

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

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

$$bwght_kg^{\wedge} = bwght_g^{\wedge} / 1000 = (\beta^{\wedge}0 / 1000) + (\beta^{\wedge}1 / 1000) cigs + (\beta^{\wedge}2 / 1000) faminc \\ = \alpha0^{\wedge} + \alpha1^{\wedge} cigs + \alpha2^{\wedge} faminc$$

$$\Rightarrow \alpha0^{\wedge} = (\beta^{\wedge}0 / 1000), \alpha1^{\wedge} = (\beta^{\wedge}1 / 1000), \alpha2^{\wedge} = (\beta^{\wedge}2 / 1000)$$

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

$$bwght_g = \beta^{\wedge}0 + \beta^{\wedge}1 cigs + \beta^{\wedge}2 / 1000 faminc_usd \\ = \beta^{\wedge}0 + \beta^{\wedge}1 cigs + \theta^{\wedge}2 faminc_usd$$

$$\Rightarrow \theta^{\wedge}2 = \beta^{\wedge}2 / 1000$$

in other words $\theta^{\wedge}2$ = impact of 1 USD increase in income

$\beta^{\wedge}2$ = impact of 1000 USD increase in income

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

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

$$bwght_kg = \beta^{\wedge}0 / 1000 + \beta^{\wedge}1 / 1000 cigs + \beta^{\wedge}2 / 1000 faminc_thb$$

2 More on functional forms

- Logarithmic Functional Form

$$\Delta y = Y1 - Y2$$

$$\Delta x1 = X11 - X12$$

usually means natural log

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

$$\beta_1 = d \log(y) / d \log(x)$$

$$= (1/y)dy / (1/x_1)dx_1$$

$$= (1/y)\Delta y / (1/x_1)\Delta x_1$$

$$= [100(1/y)]\Delta y / [100(1/x_1)]\Delta x_1$$

$$= \%y/\%x$$

with the log y & log x format, the coefficient is going to be the elasticity.
(x1 elasticity of y)
price demand

$$\beta_2 = d \log(y) / dx_2$$

$$= (1/y)dy / dx_2$$

$$= (1/y)\Delta y / \Delta x_2$$

-> If we want the upper term to be % change then

$$100 \beta_2 = [100(1/y)]\Delta y / \Delta x_2$$

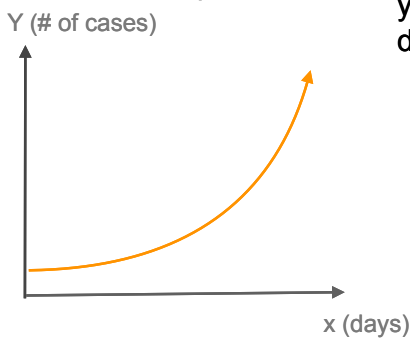
$$100 \beta_2 = \% \Delta y / \Delta x_2$$

100 β2 = %Δ in y given that x2 increases by 1 unit

- Models with Quadratics

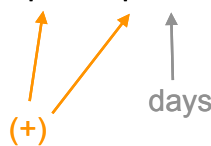
-> capture increasing / decreasing marginal effects (slope of the relationship between x & y is not constant)

COVID-19 example

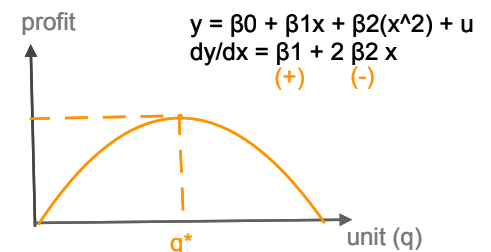


$$y = \beta_0 + \beta_1 x + \beta_2 (x^2) + u$$

$$dy/dx = \beta_1 + 2 \beta_2 x$$



Decreasing returns



$$y = \beta_0 + \beta_1 x + \beta_2 (x^2) + u$$

$$dy/dx = \beta_1 + 2 \beta_2 x$$

(+) (-)

$$\pi = (p - mc) q$$

$$\pi = (100 - q(\text{hat}) - 10) q$$

$$F.O.C \ d\pi/dq = 90 - 2Q$$

Assume mc=10
Demand: P=100-q

Example : Effects of Pollution on Housing Prices

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

page 76

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

↑
with the log y & log x format, the coefficient is going to be the elasticity (x₁ elasticity of y)
(price) (demand)

$$\bullet \beta_2 = \frac{d(\log(y))}{dx_2} = \frac{\frac{1}{y} dy}{dx_2} = \frac{\frac{1}{y} \Delta y}{\Delta x_2}$$

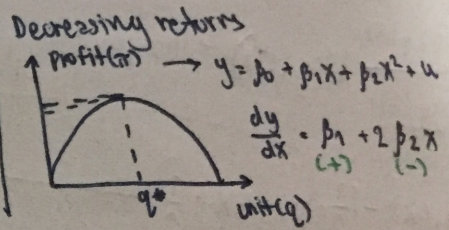
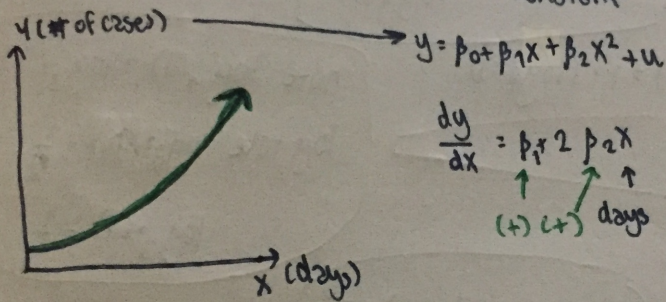
→ if we want the upper term to be % change,

then $100 \beta_2 = \frac{100 \frac{1}{y} \Delta y}{\Delta x_2}$

$$100 \beta_2 = \frac{\% \Delta y}{\Delta x_2} \quad \left| \quad 100 \beta_2 = \% \Delta \text{ in } y \text{ given that } x_2 \text{ increases by 1 unit} \right.$$

→ capture increasing / decreasing marginal effects (slope of the relationship between x & y is not constant)

COVID-19 example



$$\pi = (p - mc)q$$

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

Assume $mc = 10$
demand $p = 100 - q$

$$FOC \frac{d\pi}{dq} = 90 - 2q \rightarrow \beta_2 \text{ is } (-)$$

β_1 is (+)

when adding quadratics to the same x this will make the relationship btw x & y be non linear and can capture the increasing and decreasing marginal effects

where

- 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]	
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 variables are significant

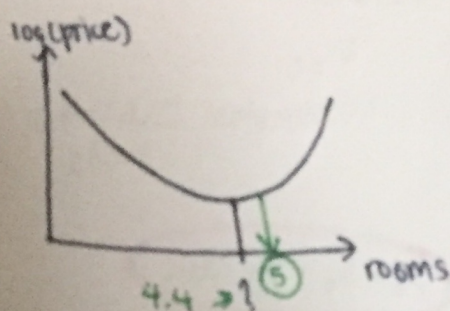
all < 0,05

When t statistic reject null P would also reject

Consider the effect of "room"

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

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



at how many rooms does 1 additional rooms has a positive impact on $\log(\text{price})$

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

$$\text{rooms} = 4.4$$

Answer → at 4.4 rooms or more
 \approx 5 rooms or more

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

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

$$= 100 \times 0.067 = 6.77 \text{ increase}$$

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

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

$$= 19.1\%$$

when #rooms ↑ from 5 → 7
 total % Δ in price 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

$$price = \beta_0 + \beta_1 \underset{X1}{sqr\ ft} + \beta_2 \underset{X2}{bdrms} + \beta_3 \overset{X3}{sqr\ ft \times bdrms} + \beta_4 \underset{X2}{bthrms} + u$$

where

price = housing price

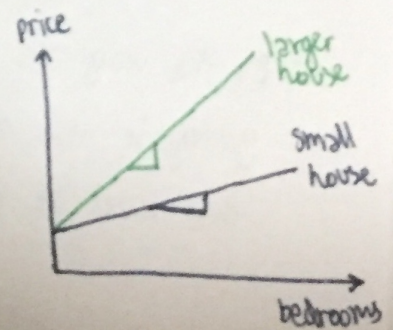
sqr ft = house size (square feet)

bdrms = number of bedrooms

bthrms = number of bathrooms

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

→ if $\beta_2 > 0$, then an additional bedroom 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 R^2 always increase

But if we lose the “degree of freedom”

(d.f. = free data point used to estimate the parameter)

→ 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 R^2 when we want to punish adding too many regressors

$$R^2 = 1 - (SSR/SST) = 1 - (SSR/n / SST/n)$$

$$\text{Adj. } R^2 = 1 - [SSR/(n-k-1) / SST / (n-1)]$$

If we have more k, d.f. = n-k-1 decrease

SSR/ (n-k-1) increases, adj- R^2 decrease

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

Consider Model 1

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

Consider Model 2

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

27.5% of variation in Y is explained
So, this model is better

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

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

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.

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

$\downarrow = 0$
(assumption MLR1-4)
holds

② Thus,

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

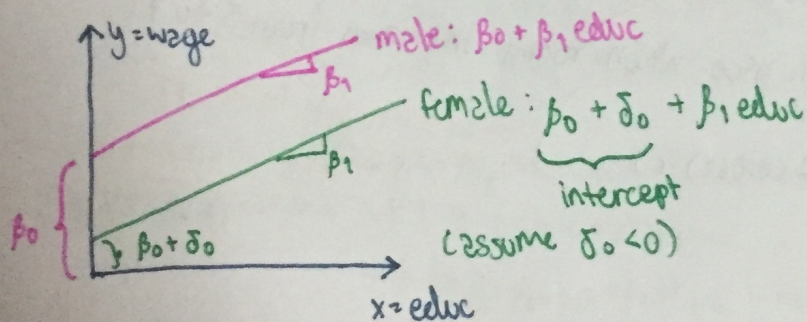
women ♂ : $E(\text{wage} | \text{female} = 0, \text{educ}) = \beta_0 + \delta_0(0) + \beta_1 \text{educ} = \beta_0 + \beta_1 \text{educ}$

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

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

* given the same value of educ (same education level),

δ_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.

$$\text{wage} = \beta_0 X_0 + \delta_0 \text{female} + \beta_1 \text{educ} + \delta_1 \text{male} + u$$

For example:

$$\begin{aligned} & \uparrow \\ & \text{intercept} = 1 \\ & x_0 = x_1 + x_3 \\ & 1 = \text{female} + \text{male} \\ & \text{female} = \text{male} + 1 \end{aligned}$$

or

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

$$\begin{aligned} 1 &= \text{winter} + \text{spring} + \text{summer} + \text{fall} \\ \text{winter} &= 1 - \text{spring} - \text{summer} - \text{fall} \end{aligned}$$

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

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

etc.

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

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

Source	SS	df	MS	Number of obs = 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

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 if female
0 if otherwise 1 if married
0 if 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		
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

$\hat{\beta}$ {

Comments:

1) δ_0 measures the expected difference between female & male workers given the same marital status and other factors

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

• female workers are expected to earn less than male workers by

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

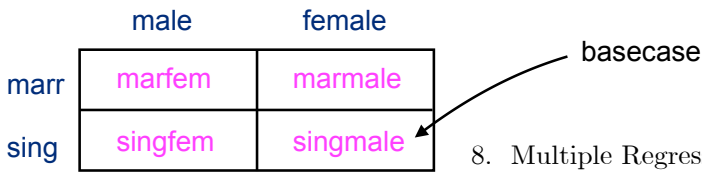
29.02%, holding other factors same

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

2) δ_0 measures 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



Consider a model which includes dummy variables for each gender/marital status combination– *marrmale*, *marrfem* and *singfem*. (*singmale* <- used as the basecase)

$$\log(wage) = \beta_0 + \delta_0marrmale + \delta_1marrfem + \delta_3singfem + \beta_1educ + \beta_2exper + \beta_3exper^2 + \beta_4tenure + \beta_5tenure^2 + u. \quad (8.1)$$

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

$\hat{\beta}$ {

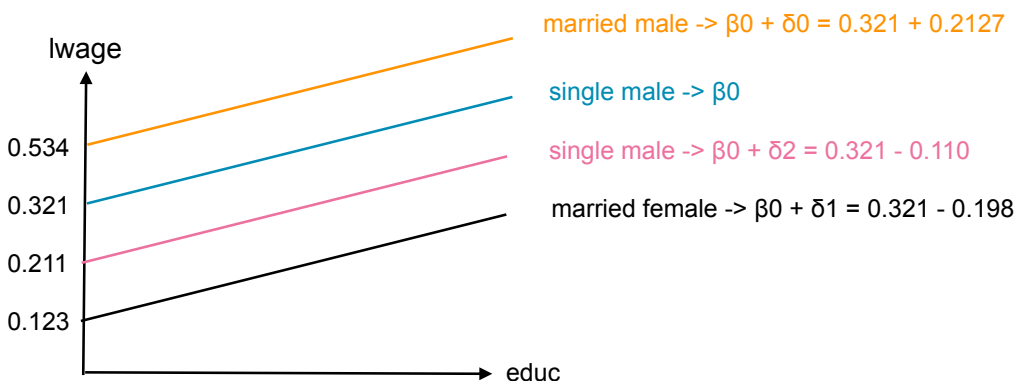
Comments:

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)

- δ_0 measures 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 -> 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 variable than the "value" itself. eg age may be explain the model better if split into generations young (0-15) genz (16 - 29) etc.

```
. 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

the baseline is ranking 61st and worse

Comments:

1) δ_0 measures the difference in expected log (salary) of a law-school graduate from a top 10 university compared to expected log (salary) of those who graduated from the school ranked 61st and worse

2) δ_1 → use the same rationale