

# 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

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

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 ??

$$1 \text{ kg.} = 1,000 \text{ g.}$$

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

$$= \hat{\alpha}_0 + \hat{\alpha}_1 cigs + \hat{\alpha}_2 faminc.$$

$$\rightarrow \hat{\alpha}_0 = \frac{\hat{\beta}_0}{1,000}, \hat{\alpha}_1 = \frac{\hat{\beta}_1}{1,000}, \hat{\alpha}_2 = \frac{\hat{\beta}_2}{1,000}$$

- What if we use  $faminc$  in USD (instead of 1,000 USD)

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

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

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

$$\rightarrow \hat{\theta}_2 = \frac{\hat{\beta}_2}{1,000}$$

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

$$\hat{\beta}_2 = \underbrace{\quad}_{1,000} \text{ USD} \uparrow \text{ in income}$$

- What if we use  $bwght$  in Kg & income in THB

$$\widehat{bwght}_{kg} = \frac{\hat{\beta}_0}{1,000} + \frac{\hat{\beta}_1}{1,000} cigs + \left( \frac{\hat{\beta}_2}{1,000} \right) faminc_{THB}$$

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

$$\frac{d \ln x}{dx} = \frac{1}{x} \rightarrow d \ln x = \frac{1}{x} dx$$

2 More on functional forms

- Logarithmic Functional Form

$Y$  nonlog  $X$  log

usually mean natural log

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

$\Delta Y = Y_1 - Y_2$   
 $\Delta X = X_{11} - X_{12}$

$$\beta_1 = \frac{d \log(y)}{d \log(x_1)} = \frac{\frac{1}{y} dy}{\frac{1}{x_1} dx_1} = \frac{\frac{1}{y} \Delta Y}{\frac{1}{x_1} \Delta X} = \frac{100 \times \frac{1}{y} \Delta Y}{100 \times \frac{1}{x_1} \Delta X} = \frac{\% \Delta Y}{\% \Delta X}$$

with the log y & log x format, the coefficient is going to be the elasticity! ( $x_1$  elasticity of y)  
(price) (demand)

$$\beta_2 = \frac{d \log(Y)}{d x_2} = \frac{\frac{1}{Y} dY}{d x_2} = \frac{\frac{1}{Y} \Delta Y}{\Delta X_2}$$

→ if we want the upper term to be 1% change, then

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

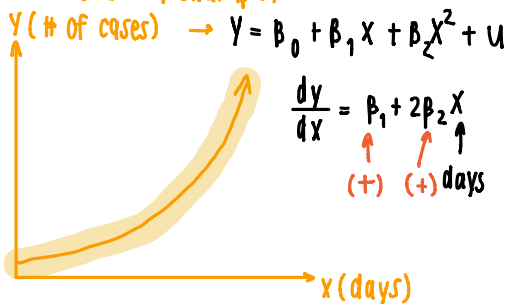
$$100 \beta_2 = \frac{\% \Delta Y}{\Delta X_2}$$

$100 \beta_2 = \% \Delta$  in Y given that  $X_2$  increases by 1 unit.

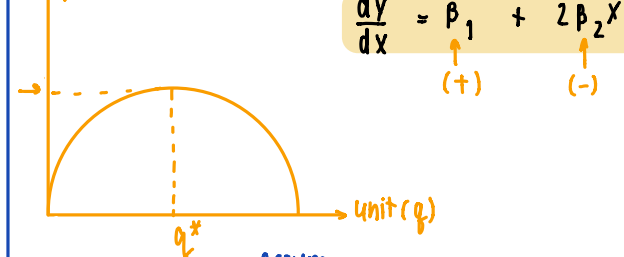
- Models with Quadratics (squares)

→ Capture increasing / decreasing marginal effects (slope of the relationship between X & Y is not constant.)

COVID-19 example.



Decreasing returns →  $y = \beta_0 + \beta_1 X + \beta_2 X^2 + u$



Assume  $\pi = (p - mc)q$ ;  $mc = 10$   
 $\pi = (100 - q - 10)q$  demand:  $p = 100 - q$   
F.O.C.  $\frac{\partial \pi}{\partial q} = 0 = 90 - 2q$   $\beta_1$  is positive  $\beta_2$  is (-)

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

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

In the US or many other countries, students can apply to schools in the area without having to take any test. So, 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	
Total	84.582225	505	.167489554	Prob > F =	0.0000	
				R-squared =	0.6088	
				Adj R-squared =	0.6049	
				Root MSE =	.25725	

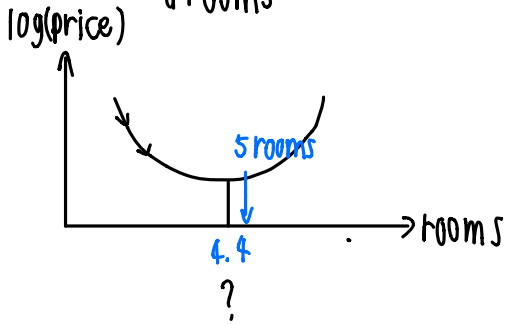
	lprice	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
$\beta_1$	lnox	-.9767545	.0995938	-9.81	0.000	[-1.172429, -.7810806]
$\beta_2$	dist	-.0321972	.0094013	-3.42	0.001	[-.050668, -.0137264]
$\beta_3$	rooms	-.5528032	.1612965	-3.43	0.001	[-.8697056, -.2359007]
$\beta_4$	rooms_sq	.0624697	.0124867	5.00	0.000	[.0379368, .0870025]
$\beta_5$	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  $< 0.05$

~ all variable 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 has a positive impact on log(price)??

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

$$\text{rooms} = 4.4$$

Answer  $\rightarrow$  at 4.4 rooms or more

at  $\sim 5$  rooms or more.

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

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

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

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

What about % change in price when # rooms increases from 5 to 7? ??

$$\% \Delta \text{price} = 100(-0.553 + 2(0.062) \cdot 6) = 19.1\%$$

Total %  $\Delta$  in price when # rooms  $\uparrow$  from 5 to 7 is  $6.7 + 19.1\% = 25.8\%$

### 3 Models with Interaction Terms → Used when the impact of one variable depends on the value (level) of another variable

Consider

$$price = \beta_0 + \beta_1 \underset{X_1}{sqr\ ft} + \beta_2 \underset{X_2}{bdrms} + \beta_3 \overset{X_3}{\underbrace{sqr\ ft \times bdrms}} + \beta_4 \underset{X_4}{bthrms} + u$$

where

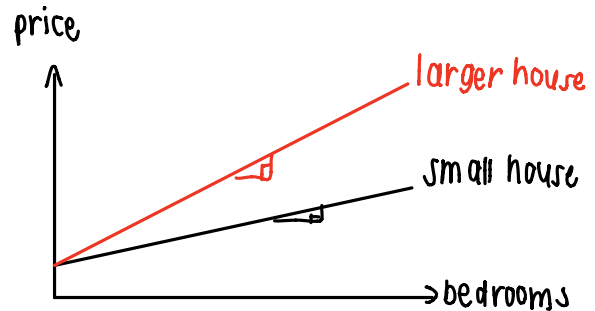
*price* = housing price

*sqr ft* = house size (square feet)

*bdrms* = number of bedrooms

*bthrms* = number of bathrooms

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



→ If  $\beta_2 > 0$  then, an additional bedrooms would increase price more for a larger house!

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

- Adding more regressors ALWAYS improve fit  $\rightarrow R^2$  always  $\uparrow$
- But we lose the "degree of freedom"  
(d.f. = free data point used to estimate the parameter)  
 $\leadsto$  1 data point is sacrificed every time we estimate a parameter
- using  $R^2$  would not punish "having too many regressors"
- we use adjusted  $R^2$  or  $\bar{R}^2$  when we want to punish adding too many regressors.

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

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

$\rightarrow$  If d.f =  $n-k-1 \downarrow$ ,  $SSR/(n-k-1) \uparrow$ ,  $\text{adj-}R^2 \downarrow$

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

Consider Model 1

$$\begin{aligned} \widehat{\text{salary}} &= 830.63 + 0.0163\text{sales} + 19.63\text{roe} \\ & \quad (223.90) \quad (0.0089) \quad (11.08) \\ n &= 209, \quad R^2 = 0.029, \quad \bar{R}^2 = 0.020 \end{aligned}$$

Consider Model 2

$$\begin{aligned} \log(\widehat{\text{salary}}) &= 4.36 + 0.2751 \log(\text{sales}) + 0.0179\text{roe} \\ & \quad (0.29) \quad (0.033) \quad (0.004) \\ n &= 209, \quad R^2 = 0.282, \quad \bar{R}^2 = 0.275 \end{aligned}$$

$\hookrightarrow$  27.5% of variation in  $y$  is explained.  
So, this model is better!



# Multiple Regression Analysis with Qualitative Information:

info that represent quality  
 e.g. gender, season, month

## 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

only effect value of intercept

Consider



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

where

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

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.

$$\begin{aligned} \textcircled{1} E(\text{wage} | \text{female}, \text{educ}) &= E(\beta_0 + \delta_0 \text{female} + \beta_1 \text{educ} + u | \text{female}, \text{educ}) \\ &= \beta_0 + \delta_0 \text{female} + \beta_1 \text{educ} + E(u | \text{female}, \text{educ}) \\ &= \beta_0 + \delta_0 \text{female} + \beta_1 \text{educ} \end{aligned}$$

↓  
= 0  
(assump. MLR1-4  
holds)

② thus

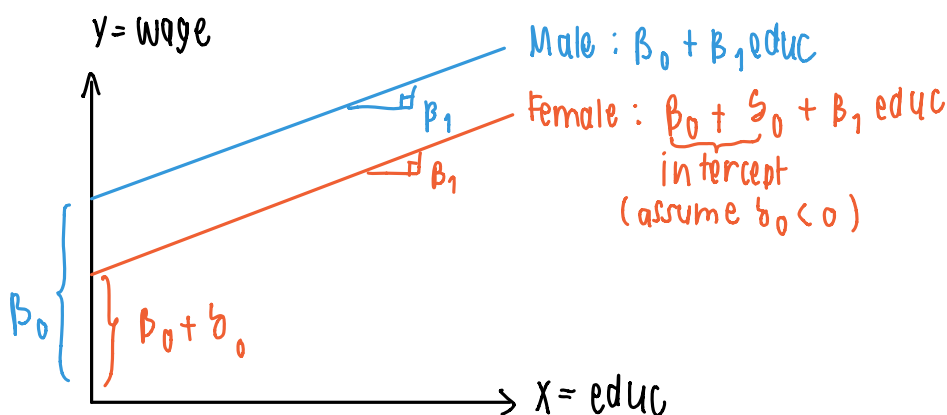
$$\textcircled{f} : E(\text{wage} | \text{female}=1, \text{educ}) = \beta_0 + \delta_0(1) + \beta_1 \text{educ} = \beta_0 + \delta_0 + \beta_1 \text{educ}$$

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

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

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

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



→ 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 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.

$$wage = \beta_0 X_0 + \beta_1 female + \beta_2 educ + \beta_3 male + u$$

$(X_1)$                        $(X_2)$                        $(X_3)$

For example:

↑ intercept + 1

$$X_0 = X_1 + X_3$$

$$1 = female + male$$

$$female = male + 1$$

Id	female	male	$X_0$
1	1	0	1
2	1	0	1
3	1	0	1
4	0	1	1
⋮	⋮	⋮	⋮
99	1	0	1

or

If there are "n" categories, we omit "1" category to avoid multicollinearity

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

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

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

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

etc

id	winter	spring	summer	fall	$X_0$
1	1	0	0	0	1
2	1	0	0	0	1
3	0	1	0	0	1
4	0	1	0	0	1
⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮

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

$\begin{cases} 1 & \text{if female} \\ 0 & \text{if male} \end{cases}$        $\begin{cases} 1 & \text{if male} \\ 0 & \text{if female} \end{cases}$

↑ in this case, male

. 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	<b>-0.3251146</b>	.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 compare 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.

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

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

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
Total	148.329751	525	.28253286	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

2.  $\delta_1$  measure the impact of be married.

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

Comments:

1.)  $\delta_0$  measures the expected difference between female & male workers given the same marital status and other factors

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

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

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

• female worker are expected to earn less than male workers by 29.02%, holding other factors the same.

	♀	♂
marr	marrfem	marmale
sing	singfem	singmale

Base case

Consider a model which includes dummy variables for each gender/marital status combination- marmale, marrfem and singfem. (Or singmale used as the base case)

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

```
regress lwage marmale 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	

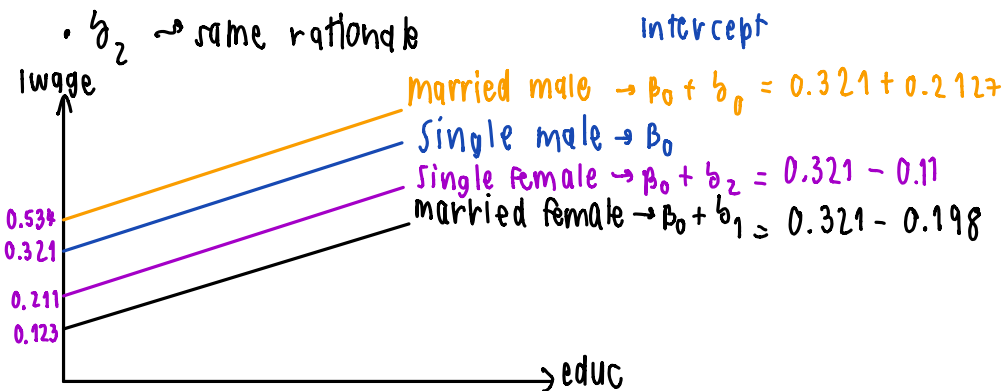
  

lwage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
marmale	.2126757	.0553572	3.84	0.000	.103923	.3214284
marrfem	-.1982676	.0578355	-3.43	0.001	-.311889	-.0846462
singfem	-.1103502	.0557421	-1.98	0.048	-.219859	-.0008414
educ	.0789103	.0066945	11.79	0.000	.0657585	.092062
exper	.0268006	.0052428	5.11	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. It uses "single male" as the base group. (The previous one use male & single as 2 base groups.)

- $b_0$  measure the expected diff. in wage of married male as compared with single males, holding other factor constant.
- $b_1$  measure the expected diff. in wage of married female as compared with single males, holding other factor constant.

•  $b_2$  same rationale



**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" serve as a better explanatory variables 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 Young 0-15 gen z 16-19 etc.

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
				R-squared =	0.8833
				Adj R-squared =	0.8759
Total	10.3763518	135	.076861865	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

the baseline is ranking 61<sup>th</sup> and worse

Comments:

rank	top10	r11-25	r26-40	
1	}	0	0	
2		0	0	
3		0	0	
⋮				
10			0	
11		0	}	0
⋮		0		0
25		0		0
⋮				
40		0		0
⋮				
⋮		0		
⋮		0		

1)  $\delta_0$  measures the difference 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 61<sup>th</sup> and worse

2)  $\delta_1$  use the same rationale