}{wage} = 100 \times -0.29$$

$$\text{to earn less than male workers by } 29.02\% \text{, holding other factors the same.}$$

Comment: since |t| < 1.96 OR P > 0.05, we don't reject H0 if no impact

