

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 *with 5% sig. level* 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(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

sales
other company performance
CEO characteristics

simple regression with 12

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

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

where

$bwght$ = 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 $bwght$ in kilograms?

$$\widehat{bwght}_{kg} = \frac{1kg = 1000g}{1,000} \widehat{bwght}_g = \frac{\widehat{\beta}_0}{1,000} + \frac{\widehat{\beta}_1 cigs}{1,000} + \frac{\widehat{\beta}_2 faminc}{1,000} = \widehat{\alpha}_0 + \widehat{\alpha}_1 cigs + \widehat{\alpha}_2 faminc$$

$$\rightarrow \widehat{\alpha}_j = \frac{\widehat{\beta}_j}{1,000}$$

2p. 35

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

$$bwght_g = \widehat{\beta}_0 + \widehat{\beta}_1 cigs + \frac{\widehat{\beta}_2 faminc_{USD}}{1,000}$$

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

$$\Rightarrow \widehat{\theta}_2 = \frac{\widehat{\beta}_2}{1,000}$$

the value of this variable is going to be 1000 times larger than $faminc$

in other words $\widehat{\theta}_2$ = impact of 1USD ↑ in income

$\widehat{\beta}_2$ = impact of 1000USD ↑ in income

• What if we use $bwght$ in kg & income in THB

$$bwght_{kg} = \frac{\widehat{\beta}_0}{1,000} + \frac{\widehat{\beta}_1}{1,000} cigs + \frac{\widehat{\beta}_2}{30,000} faminc_{THB}$$

↪ This value is going to be 30,000 times more than $faminc$

where

- 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

In the US or many other countries, students can apply to schools in the area without competition; taking any test. So, the lower stratio, the better the school.

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	
Total	84.582225	505	.167489554	Prob > F =	0.0000	
				R-squared =	0.6088	
				Adj R-squared =	0.6049	
				Root MSE =	.25725	

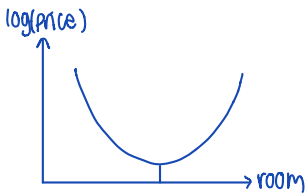
	lprice	lnox	dist	rooms	rooms_sq	stratio	_cons
Coef.		-.9767545	-.0321972	-.5528032	.0624697	-.0486667	13.59154
Std. Err.		.0995938	.0094013	.1612965	.0124867	.0058131	.5650901
t		-9.81	-3.42	-3.43	5.00	-8.37	24.05
P> t		0.000	0.001	0.001	0.000	0.000	0.000
[95% Conf. Interval]		-1.172429	-.050668	-.8697056	.0379368	-.0600879	12.4813
		-.7810806	-.0137264	-.2359007	.0870025	-.0372455	14.70178

(log(price))
(log(nox))

$|t| > 1.96$
 \rightarrow All variable are significant

Consider the effect of "room"

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



• at how many room dose 1 additional room has a positive impact on log(price)?
 $0 = -0.553 + 2(0.062)(\text{rooms})$
 $\text{room} = 4.4$
 answer \Rightarrow at ~~4.4~~ ^{5; round up} room or more

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

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

$$\frac{100 \cdot \frac{1}{\text{price}} \Delta \text{price}}{\Delta \text{room}} = 100(-0.553 + 2(0.062)(5))$$

$$= 100 \times 0.067 = 6.7\% \text{ increase.}$$

What about the Δ in price when #room increase from 5 to 7.

$$\text{for } 6-7 \rightarrow 100(-0.553 + 2(0.062)(6)) = 19.1\%$$

$$\text{total } \% \Delta \text{ in price} = 6.7\% + 19.1\% = 25.8\%$$

3 Models with Interaction Terms \Rightarrow used when the impact of one variable depends on the value (level) of another variable.

Consider

$$price = \beta_0 + \beta_1 \overset{\lambda_1}{sqr\ ft} + \beta_2 \overset{\lambda_2}{bdrms} + \beta_3 \overset{\lambda_3}{sqr\ ft \times bdrms} + \beta_4 \overset{\lambda_2}{bthrms} + u$$

where

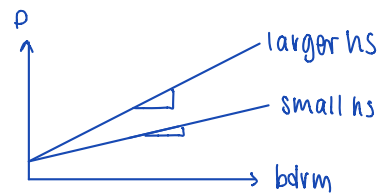
$price$ = housing price

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

$bdrms$ = number of bedrooms

$bthrms$ = number of bathrooms

$$\frac{d\ price}{d\ bdrms} = \beta_2 + \beta_3\ sqr\ ft$$



\Rightarrow if $\beta_3 > 0$ then, an additional bedroom would increase price more for a large 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"
 - (d.f. = free data point used to estimate the parameter)
 - \rightarrow 1 data point is sacrificed every time we estimate a parameter
- using just R^2 would not punish "having too many regressor"
- We use adjusted R^2 or \bar{R}^2 when we want to punish adding too many regressors.

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

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

\rightarrow If we have k , d.f. = $n-k-1 \downarrow$, $SSR/n-k-1 \uparrow \rightarrow adj_R^2 \downarrow$

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

Consider Model 1

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

$$= (223.90) \quad (0.0089) \quad (11.08)$$

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

Consider Model 2

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

$$= (0.29) \quad (0.033) \quad (0.004)$$

$$n = 209, \quad R^2 = 0.282, \quad \bar{R}^2 = 0.275 \quad \begin{array}{l} 27.5\% \text{ of variation in } Y \\ \text{is explained} \\ \text{so, this model is better!} \end{array}$$

Multiple Regression Analysis with Qualitative Information:

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

$$wage = \beta_0 + \delta_0 female + \beta_1 educ + u. \quad (1)$$

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.

$$\begin{aligned} 1) E(wage | female, educ) &= E(\beta_0 + \delta_0 female + \beta_1 educ + u | female, educ) \\ &= \beta_0 + \delta_0 female + \beta_1 educ + E(u | female, educ) \\ &= \beta_0 + \delta_0 female + \beta_1 educ \quad \downarrow = 0 \quad \text{MLR 1-4 satisfied} \end{aligned}$$

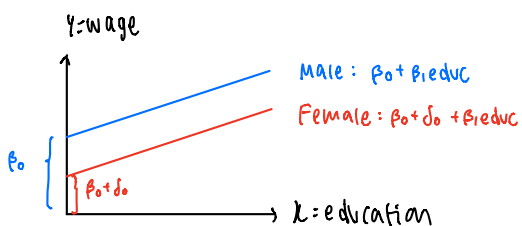
2) Thus

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

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

$$\begin{aligned} \delta_0 &= E(wage | female=1, educ) - E(wage | female=0, educ) \\ &\text{or } E(wage | female, educ) - E(wage | male, educ) \end{aligned}$$

* given the same value of educ (same education level), δ_0 is the different in the expected wage of female & male.



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

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

For example:

$$wage = \beta_0 + \beta_1 female + \beta_2 educ$$

↑
intercept x 1

$$\beta_0 = \beta_1 + \beta_3$$

$$1 = female + male$$

$$female = male + 1$$

id	female	male	β_0
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 "1" category, to avoid multi collinearity

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

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

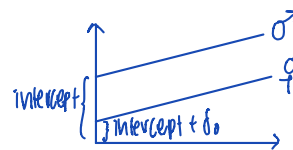
$$winter = \begin{cases} 1 & \text{if winter} \\ 0 & \text{otherwise} \end{cases}$$

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

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

in this case, male

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



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

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*.

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

} 1 female
} 0 otherwise

`regress lwage female married educ exper expersq tenure tenursq`

Source	SS	df	MS	Number of obs =	526
Model	65.6482326	7	9.37831895	F(7, 518) =	58.76
Residual	82.6815188	518	.159616832	Prob > F =	0.0000
				R-squared =	0.4426
				Adj R-squared =	0.4351
Total	148.329751	525	.28253286	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 difference between female & male workers given the same marital status and other factors

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

$$100 \cdot \frac{\frac{1}{\text{wage}} \Delta \text{wage}}{dfemale} = 100 \cdot -0.29$$

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

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

2) δ_0 measure the impact of being married (marriage premium)

But since $|t| < 1.96$ or $p > 0.05$, we do not reject H_0 of no impact

	♀	♂
ma	marrfem	marrmale
si	singfem	singlemale

Consider a model which includes dummy variables for each gender/marital status combination- *marrmale*, *marrfem* and *singlemale*. *singlemale* ← used as the basecase

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

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	

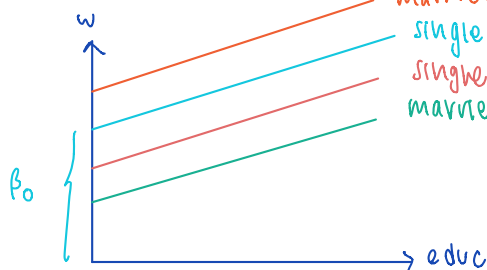
lwage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
marrmale	.2126757	.0553572	3.84	0.000	.103923	.3214284
marrfem	-.1982676	.0578355	-3.43	0.001	-.311889	-.0846462
singfem	-.1103502	.0557421	-1.98	0.048	-.219859	-.0008414
educ	.0789103	.0066945	11.79	0.000	.0657585	.092062
exper	.0268006	.0052428	5.11	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

This regression is not the same as the previous one as it uses "singlemale" as the basegroup,
 Comments: the previous one use "male and single" as the base group

• δ_0 measure the expect diff in wage of married male as compared to single male, holding other factors constant

• δ_1 measure the expected diff " female " female, "

• δ_2 → same rationale intercepts
 married male $\beta_0 + \delta_0 = 0.321 + 0.2127$
 single male β_0
 single fem $\beta_0 + \delta_2 = 0.321 - 0.110$
 married fem $\beta_0 + \delta_1 = 0.321 - 0.198$



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 cases the "range of value" serves as a better explanatory variable than the "value" itself

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

↳ e.g. age may explain the model better if split into generations

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:

rank	top 10	11-25
1	1	0
2	1	0
3	1	0
...
10	1	0
11	0	1
...
25	0	1
...
40	0	1
...

1) do measure the diff in expected $\log(\text{salary})$ of law-school graduate from a top10 university compared to expected $\log(\text{salary})$ of those who graduated from school ranked 61th and worse

2) $\delta_i \rightarrow$ same rational