

Athicha Korkietsatean

6104640062



4 Testing Hypotheses about a Single Linear Combination of the Parameter

Consider

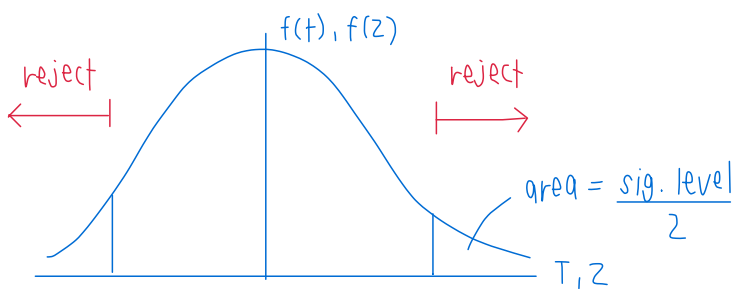
maybe non-linear
so, take log to be more linear

$$\log(\text{wage}) = \beta_0 + \beta_1 jc + \beta_2 univ + \beta_3 \text{exp er} + u$$

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

We want to test whether $\beta_1 = \beta_2$. If the return from 1 more year of education at a JC is the same as that of the university.

$$\left. \begin{array}{l} H_0 : \beta_1 = \beta_2 \rightarrow H_0 : \beta_1 - \beta_2 = 0 \\ H_a : \beta_1 \neq \beta_2 \rightarrow H_a : \beta_1 - \beta_2 \neq 0 \end{array} \right\} \text{Two-tailed test}$$



$$t = \frac{(\hat{\beta}_1 - \hat{\beta}_2) - 0}{\text{s.e.}(\hat{\beta}_1 - \hat{\beta}_2)}$$

Then compute t and compare w/ the critical value.

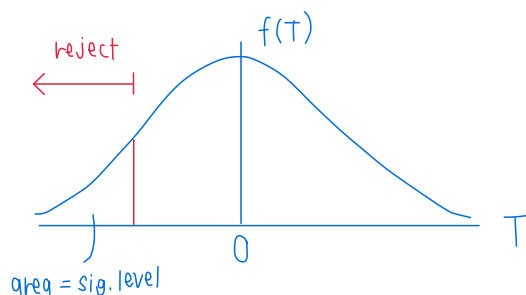
where $\text{s.e.}(\hat{\beta}_1 - \hat{\beta}_2) = \sqrt{\text{var}(\hat{\beta}_1 - \hat{\beta}_2)}$
 (use variable transformation)
 $= \sqrt{\text{var}(\hat{\beta}_1) + \text{var}(\hat{\beta}_2) + 2\text{cov}(\hat{\beta}_1, \hat{\beta}_2)}$

another possible hypothesis test (one-tailed alternative)

$$H_0: \beta_1 = \beta_2, \quad \beta_1 - \beta_2 = 0$$

$$H_a: \beta_1 < \beta_2, \quad \beta_1 - \beta_2 < 0$$

It's assume that β_1 would not be more than β_2
(return to a 2-yr will never be more than university)

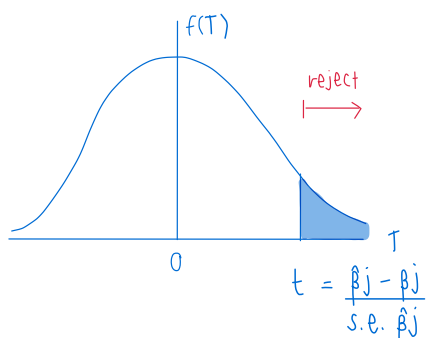


$$t = \frac{(\hat{\beta}_1 - \hat{\beta}_2) - 0}{\text{s.e.}(\hat{\beta}_1 - \hat{\beta}_2)}$$

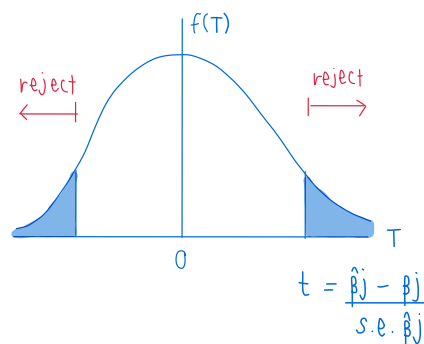
go to extra Note.

5 Computing p-Values for t-Tests

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



The shaded area in the rejection region is the p-value.



Total area from 2-sides is the p-value

- p-value : $P(|T| > |t|)$

T = t-distributed random variable w/ d.f. = $n - k - 1$

t = computed t-statistic.

p-value = prob. that a random T value will be greater
(in the $| |$ term) than our t in the H_0 test

In class exercise.

Consider the multiple regression, assume that 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 the $H_0: \beta_1 - 3\beta_2 = 1$

H_a : otherwise

1st write t-statistic for testing H_0 : $t = \frac{(\hat{\beta}_1 - 3\hat{\beta}_2) - 1}{\text{s.e.}(\hat{\beta}_1 - 3\hat{\beta}_2)}$

2nd Define $\hat{\theta}_1 = \hat{\beta}_1 - 3\hat{\beta}_2 \longrightarrow H_0: \theta_1 = 1$ $t = \frac{\hat{\theta}_1 - 1}{\text{s.e.}(\hat{\theta}_1)} \longrightarrow$ we need our regression to have θ_1 in it.
 $H_a: \theta_1 \neq 1$ so, STATA or OLS estimation will automatically give $\hat{\theta}_1, \text{s.e.}(\hat{\theta}_1)$

Now, if rearrange $\theta_1 = \beta_1 - 3\beta_2$, we have, $\hat{\beta}_1 = \hat{\theta}_1 + 3\hat{\beta}_2$
 $\beta_1 = \theta_1 + 3\beta_2$

Then, sub. in the main regression and get

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

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

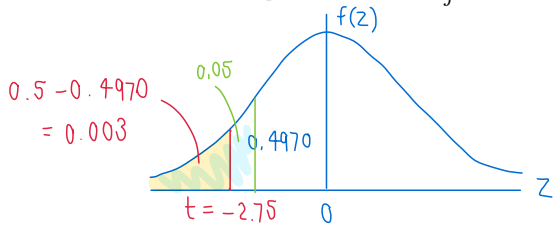
$$Y = \beta_0 + \theta_1 X_1 + \beta_2 (X_2 + 3X_1) + \beta_3 X_3 + u$$

so, we can calculate

$$t = \frac{\hat{\theta}_1 - 1}{\text{s.e.}(\hat{\theta}_1)}$$

Now the explanatory variables are $X_1, X_2 + 3X_1$ and X_3

Example 1: $H_0 : \beta_j \geq 0$, $H_a : \beta_j < 0$, d.f. = 140. \rightarrow z-table



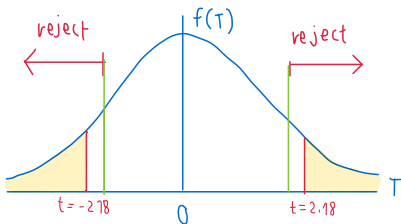
p-value = what should be the sig. level given the critical value of -2.75

find the shaded area = 0.003

suppose the calculated $t_{\hat{\beta}_j} = -2.75$

- 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**
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. (t-table)



suppose the calculated $t_{\hat{\beta}_j} = -2.18$

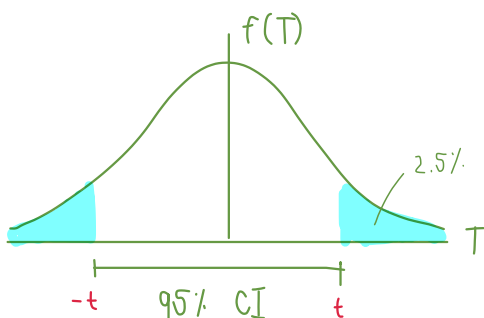
- From the t-table, the value -2.18 corresponds to area = 0.02 - 0.05
- Thus, p-value = btw 0.02 and 0.05
- Would we reject H_0 if we use the significance level = 5%?
yes, reject H_0 because the area is less than 0.05 or p-value < 0.05

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 95% chance.

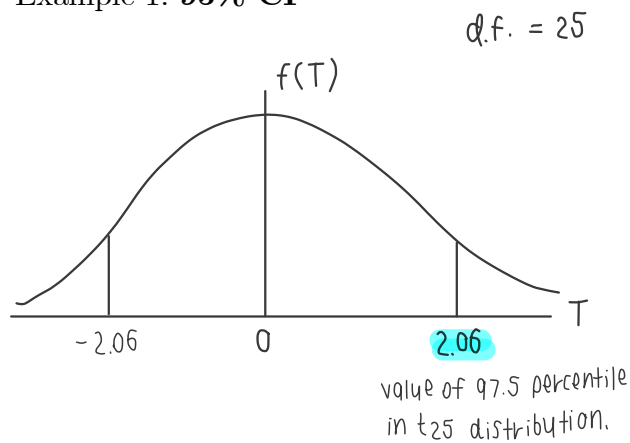
- A 95% CI of β_j is given by



$$CI \rightarrow \hat{\beta}_j \pm c \cdot s.e.(\hat{\beta}_j)$$

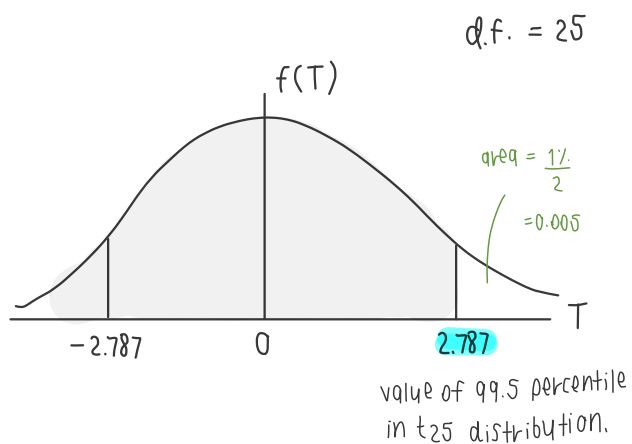
c is the 97.5 percentile in the t-distribution w/ n-k-1 d.f.

Example 1: 95% CI

The 95% CI for $\hat{\beta}_j$

$$= [\hat{\beta}_j - 2.06 \cdot \text{s.e.}(\hat{\beta}_j), \hat{\beta}_j + 2.06 \cdot \text{s.e.}(\hat{\beta}_j)]$$

Example 2: 99% CI

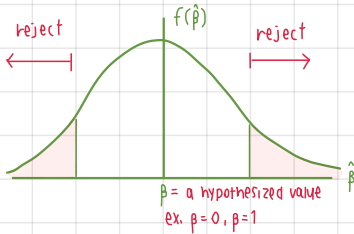
The 99% CI for $\hat{\beta}_j$

$$= [\hat{\beta}_j - 2.787 \cdot \text{s.e.}(\hat{\beta}_j), \hat{\beta}_j + 2.787 \cdot \text{s.e.}(\hat{\beta}_j)]$$

Inference \rightarrow Hypothesis testing about " β ", the true parameter

$$\text{wage} = \beta_0 + \beta_1 \text{edu} + \beta_2 \text{exp} + \dots + u$$

we want to test the hypothesis about true impact (β) of each X variables (edu, exp) on the dependent 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.

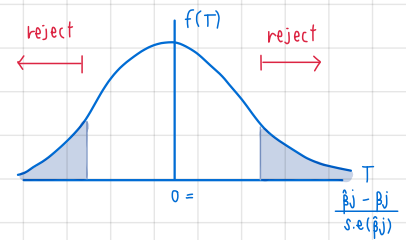


1 Test if $\beta =$ same number

e.g. $\beta_j = 0$, X_j has no impact on Y .

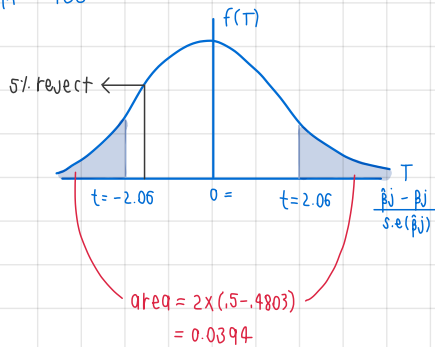
$\beta_j = 1$, 1 unit increase in X_j corresponded to 1 unit increase in Y .

$$t\text{-test} : \frac{\hat{\beta}_j - \beta_j}{\text{s.e.}(\hat{\beta}_j)} \sim t \text{ d.f.}$$



Significant level = total area in the rejection region

d. f. = 100



• suppose, we calculate a t-statistic = $\frac{\hat{\beta}_j - \beta_j}{\text{s.e.}(\hat{\beta}_j)} = 5.78$

• suppose, we are testing $H_0: \beta_j = 0$, $H_a: \beta_j \neq 0$

• p-value = total shaded area
 = significant level or prob which we'll reject H_0 .

If p-value < significant level \rightarrow Reject H_0 .

Ftest

test the sig. of a group of hypotheses (multiple hyp.)

$$\text{Grade}_{325} = \beta_0 + \beta_1 \# \text{time_front} + \beta_2 \# \text{time_back} + \beta_3 \text{hr_study} + \beta_4 \text{past_GPA} + \beta_5 \text{gender} + \epsilon$$

H_0 : seat position doesn't have impact on GPA

$$\beta_1 = 0, \beta_2 = 0; \beta_1 = \beta_2 = 0$$

H_a : seat position matters

$$\left. \begin{array}{l} \beta_1 \neq 0, \beta_2 \neq 0 \\ \text{or } \beta_1 \neq 0, \beta_2 = 0 \\ \text{or } \beta_1 = 0, \beta_2 \neq 0 \end{array} \right\} \text{at least 1 of } \beta_1, \beta_2 \neq 0$$

7 Testing Multiple Linear Restrictions: The F-test

Suppose the model is specified by

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

$$H_0 : \beta_2 = 0 \text{ and } \beta_3 = 0 \text{ if both } X_2 \text{ \& } X_3 \text{ don't have impact on } Y$$

$$H_1 : H_0 \text{ is not true \# otherwise}$$

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

1. Our full model is called the "unrestricted" model (ur). ^{big model: contain all variables} Suppose it can be expressed as:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + u \text{ is true (reject } H_0)$$

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

2. The model which takes out x (which we think its associated $\beta = 0$) is called the restricted model (r).

^{small model}

$$y = \beta_0 + \beta_1 x_1 + u \text{ is true (not reject } H_0)$$

suppose there're q no. of β that we would like to perform a joint test of $= 0$

e.g. $q = 2$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{k-q} x_{k-q}$$

$$H_0 : \beta_{k-q+1} = \beta_{k-q+2} = \dots = \beta_k = 0, \text{ last } q \text{ } \beta\text{'s} = 0$$

H_a : H_0 is not true.

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

$\underbrace{\hspace{150px}}_r$
 $\underbrace{\hspace{150px}}_{k+4}$

$\underbrace{\hspace{300px}}_{ur}$

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

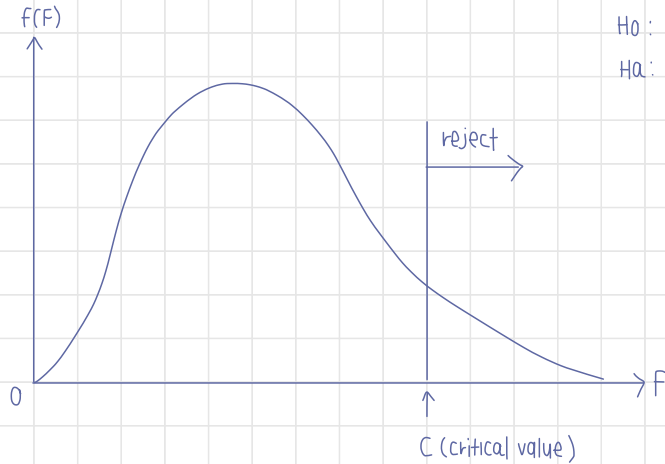
d.f. of ur model

Always positive as $SSR_{ur} < SSR_r$ b/c more explanatory var., less residual. } $SSR \downarrow, r^2$ improve
 so, everytime you add 1 more X , the model will be better explained.
 so, why don't keep adding X variables AMAP
 B/c, everytime we add 1 more X , $var(\hat{\beta}_s)$ will increase making the prediction of β less precise. so, we only keep the addition of X s if it/they can improve the model enough
 can significantly $\downarrow SSR$ and $\uparrow r^2$

SSR_r = some of sq. residual in restricted model
 SSR_{ur} = some of sq. residual in unrestricted model
 q = no. of explanatory variables we want to test = 0

F-test always have 1 tail

if F statistic is close to zero then we fail to reject H_0
if F statistic is far from zero then we reject H_0



$H_0: \beta_2 = \beta_3 = \dots = 0$
 $H_A: H_0 \text{ not true}$

$F \sim F_{q, n-k-1}$ d.f. of ur Model
/
no. of joint hyp being tested

we reject H_0 of jointly no effect if $F > C$

3. Some useful facts

1. $R^2_{ur} > R^2_r$ # better b/c additional X will increase R^2 (improve fit) $\rightarrow SSR_{ur} < SSR_r$
2. Add more X the model is certainly explained. However we would like to reject H_0 if the inclusion of extra variables doesn't improve the model enough.

4. Other ways to calculate the F-statistics: in terms of R^2

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

We have $F = \frac{(R^2_r - R^2_{ur})}{q - \text{no. of } \beta \text{ that set to zero}} \cdot \frac{1 - R^2_{ur}}{n - k - 1}$

$\underbrace{\hspace{10em}}_{\text{obs. coefficient}}$
 $\underbrace{\hspace{1em}}_{\text{slope}}$
 $\underbrace{\hspace{1em}}_{\text{intercept}}$

W. O. test if $\beta = 0$
 \rightarrow don't know that X putting in the regression is making sense.
 \rightarrow test whether any β are non zero.

If we want to test the overall sig. of the model

$H_0: \beta_1 = \beta_2 = \beta_3 = \dots = \beta_K = 0$

H_A : otherwise



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

R^2 \leftarrow R of ur that the "r" model has no X at all.

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

- | | | |
|---|----|---|
| { | ur | y <i>salary</i> = season salary |
| | | y <i>years</i> = years in major leagues |
| | | y <i>gamesyr</i> = games per year in the league |
| | | y <i>bavg</i> = career batting average |
| | | y <i>hrunsyr</i> = homeruns per year |
| | | y <i>rbisyr</i> = runs batted in per year |

If we want to test whether performance has any impact on salary.

$H_0: \beta_{bavg} = \beta_{hrunsyr} = \beta_{rbisyr} = 0$

H_A : otherwise is true

- the unrestricted model (ur) is defined by

y

```
. regress log_salary years gamesyr bavg hrunsyr rbisyr # 4r model
```

Source		SS	df	MS	
SSE	Model	308.989208	5	61.7978416	Number of obs = 353
SSR	Residual	183.186327	347	.527914487	F(5, 347) = 117.06
					Prob > F = 0.0000
					R-squared = 0.6278
SST	Total	492.175535	352	1.39822595	Adj R-squared = 0.6224
					Root MSE = .72658

more explanatory variables = more r²

use in Multiple Regression

lower S.E.

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	.0073464 .0177578
bavg	.0009786	.0011035	0.89	0.376	-.0011918 .003149
hrunsyr	.0144295	.016057	0.90	0.369	-.0171518 .0460107
rbisyr	.0107657	.007175	1.50	0.134	-.0033462 .0248776
_cons	11.19242	.2888229	38.75	0.000	10.62435 11.76048

• the restricted model (r) is defined by

when considering each of the performance X one-by-one. None of them has a significant impact at 5%. However, there is a small impact. But when performing a F-test, performances have joint impact.

```
# r model . regress log_salary years gamesyr
```

Source		SS	df	MS	
SSE	Model	293.864058	2	146.932029	Number of obs = 353
SSR	Residual	198.311477	350	.566604221	F(2, 350) = 259.32
					Prob > F = 0.0000
					R-squared = 0.5971
SST	Total	492.175535	352	1.39822595	Adj R-squared = 0.5948
					Root MSE = .75273

higher S.E.

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	.108312	103.62	0.000	11.01078 11.43683

Now, our H_0 and H_a becomes

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

no. of joint hyp. we are testing

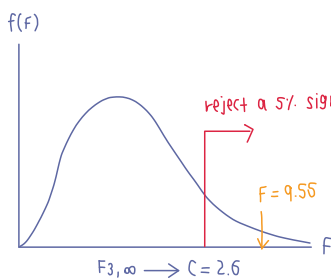
$$= \frac{(198.311 - 183.186) / 3}{183.186 / (353 - 5 - 1)}$$

$$= 9.55$$

HW.

$$F = \frac{R^2 / q}{(1 - R^2) / (n - k - 1)}$$

$$= ?$$



Since our F is 9.55 > 2.6, we reject H_0 at 5% significant level. And conclude that performances have joint effects on salary.

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)

⇒ These t – *statistics* are to test $H_0 : \beta_i = 0$

⇒ If the d.f. > 30 , then when $t > 1.96$, we can reject H_0 w/ 5% sig. ↖ z-table

⇒ **When $t > 1.96$** , we can say that β_i is **statistically significant** at 5% level. (value of $\beta_i \neq 0$)

⇒ **When $t < 1.96$** we can say that β_i is **not statistically significant** at 5% level.

⇒ If $t < 1.96$ we can drop x_i from the model

⇒ After we drop x_i , we estimate the new regression function and obtain a new set of $\hat{\beta}$.

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

e.g. $H_0 : \beta_i = \beta_j$

or $H_0 : \beta_i = 5$ 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: $\log(\text{salary})$			
Independent Variables	(1) <i>very restricted</i>	(2) <i>more restricted</i>	(3) <i>unrestricted</i>
<i>sales</i> log(<i>sales</i>)	.224 (.027)	.158 (.040)	.188 (.040)
<i>company performance</i> log(<i>mktval</i>) <i>profmarg</i>	—	.112 (.050)	.100 (.049)
	—	-.0023 (.0022)	-.0022 (.0021)
<i>CEO characteristics</i> <i>ceoten</i> <i>comten</i>	—	—	.0171 (.0055)
	—	—	-.0092 (.0033)
<i>intercept</i>	4.94 (0.20)	4.62 (0.25)	4.57 (0.25)
Observations	177	177	177
R-squared	.281	.304	.353

simple regression
w/ 1x variable

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

$$\widehat{bweght} = \widehat{\beta}_0 + \widehat{\beta}_1 cigs + \widehat{\beta}_2 faminc,$$

where

$bweght$ = 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 $bweght$ in kilograms ; $1kg = 1000g$

$$\begin{aligned} \widehat{bweght}_{kg} &= \frac{\widehat{bweght}_g}{1000} = \frac{\widehat{\beta}_0}{1000} + \frac{\widehat{\beta}_1}{1000} cigs + \frac{\widehat{\beta}_2}{1000} faminc \\ &= \widehat{\alpha}_0 + \widehat{\alpha}_1 cigs + \widehat{\alpha}_2 faminc \end{aligned}$$

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

$$bweght_g = \widehat{\beta}_0 + \widehat{\beta}_1 cigs + \frac{\widehat{\beta}_2}{1000} faminc_{USD} \leftarrow \text{will be 1000 times larger than actual faminc.}$$

$$bweght_g = \widehat{\beta}_0 + \widehat{\beta}_1 cigs + \widehat{\theta}_2 faminc_{USD}$$

in other words, $\widehat{\theta}_2 = \text{impact of 1 USD } \uparrow \text{ in income.}$

$\widehat{\beta}_2 = \text{impact of 1000 USD } \uparrow \text{ in income.}$

- What if we use $bweght$ in kg and income in THB

$$bweght_{kg} = \widehat{\beta}_0 + \frac{\widehat{\beta}_1}{1000} cigs + \frac{\widehat{\beta}_2}{30,000} faminc_{THB} \leftarrow \text{will be 30,000 times more than actual faminc.}$$

2 More on functional forms

• Logarithmic Functional Form

$$\# \frac{d \log(x)}{d(x)} = \frac{1}{x} \rightarrow d \ln(x) = \frac{1}{x} dx$$

$$\Delta Y = Y_1 - Y_2$$

$$\Delta X_1 = X_{11} - X_{12}$$

natural log

$$\log(y) = \beta_0 + \beta_1 \log(x_1) + \beta_2 x_2 + u$$

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

with the $\log(Y)$ and $\log(X_1)$, the coefficient is going to be the elasticity. (X1 elasticity of Y) (price) (demand)

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

we want this to be %Δ, then:

$$100 \beta_2 = \frac{100 \frac{1}{Y} \Delta Y}{\Delta X_2}$$

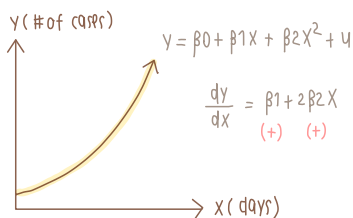
∴ so, $100\beta_2 = \% \Delta Y$ given that X_2 increases by 1 unit

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

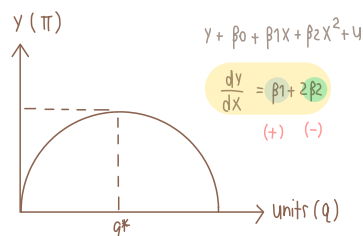
• Models with Quadratics (squares)

capture inc/dec marginal effects. (slope of the relationship btw X,Y is not constant.)

COVID-19



Decreasing returns.



Assume that $MC = 10$, $P = 100 - Q$

$$\pi = (P - MC) Q$$

$$\pi = (100 - Q - 10) Q$$

$$FOC: \frac{d\pi}{dQ} = 0 = 90 - 2Q$$

Example : Effects of Pollution on Housing Prices

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

where

price = housing price

nox = level of pollution

dist = distance from downtown

rooms = number of rooms

stratio = average student per teacher ratio (in us or many other country, student can apply to school in the area w/o any competition. The lower stratio, the better the school.)

The estimation result is given by

regress lprice lnox dist rooms rooms_sq stratio

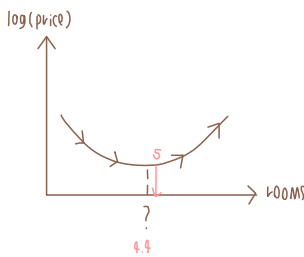
Source	SS	df	MS			
Model	51.4933152	5	10.298663	Number of obs =	506	
Residual	33.0889098	500	.06617782	F(5, 500) =	155.62	
				Prob > F =	0.0000	
				R-squared =	0.6088	
				Adj R-squared =	0.6049	
				Root MSE =	.25725	
lprice	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
lnox	-.9767545	.0995938	-9.81	0.000	-1.172429	-.7810806
dist	-.0321972	.0094013	-3.42	0.001	-.050668	-.0137264
rooms	-.5528032	.1612965	-3.43	0.001	-.8697056	-.2359007
rooms_sq	.0624697	.0124867	5.00	0.000	.0379368	.0870025
stratio	-.0486667	.0058131	-8.37	0.000	-.0600879	-.0372455
_cons	13.59154	.5650901	24.05	0.000	12.4813	14.70178

$|t| > 1.96$
all are significant $p < 0.05$

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 has a positive impact on log(price)

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

$$\text{rooms} = 4.4 \leftarrow \text{round up}$$

\therefore At 5 rooms or more

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

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

$$100 \cdot \frac{1}{\text{price}} \frac{d \text{price}}{d \text{rooms}} = 100 [-0.553 + 2(0.062) \cdot 5]$$

$$= 100 \cdot 0.067$$

$$= 6.7\% \text{ increase}$$

what abt % change in price when #rooms increase from 5 to 7

$$\% \Delta \text{price} = 100 [-0.553 + 2(0.062) \cdot 6]$$

$$= 100 (0.191)$$

$$= 19.1\% \text{ increase}$$

\therefore Total % Δ price when #rooms increase from 5 to 7 is $6.7 + 19.1 = 25.8\%$

3 Models with Interaction Terms

use when the impact of one variable depends on the value (level) of another variable.

Consider

$$price = \beta_0 + \beta_1 \overset{x_1}{sqr\ ft} + \beta_2 \overset{x_2}{bdrms} + \beta_3 \overset{x_3 = x_1 x_2}{sqr\ ft \times bdrms} + \beta_4 \overset{x_4}{bthrms} + u$$

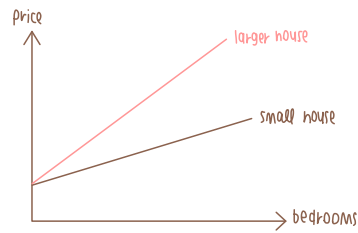
where

$price$ = housing price

$sqr\ ft$ = house size (square feet)

$bdrms$ = number of bedrooms

$bthrms$ = number of bathrooms



$$\frac{dprice}{dbdrms} = \beta_2 + \beta_3 \text{ sqrft}$$

if $\beta_3 > 0$ then, an additional bedroom would inc. price more for a larger house.

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

- Adding more regressors ALWAYS improve fit R^2 always \uparrow

• But we lose the "degree of freedom"
 (free data point used to estimate the parameter
 1 data point is sacrificed everytime we estimate a parameter.

- using R^2 would not punish having too many regressors

- use adjusted R^2 or \bar{R}^2 when we want to punish adding too many regressors.

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

Adjusted $R^2 \rightarrow \bar{R}^2 = 1 - \frac{SSR/n-k-1}{SST/(n-1)}$

If we have 1 more k , the d.f. \downarrow .
 making $SSR/(n-k-1) \uparrow$
 result in the dec in \bar{R}^2

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

Consider Model 1

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

$$n = 209, R^2 = 0.029, \bar{R}^2 = 0.020$$

Consider Model 2

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

$$n = 209, R^2 = 0.282, \bar{R}^2 = 0.275$$

non-restricted model.

27.5% of the variation in y is explained. So, this model is better.

Multiple Regression Analysis with Qualitative Information:

dummy variables

1 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

2 Describing Qualitative Information

- "Female" and "Married" are qualitative variable.
- We arbitrarily assign a dummy variable to describe them.

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

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

<i>person</i>	<i>wage</i>	<i>educ</i>	<i>exper</i>	<i>female</i>	<i>married</i>
1	3.10	11	2	1	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

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.

$$\begin{aligned} 1. E(wage | female, educ) &= E(\beta_0 + \delta_0 female + \beta_1 educ + u | female, educ) \\ &= \beta_0 + \delta_0 female + \beta_1 educ + E(u | female, educ) \quad \text{Assume that MLR 1-4 hold} \\ &= \beta_0 + \delta_0 female + \beta_1 educ \quad \begin{array}{c} || \\ 0 \end{array} \end{aligned}$$

2. Thus

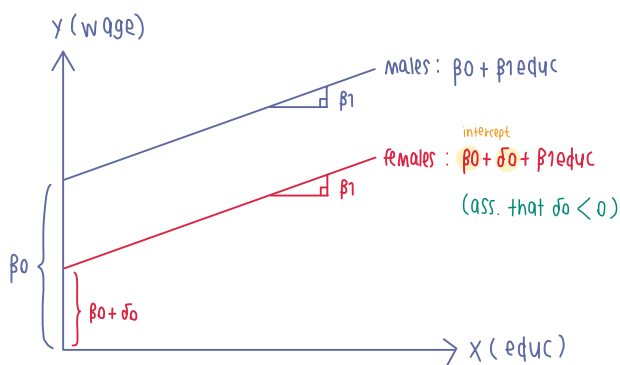
$$\text{♀} : E(wage | female = 1, educ) = \beta_0 + \delta_0(1) + \beta_1 educ = \beta_0 + \delta_0 + \beta_1 educ$$

$$\text{♂} : E(wage | female = 0, educ) = \beta_0 + \delta_0(0) + \beta_1 educ = \beta_0 + \beta_1 educ$$

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

* given the same value of educ (same education level), δ_0 is the diff.

in the expected wage of females and males.



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

4 It is not possible to include all of the dummy alternatives in the same model as long as there is an intercept in the model.

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

$$wage = \beta_0 x_0 + \beta_1^{x_1} female + \beta_2 educ + \beta_3^{x_3} male + \epsilon$$

For example:

$$x_0 = \text{intercept} = 1$$

$$x_0 = x_1 + x_3$$

$$1 = female + male$$

$$female = male + 1$$

or

$$1 = \cancel{winter} + \cancel{spring} + \cancel{summer} + \cancel{fall}$$

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

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

If there are n categories we omit 1 categories to avoid multi collinearity

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

1 if female
0 if male } 1 if male
0 if female

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

Source	SS	df	MS	Number of obs = 526		
Model	54.3265253	4	13.5816313	F(4, 521) = 75.27		
Residual	94.0032262	521	.180428457	Prob > F = 0.0000		
Total	148.329751	525	.28253286	R-squared = 0.3663		
				Adj R-squared = 0.3614		
				Root MSE = .42477		

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	.4690918	.1040575	4.51	0.000	.264668	.6735156

Female workers are expected to have less wage compared to male workers

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*.

$$\log(\text{wage}) = \beta_0 + \delta_0 \text{female} + \delta_1 \text{married} + \beta_1 \text{educ} + \beta_2 \text{exper} + \beta_3 \text{exper}^2 + \beta_4 \text{tenure} + \beta_5 \text{tenure}^2 + u$$

$\left. \begin{matrix} 1 \text{ if female} \\ 0 \text{ if otherwise} \end{matrix} \right\}$
 $\left. \begin{matrix} 1 \text{ if married} \\ 0 \text{ if otherwise} \end{matrix} \right\}$

`regress lwage female married educ exper expersq tenure tenursq`

Source	SS	df	MS			
Model	65.6482326	7	9.37831895	Number of obs = 526		
Residual	82.6815188	518	.159616832	F(7, 518) = 58.76		
Total	148.329751	525	.28253286	Prob > F = 0.0000		
				R-squared = 0.4426		
				Adj R-squared = 0.4351		
				Root MSE = .39952		

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	-.0271456	.1329894
educ	.0791547	.0068003	11.64	0.000	.0657952	.0925143
exper	.0269535	.0053258	5.06	0.000	.0164907	.0374163
expersq	-.0005399	.0001122	-4.81	0.000	-.0007603	-.0003196
tenure	.0312962	.0068482	4.57	0.000	.0178426	.0447499
tenursq	-.0005744	.0002347	-2.45	0.015	-.0010355	-.0001134
_cons	.4177837	.0988662	4.23	0.000	.2235557	.6120116

Comments:

1. δ_0 measures the expected diff. btw. female & male workers given same marital status and other factors.

$$\frac{d \log(\text{wage})}{d \text{female}} = \frac{1}{\text{wage}} \frac{d \text{wage}}{d \text{female}} = -0.29$$

∴ female workers are expected to earn less than male by 29.02%, holding other factors to be constant.

$$= 100 \frac{1}{\text{wage}} \frac{d \text{wage}}{d \text{female}} = -0.29 (100)$$

$$\frac{\% \Delta \text{wage}}{\% \Delta \text{female}} = 29.02 \%$$

2. δ_1 measure the impact of being married (marriage premium)

can't conclude that have impact

But since $|t| < 1.96$ or $p > 0.05$, we don't reject H_0 of no impact.

	♀	♂	
marr	marrfem	marrmale	
sing	singfem	singmale	→ base case

8. Multiple Regression Analysis with Qualitative Information: 85

Consider a model which includes dummy variables for each gender/marital status combination— *marrmale*, *marrfem* and *singfem*. Or *singmale* → base case

$$\log(wage) = \beta_0 + \delta_0 marrmale + \delta_1 marrfem + \delta_2 singfem + \beta_1 educ + \beta_2 exper + \beta_3 exper^2 + \beta_4 tenure + \beta_5 tenure^2 + u. \tag{8.1}$$

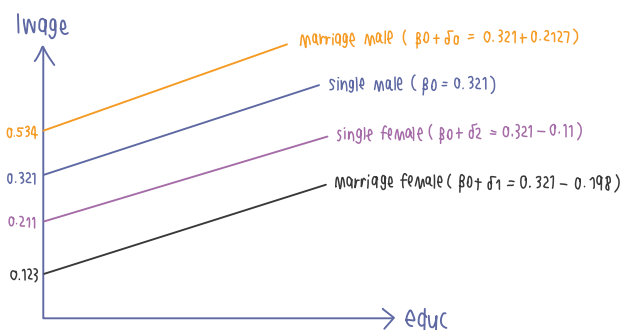
`regress lwage marrmale marrfem singfem educ exper expersq tenure tenursq`

Source	SS	df	MS	Number of obs = 526		
Model	68.3617623	8	8.54522029	F(8, 517) = 55.25		
Residual	79.9679891	517	.154676961	Prob > F = 0.0000		
Total	148.329751	525	.28253286	R-squared = 0.4609		
				Adj R-squared = 0.4525		
				Root MSE = .39329		

lwage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
marrmale	.2126757	.0553572	3.84	0.000	.103923	.3214284
marrfem	-.1982676	.0578355	-3.43	0.001	-.311889	-.0846462
singfem	-.1103502	.0557421	-1.98	0.048	-.219859	-.0008414
educ	.0789103	.0066945	11.79	0.000	.0657585	.092062
exper	.0268006	.0052428	5.11	0.000	.0165007	.0371005
expersq	-.0005352	.0001104	-4.85	0.000	-.0007522	-.0003183
tenure	.0290875	.006762	4.30	0.000	.0158031	.0423719
tenursq	-.0005331	.0002312	-2.31	0.022	-.0009874	-.0000789
_cons	.3213781	.100009	3.21	0.001	.1249041	.5178521

Comments: Not the same as the prev. one. It use "single male" as the base group
 But the prev. one use male & single as 2 base groups.

- δ_0 measure the expected diff. in wage of married male as compared with single males, holding other factors constant.
- δ_1 measures the expected diff. in wage of married female as compared with single males, holding other factors constant.
- δ_2 measures the expected diff. in wage of single female as compared with single males, holding other factors constant.



Case 2 We can use dummy variables to represent multiple categories of a variable. Consider the relationship between law school rankings and starting salaries

$$\log(\text{salary}) = \beta_0 + \delta_0 \text{top10} + \delta_1 r11_25 + \delta_3 r26_40 + \delta_4 r41_60 + \beta_1 \text{LSAT} + \beta_2 \text{GPA} + \beta_3 \log(\text{libvol}) + \beta_4 \log(\text{cost}) + u.$$

where top10 , $r11_25$, $r26_40$, $r41_60$ would be equal to 1 when the variable rank falls into the appropriate range.

** Rank below 60 would be the base case.

In many case the range of value serve as a better explanatory variable than the value itself.

```
. regress lsalary top10 r11_25 r26_40 r41_60 LSAT GPA llibvol lcost
```

Source	SS	df	MS	Number of obs =	136
Model	9.16538532	8	1.14567316	F(8, 127) =	120.15
Residual	1.2109665	127	.009535169	Prob > F =	0.0000
Total	10.3763518	135	.076861865	R-squared =	0.8833
				Adj R-squared =	0.8759
				Root MSE =	.09765

lsalary	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
top10	.5393428	.053542	10.07	0.000	.4333927 .6452928
r11_25	.4716199	.0390921	12.06	0.000	.3942637 .548976
r26_40	.2790977	.0346972	8.04	0.000	.2104383 .3477571
r41_60	.182382	.0283098	6.44	0.000	.126362 .238402
LSAT	.0060482	.0034919	1.73	0.086	-.0008616 .012958
GPA	.1305893	.0818678	1.60	0.113	-.0314122 .2925908
llibvol	.0725522	.0289213	2.51	0.013	.0153221 .1297824
lcost	.0249169	.0283224	0.88	0.381	-.031128 .0809619
_cons	8.363103	.4457314	18.76	0.000	7.481081 9.245125

Comments:

δ_0 measures the diff. in expected $\log(\text{salary})$ of a law-school graduate from a top 10 university compared to expected $\log(\text{salary})$ of those who graduated from the school ranked 61th and worse.

δ_1 → same rationale

eg. age may explain the model better if split into generations. young 0-15 genz 16-29 etc.

The baseline is ranking 61th and worse