

## F-test motivation

⇒ We want to test the significance of a group of hypotheses (multiple hypotheses)

$$\text{Grade}_{325} = \beta_0 + \beta_1 \# \text{times\_front} + \beta_2 \# \text{times\_back} \\ + \beta_3 \text{hr\_study} + \beta_4 \text{past\_GPA} + \beta_5 \text{gender} + u$$

$H_0$ : seat position doesn't have impact on GPA

$$\beta_1 = 0 \text{ and } \beta_2 = 0 \Rightarrow \beta_1 = \beta_2 = 0$$

$H_a$ : seat position matters

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

7 Testing Multiple Linear Restrictions: The F-test

Suppose the model is specified by

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

$$H_0 : \beta_2 = 0 \text{ and } \beta_3 = 0$$

$$H_a, H_1 : H_0 \text{ is not true}$$

Want to test if  $x_1$  and  $x_2$  BOTH have no impact on  $y$ .

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

Big model

1. Our full model is called the "unrestricted" model (ur). Suppose it can be expressed as:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + u \quad \text{is true} \Rightarrow \text{Reject } H_0$$

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

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

$$y = \beta_0 + \beta_1 x_1 + u \quad \text{is true} \Rightarrow \text{do not reject } H_0$$

Suppose there are "q" number of  $\beta$  that we would like to perform a joint-test of = 0  
 - e.g. in this model  $q = 2$

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

$$H_0 : \beta_{k-q+1} = \beta_{k-q+2} = \dots = \beta_k = 0$$

(the last q  $\beta_s = 0$ )

$H_a$ :  $H_0$  is not true.

$$y = \underbrace{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{k-q} x_{k-q}}_{(r)} + \beta_{k-q+1} x_{k-q+1} + \beta_{k-q+2} x_{k-q+2} + \dots + \beta_k x_k + u$$

ur

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

$SSR_{ur}$

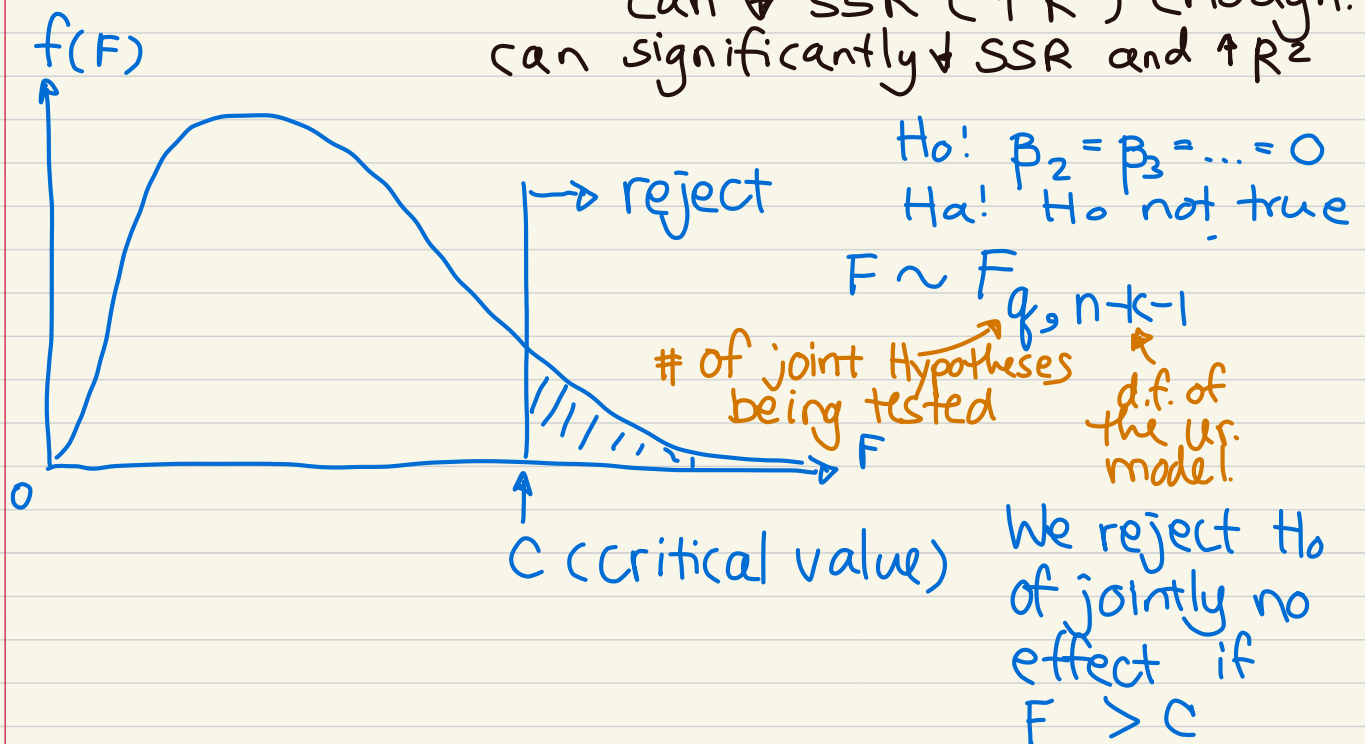
$(n - k - 1)$

This is always (+) b/c  $SSR_{ur} < SSR_r$ . Every time you add 1 more  $x$ , the model will be better explained.

d.f. of the "ur" model.

- 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??

$\Rightarrow$  Because every time we add 1 more  $X$ ,  $\text{var}(\hat{\beta}_s)$  will increase, making the prediction of  $\beta$  less precise. So, we only keep the addition  $X_s$  if it/they can improve the model enough  
 can  $\downarrow$   $SSR$  ( $\uparrow$   $R^2$ ) enough.  
 can significantly  $\downarrow$   $SSR$  and  $\uparrow$   $R^2$



3. Some useful facts

①  $R^2_{ur} > R^2_r$  because any additional X would increase  $R^2$  (improve fit).  
 $\Rightarrow SSR_{ur} < SSR_r$

② By including more X, 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  $R^2 = 1 - \frac{SSR}{SST}$   $\begin{matrix} \nearrow \text{RSS} \\ \searrow \text{TSS} \end{matrix}$

We have  $F = \frac{(R^2_{ur} - R^2_r)}{\dots}$

# of  $\beta$  that are set to "0"  $\rightarrow q$   $\frac{(1 - R^2_{ur})}{n - k - 1}$  intercept.  
 $\uparrow$  # of obs.  $\leftarrow$  # of slope  $\beta$ .

$\Rightarrow$  If we want to test the overall significance of the model

$H_0: \beta_1 = \beta_2 = \beta_3 = \dots = \beta_k = 0$  ,  $H_a = \text{otherwise}$

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

$R^2$  of the model  $\approx UR$   
 the "r" model has no X at all.

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

- $r \left\{ \begin{matrix} ur \\ \end{matrix} \right. \left\{ \begin{matrix} \text{salary} & = & \text{season salary} \\ \text{years} & = & \text{years in major leagues} \\ \text{gamesyr} & = & \text{games per year in the league} \\ \text{bavg} & = & \text{career batting average} \\ \text{hrunsyr} & = & \text{homeruns per year} \\ \text{rbisyr} & = & \text{runs batted in per year} \end{matrix} \right.$

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

$H_0: \beta_{bavg} = \beta_{hrunsyr} = \beta_{rbisyr} = 0$   
 $H_a: \text{otherwise is true}$

- the unrestricted model (ur) is defined by

ur model  
 $Y = X\beta + \epsilon$

```
. regress log_salary years gamesyr bavg hrunsyr rbisyr
```

Source	SS	df	MS	
Model	308.989208	5	61.7978416	Number of obs = 353
Residual	183.186327	347	.527914487	F( 5, 347) = 117.06
Total	492.175535	352	1.39822595	Prob > F = 0.0000

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	.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

```
. regress log_salary years gamesyr
```

Source	SS	df	MS	
Model	293.864058	2	146.932029	Number of obs = 353
Residual	198.311477	350	.566604221	F( 2, 350) = 259.32
Total	492.175535	352	1.39822595	Prob > F = 0.0000

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	.108312	103.62	0.000	11.01078 11.43683

When considering each of the performance X one-by-one, none of them has a significant impact at 5%

But when performing an F-test, performances have joint impact.

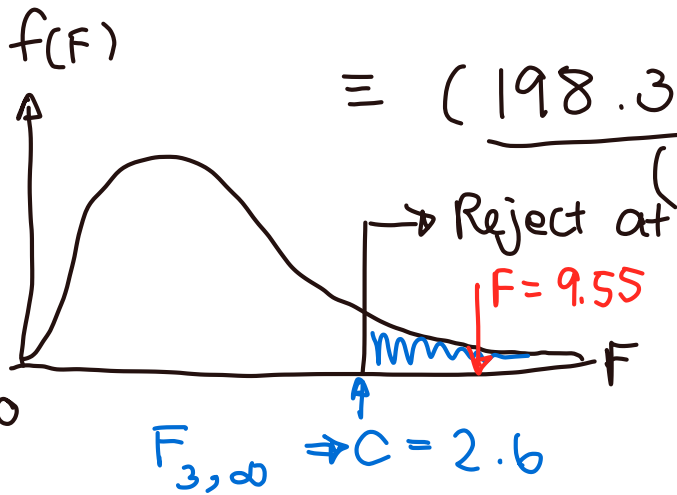
Now, our  $H_0$  and  $H_a$  becomes

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

HW:

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

$$\equiv \frac{(198.311 - 183.186) / 3}{(183.186) / (353 - 5 - 1)} \approx 9.55$$



Let's use 5% level of sig. Since  $F = 9.55 > 2.6$ , we reject  $H_0$  at 5% 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  $H_0$  *with 5% sig. level*

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

⇒ **When  $t < 1.96$**  we can say that  $\beta_i$  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)	(2)	(3)
$\log(\text{sales})$	.224 (.027)	.158 (.040)	.188 (.040)
$\log(\text{mktval})$	—	.112 (.050)	.100 (.049)
$\text{profmarg}$	—	-.0023 (.0022)	-.0022 (.0021)
$\text{ceoten}$	—	—	.0171 (.0055)
$\text{comten}$	—	—	-.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

*sales* →  
*other Company performance*  
*CEO characteristics*

↑  
*I like a simple regression with 1 X.*

# 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

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

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 *bweght* in kilograms??

1 kg. = 1,000 g.

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

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

$$\begin{aligned} bweght_g &= \hat{\beta}_0 + \hat{\beta}_1 cigs + \frac{\hat{\beta}_2}{1000} faminc_{USD} \\ &= \hat{\beta}_0 + \hat{\beta}_1 cigs + \hat{\theta}_2 faminc_{USD} \end{aligned}$$

$\Rightarrow \hat{\theta}_2 = \frac{\hat{\beta}_2}{1000}$   
 in other words  $\hat{\theta}_2$  = impact of 1 USD ↑ in income  
 $\hat{\beta}_2$  = 1000 USD ↑ in income.

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

- What if we use *bweght* in kg & income in THB

$$bweght_{kg} = \frac{\hat{\beta}_0}{1000} + \frac{\hat{\beta}_1}{1000} cigs + \left(\frac{\hat{\beta}_2}{30,000}\right) faminc_{THB}$$

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

2 More on functional forms

- Logarithmic Functional Form

$$\Delta Y = Y_1 - Y_2$$

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

usually means natural log!

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

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

with the log y & log x format, the coefficient is going to be the elasticity! (X<sub>1</sub> elasticity of Y) (price) (demand)

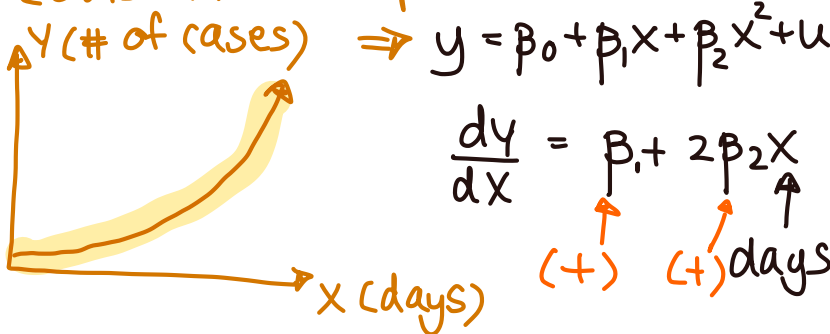
$$\beta_2 = \frac{d \log(Y)}{d X_2} = \frac{\frac{1}{Y} dY}{d X_2} = \frac{\frac{1}{Y} \Delta Y}{\Delta X_2}$$

⇒ if we want the upper term to be % change, then  $100 \beta_2 = \frac{100 \frac{1}{Y} \Delta Y}{\Delta X_2}$   $100 \beta_2 = \% \Delta$  in Y given that X<sub>2</sub> increases by 1 unit.

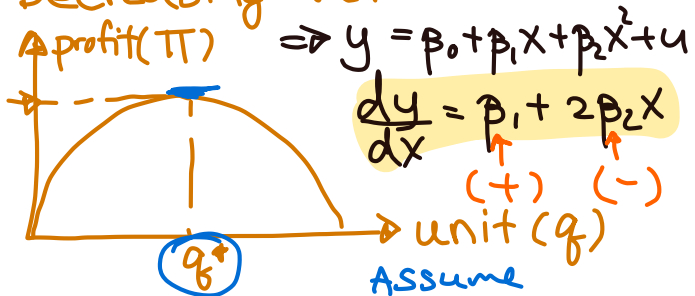
- Models with Quadratics (Squares)

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

COVID-19 example



Decreasing returns.



Assume  $\pi = (p - mc)q$ ;  $mc = 10$   
 Demand:  $P = 100 - q$   
 $\pi = (100 - q - 10)q$   $\beta_1$  is positive

F.o.c  $\frac{\partial \pi}{\partial q} = 0 = 90 - 2q$   $\beta_2$  is (-)

Example : Effects of Pullution on Housing Prices

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

## Derivatives of exponential and logarithmic functions [edit]

$$\frac{d}{dx} (c^{ax}) = ac^{ax} \ln c, \quad c > 0$$

the equation above is true for all  $c$ , but the derivative for  $c < 0$  yields a complex number.

$$\frac{d}{dx} (e^{ax}) = ae^{ax}$$

$$\frac{d}{dx} (\log_c x) = \frac{1}{x \ln c}, \quad c > 0, c \neq 1$$

the equation above is also true for all  $c$ , but yields a complex number if  $c < 0$

$$\frac{d}{dx} (\ln x) = \frac{1}{x}, \quad x > 0.$$

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

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

$$\frac{d}{dx} (\ln |x|) = \frac{1}{x}.$$

$$\frac{d}{dx} (x^x) = x^x (1 + \ln x).$$

$$\frac{d}{dx} (f(x)^{g(x)}) = g(x) f(x)^{g(x)-1} \frac{df}{dx} + f(x)^{g(x)} \ln(f(x)) \frac{dg}{dx}, \quad \text{if } f(x) > 0, \text{ and if } \frac{df}{dx} \text{ and } \frac{dg}{dx} \text{ exist.}$$

$$\frac{d}{dx} (f_1(x)^{f_2(x)^{\dots} f_n(x)}) = \left[ \sum_{k=1}^n \frac{\partial}{\partial x_k} (f_1(x_1)^{f_2(x_2)^{\dots} f_n(x_n)}) \right] \Big|_{x_1=x_2=\dots=x_n=x}, \text{ if } f_{i < n}(x) > 0 \text{ and } \frac{df_i}{dx} \text{ exists.}$$

### Logarithmic derivatives [edit]

The **logarithmic derivative** is another way of stating the rule for differentiating the **logarithm** of a function (using the chain rule):

$$(\ln f)' = \frac{f'}{f} \quad \text{wherever } f \text{ is positive.}$$

**Logarithmic differentiation** is a technique which uses logarithms and its differentiation rules to simplify certain expressions before

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 having to take 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	

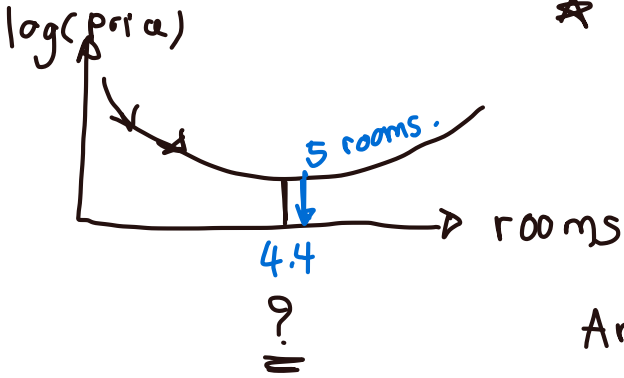
  

	price	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
log(price)	price					
	lnox	$\beta_1$ -0.9767545	.0995938	-9.81	0.000	-1.172429 -0.7810806
log(nox)	dist	$\beta_2$ -0.0321972	.0094013	-3.42	0.001	-.050668 -.0137264
	rooms	$\beta_3$ -0.5528032	.1612965	-3.43	0.001	-.8697056 -.2359007
	rooms_sq	$\beta_4$ 0.0624697	.0124867	5.00	0.000	.0379368 .0870025
	stratio	$\beta_5$ -0.0486667	.0058131	-8.37	0.000	-.0600879 -.0372455
	_cons	13.59154	.5650901	24.05	0.000	12.4813 14.70178

$|t| > 1.96$   $\uparrow$   $\uparrow$  all  $< 0.05$   
 $\rightarrow$  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 dose 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  $\rightarrow$  5 rooms or more.

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

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

total %  $\Delta$  in price

when # rooms  $\uparrow$  from 5 to 7 is 6.7 + 19.1%
- $100 \cdot \frac{1}{\text{price}} \frac{d \text{price}}{d \text{rooms}} = 100 (-0.553 + 2(0.062) \cdot 5) = 25.8\%$
- $= 100 \times 0.067 = 6.7\%$  increase.
- $\rightarrow$  What about % in price when # rooms increases from 5 to 7??
- $\% \Delta \text{ price} = 100 (-0.553 + 2(0.062) \cdot 6) = 19.1\%$

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

Consider

$$\text{price} = \beta_0 + \beta_1 \underset{X_1}{\text{sqrft}} + \beta_2 \underset{X_2}{\text{bdrms}} + \beta_3 \underset{X_1 \cdot X_2}{\text{sqrft} \times \text{bdrms}} + \beta_4 \underset{X_2}{\text{bthrms}} + u$$

where

*price* = housing price

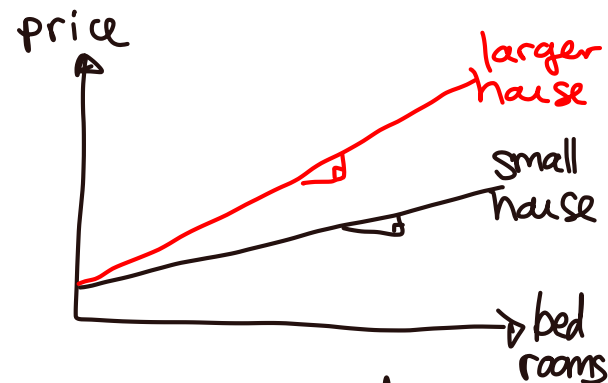
*sqrft* = house size (square feet)

*bdrms* = number of bedrooms

*bthrms* = number of bathrooms

$$\bullet \frac{\partial \text{price}}{\partial \text{bdrms}} = \beta_2 + \beta_3 \text{sqrft}$$

$\Rightarrow$  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  $\rightarrow R^2$  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  $R^2$  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.

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

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

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

Consider Model 1

If we have more  $k$ , d.f. =  $n-k-1 \downarrow$ ,  $SSR/(n-k-1) \uparrow$ ,  $adj-R^2 \downarrow$

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

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

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

Consider Model 2

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

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

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

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