

Probability & Statistics "Review"

* Goal: use data to answer economic questions

↳ almost never can we observe & measure every relevant data point

→ If we're interested in the effect of each year of schooling on income, it's extremely unlikely we'd get to observe every single person's income, much less years of schooling.

↳ And even that wouldn't be the ideal data set.

↳ Statistical Inference allows us to use a sample of the overall population

→ But in order to do this we need to go down the rabbit hole of theoretical probability & statistics for day.

* We often model real world outcomes using random variables.

↳ A variable whose possible values are outcomes of random phenomena.

Examples:

- Your grade in a given class is modeled as a Normal random variable
- The number of baskets you can make is ^{modeled as} a Binomial random variable.
- The number of inputs available to a firm ^{modeled as} is a Poisson random variable.

④ The probability of each value the RV can take on is called the probability mass (density) function

→ sometimes you can enumerate each outcome & its probability.

Example: Dice

X = the result of the roll of a single fair die.

$\sim \text{Uniform}(\frac{1}{6})$

X	$P(X=x)$
1	$\frac{1}{6}$
2	$\frac{1}{6}$
3	$\frac{1}{6}$
4	$\frac{1}{6}$
5	$\frac{1}{6}$
6	$\frac{1}{6}$

→ Sometimes you can't & you have to use a function.

Example: The time until you receive your next email is modeled as an expon. random var.
 $x \in [0, \infty)$

$$\text{pdf: } f(x) = \lambda e^{-\lambda x} = P(X=x)$$

✳ Since it isn't always particularly illuminating to write down a pdf & look at it to understand what's going on with a random variable,

↳ we often talk about:

- the expected value: $E[X]$
(the value we expect the variable to take on)
- the variance: $V[X]$
(how spread out the variable's values are)

Discrete Random Variables:

$$E[X] = \sum_{x=x} x P[X=x]$$

$$V[X] = \sum_{x=x} (x - E[X])^2 P[X=x]$$

Continuous Random Variables:

$$E[X] = \int x f(x) dx$$

$$V[X] = \int_{x \in \mathbb{R}} (x - E[X])^2 f(x) dx$$

Examples:

$\underline{x}_1 =$ Result of a single roll of a fair die.
 $E[x_1] ? (t)$

$$E[x_1] = \sum_{x_1=x} x P[x_1=x] \quad (\text{Always write the formula})$$

$$= 1(Y_6) + 2(Y_6) + 3(Y_6) + 4(Y_6) \\ + 5(Y_6) + 6(Y_6)$$

$$= \frac{1}{6} + \frac{2}{6} + \frac{3}{6} + \frac{4}{6} + \frac{5}{6} + \frac{6}{6} \\ = \frac{21}{6} = \underline{\underline{3.5}}$$

$$V[x_1] ? (i)$$

$$V[x_1] = \sum_{x_1=x} (x - E[x_1])^2 P[x_1=x]$$

$$= (1-3.5)^2(Y_6) + (2-3.5)^2(Y_6) + (3-3.5)^2(Y_6) \\ + (4-3.5)^2(Y_6) + (5-3.5)^2(Y_6) + (6-3.5)^2(Y_6)$$

$$= (6.25)(Y_6) + (2.25)(Y_6) + (0.25)(Y_6) \\ + (0.25)(Y_6) + (2.25)(Y_6) + (6.25)(Y_6) \\ = 2.91\overline{6}$$

X_2 = the time until your next email
 $\sim \text{Exp}(\lambda)$

$$E[X_2] = \int_{x \in \mathbb{X}} x f(x) dx$$

$$= \int_0^\infty x \lambda e^{-\lambda x} dx = \lambda \int_0^\infty x e^{-\lambda x} dx$$

$$= \frac{e^{-\lambda x} (\lambda x + 1)}{\lambda} \Big|_0^\infty$$

$$= \lim_{a \rightarrow \infty} \frac{e^{-\lambda a} (\lambda a + 1)}{\lambda} - \frac{e^{-\lambda(0)} (\lambda(0) + 1)}{\lambda}$$

$$= \lambda \left[\lim_{a \rightarrow \infty} e^{-\lambda a} (\lambda a + 1) \right] - \frac{(1)(1)}{\lambda}$$

L'Hopital

$$= \frac{1}{\lambda} - 0 = \frac{1}{\lambda}$$

Remember These E.V. Rules:

$$E[c] = c$$

$$E[aX] = aE[X]$$

$$E[aX+bY] = aE[X] + bE[Y]$$

* It's one thing to say I think student SAT scores are distributed according to some Distribution with mean μ & variance σ^2 .

$$Y \sim F(\mu, \sigma^2)$$

→ it's a whole other thing to say what we think μ is.

→ so we need to estimate μ using data.

* An estimator is a function of observed data used to estimate a statistical parameter.

Ex: We observe a bunch of SAT scores

hat's mean estimator $(y_1, y_2, y_3, \dots, y_n)$ } each a realization of random vars $y_i \sim F(\mu)$ "iid"

$$\hat{\mu} = f(y_1, y_2, \dots, y_n)$$

→ There are all sorts of estimators we could use.

$$\hat{\mu}' = \min\{y_1, y_2, \dots, y_n\}$$

↳ It wouldn't be a very good one but it would be a one.

* An unbiased estimator is one whose expectation is the true parameter value.

→ those little y 's are realizations of random variables

↳ functions of random variables are themselves random variables

↳ estimators are random variable & have expected values.

* An estimator, \hat{u} , of a true parameter u , is unbiased if

$$E[\hat{u}] = u.$$

* The estimator of u we usually use is the sample mean, \bar{x} .

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n y_i$$

Is \bar{x} an unbiased estimator of u .

→ need to show $E[\bar{x}] = u$.

$$E[\bar{x}] = E\left[\frac{1}{n}(y_1 + y_2 + \dots + y_n)\right]$$

$$= \frac{1}{n} E[y_1 + y_2 + \dots + y_n]$$

$$\begin{aligned}
 &= \frac{1}{n} \left(E[y_1] + E[y_2] + \dots + E[y_n] \right) \\
 &= \frac{1}{n} \underbrace{(\mu + \mu + \dots + \mu)}_{n} \\
 &= \frac{n\mu}{n} = \mu. \text{ Yes!}
 \end{aligned}$$

* \bar{X} is itself a random variable
it has a distribution.

Central Limit Theorem:

The sample mean is distributed according to Normal Distribution with mean μ and variance σ^2/n

$$\bar{X} \sim N(\mu, \sigma^2/n)$$

the associated standard deviation, σ/\sqrt{n} is called the standard error.

* Thing standard error is not.

- Sample standard deviation, s .
- the sample standard error, s/\sqrt{n}

*) Reminder of how to calculate S:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

$$S = \sqrt{S^2}$$

*) Keeping all our S's straight:

Population Variance: σ^2

↳ Std. Deviation: σ

Sample Variance: s^2

↳ Sample Std. Deviation: s

Population Standard Error: σ/\sqrt{n}
(theoretical)

Sample Standard Error: s/\sqrt{n}
(usually what you have)

T-P-S

Book price
data

→ differences
in prices across
gender?

→ genre play
a role?

Hypothesis Testing

Suppose we average the SAT^(Math) scores of all the students in our class, and get 620. (\pm std. dev. = 15)

↳ That's higher than the SMC average (610).

↳ That difference could be just due to randomness in sampling.
(you're just $\approx 20^{16}$ realizations of the Saint Mike's SAT random variable)

↳ Or it could be that the Economics majors' SAT scores are really higher than average.

* Hypothesis Testing is a process that allows us to say whether some characteristic of the data we observe (mean, std.dev, min, max, etc.) are likely due to randomness or not.

↳ First you specify a null hypothesis

$\bar{x} = \mu$ (usually the thing
 $s = \sigma$ you think isn't true
 $B = 0$ ↳ trying to show
 isn't true)

$$\text{AVG(SAT)}_{\text{ECON}} = 610$$

→ The Hypothesis test tells you whether it's likely (you get to pick your defn of likely) you'd see this data if the null hypothesis is true.

Null Hypothesis: $H_0: \bar{X}_{\text{ECON}}^{\text{SAT}} = 610$

Alternative Hyp: $H_a: \bar{X}_{\text{ECON}}^{\text{SAT}} \neq 610$

num obs. = 20 16 t-distribution

Test Statistic: $\frac{\bar{X}_{\text{ECON}}^{\text{SAT}} - 610}{s/\sqrt{n}} \sim t_{15}$

$$= \frac{620 - 610}{\sqrt{15}/\sqrt{20}} = \frac{10}{\sqrt{15}/4} = \frac{10}{15/4} = \frac{40}{15}$$

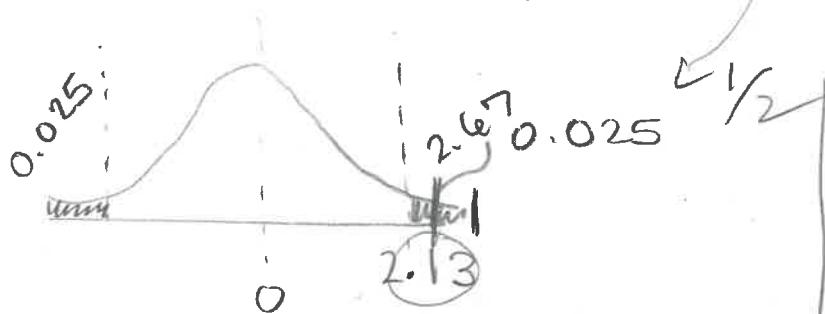
$$= \frac{620 - 610}{15/\sqrt{16}} = \frac{10}{15/4} = \frac{40}{15} = 2.67$$

Significance Level (How likely is likely).

$$\alpha = 0.05$$

$$df = 15 = n - 1$$

$$2.67 > 2.13$$



- Three Ways
1. Critical Value
 2. P-Values
 3. Confidence Intervals

Reject the Null

↳ There is a less than 5% chance that we would see this data AND the null be true

P μ = the true population mean
SAT Math score for SMC
students
= 610

s = sample std. dev \bar{x} = sample mean SAT math
score (sampled from some
Econ majors)
= 15
= 620

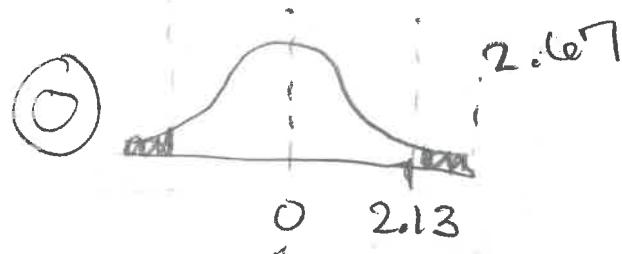
α = significance level
= 0.05

H $H_0: \bar{x} = \mu$ $H_a: \bar{x} \neq \mu$

$$\frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t_{15}$$

N Two-sided t-test at a 5% sig.
level

T test statistic = $\frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{620 - 610}{15/\sqrt{16}}$
= 2.67



M Reject the null hyp that the
sample mean SAT score equals the
popul. mean.

* Now let's use our data on book prices to do a difference of means test.

85

$$171 - 86 = 85$$

→ Question: Is there a statistically significant difference between the price of books by female authors and the price of books by male authors?

→ First: use our own individual

Open

→ Then: use our combined

SpreadSheet

(P) μ_1 = true mean price of books by female authors

in STATA

μ_2 = true mean price of books by male authors

\bar{x}_1 = sample mean price of books by female authors

\bar{x}_2 = sample mean price of books by male authors = 15.64

$$\begin{aligned} df &= (n_1 - 1) \\ &\quad + (n_2 - 2) \\ &= 168 \end{aligned}$$

$$\alpha = 0.05$$

$$s_1 = 5.20$$

$$n_1 = 84 \quad n_2 = 86$$

$$s_2 = 6.48$$

(H₀) $H_0: \mu_1 = \mu_2$ OR $H_0: \mu_1 - \mu_2 = 0$

$$\mu_1 - \mu_2 = 0$$

(A)
$$\frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{168}$$

$$\downarrow z$$

$$\begin{aligned} df &= (n_1 - 1) + (n_2 - 1) \\ &= 83 + 85 \\ &= 168 \end{aligned}$$

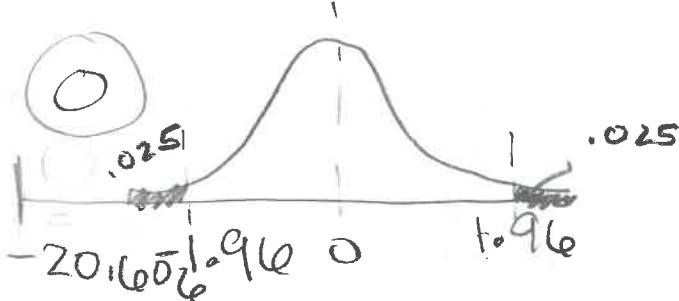
(indep)

N Two-Sided Difference of means t test

T test statistic: $\frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$

$$= \frac{\bar{x}_1 - \bar{x}_2}{S_p \sqrt{\frac{1}{84} + \frac{1}{86}}} = \frac{14.81 - 15.44}{S_p \sqrt{0.012 + 0.012}} = \frac{-0.83}{0.26 \sqrt{0.024}} = -0.83$$

$$S_p = \sqrt{\frac{s_1 + s_2}{(n_1 - 1) + (n_2 - 1)}} = \sqrt{\frac{5.20 + 6.48}{(83 + 85)}} = \sqrt{\frac{11.68}{168}} = \sqrt{0.04} = 0.21$$

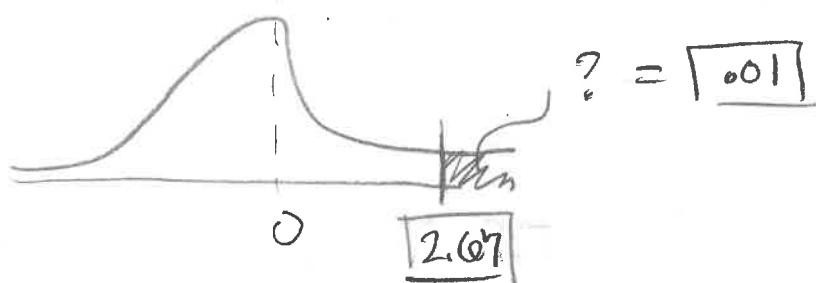


M Reject the null that $\mu_1 - \mu_2 = 0$

* Now do it again with the big data set
your dataset

(6)

* The p-value tells you exactly how likely it is.



→ There is a 1% chance that we would see this data if the null were true.

* Mnemonic Device for Hypothesis Testing

P : Parameters ($\mu, \sigma, \beta, \bar{x}, s, n, \alpha$)

H : Hypotheses

A : Assumptions

N : Name the test

T : Test Statistic

O : Outcome

M : Make Your Conclusion

→ Let's redo that test using the mnemonic.

→ Hypothesis Testing is one way to decide if our data says something statistically significant.

↳ Confidence intervals are another

Hypothesis Tests Answer: "Does this particular number constructed from this particular sample indicate statistical significance at the $\alpha\%$ significance level?"

Confidence Intervals Answer: "Suppose I have a sample of n observations, what range of numbers would indicate statistical significance at the $\alpha\%$ significance level."

(We almost always use 5%)

Ingredients for a Confidence Interval:

construct
the CI. Does my
null hyp't null
in it?

$$\bar{x} \pm t_{\alpha}(s/\sqrt{n})$$

↑ ↑
sample statistic critical value

- sample statistic
- critical value (based on α)
- std. error

• $(\mu_1 - \mu_2) = (\bar{x}_1 - \bar{x}_2) \pm t_{\alpha} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$

"I am $(1-\alpha)\%$ confident / There is a $(1-\alpha)\%$ chance that the true parameter lies in this interval."

- * Let's build a 95% Confidence Interval for the SAT Math Scores. (Together)

Question: What range of numbers can I be 95% sure contains the true mean SAT Math score for SMC students?"

$$\begin{aligned} \mu &= \bar{x} \pm t_{0.025} s/\sqrt{n} \\ &= 620 \pm (2.13) \left(\frac{15}{\sqrt{16}} \right) \\ &= 620 \pm (2.13) \left(\frac{15}{4} \right) \\ &= 620 \pm 7.99 \\ &= [612.01, 627.99] \end{aligned}$$

Is 610 in that interval? No.

↳ Reject the null at 5%..

- * Build a 95% CI for the difference in book prices using the big data set. (~~Independent Rgj~~)?

$$\mu_1 - \mu_2 = (\bar{x}_1 - \bar{x}_2) \pm t_{0.025} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$\begin{aligned}
 [] &= (14.81 - 15.64) \pm 1.96(0.04) \\
 &= -0.83 \pm 0.078 \\
 &= [-0.908, -0.752]
 \end{aligned}$$

* Does the confidence interval contain 0? No

→ Reject the null

* These are all ways to answer "Does this data tell us something or not?"

↳ The rest of the class is concerned with defining the "something."

* Build a 95% Confidence Interval for the true difference in prices

Math Tutoring Center

Fitting a Line to Data

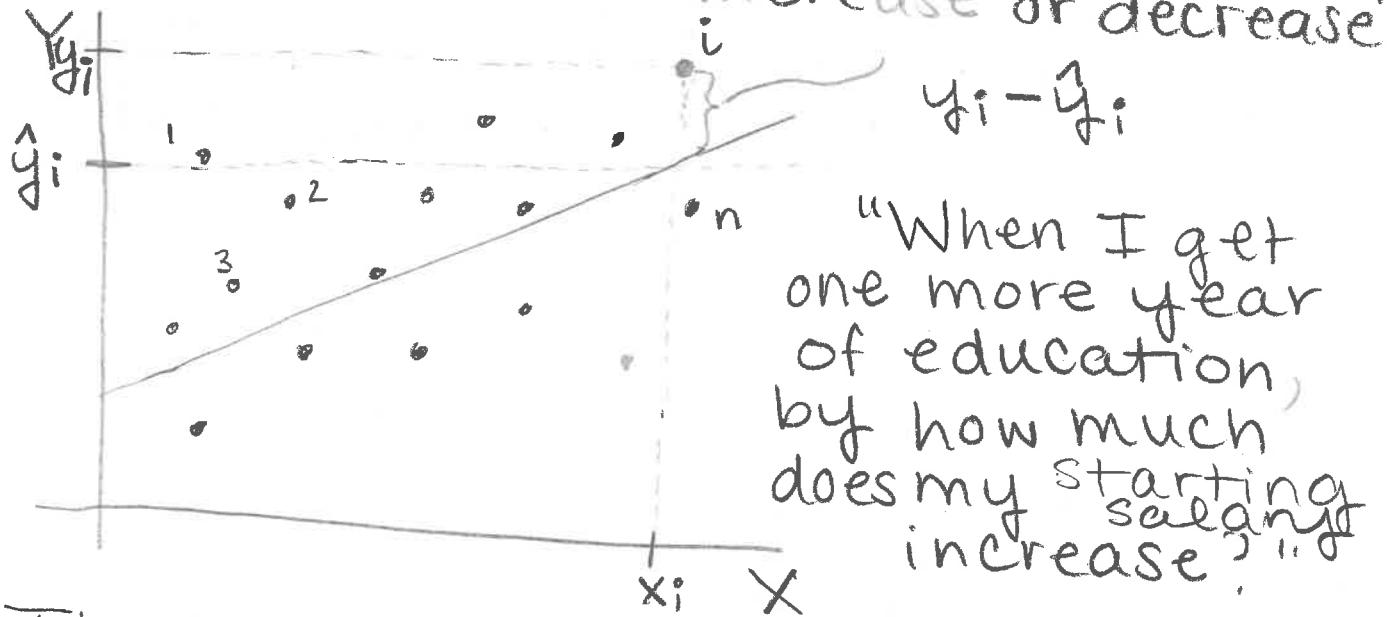
One Sided
Formula
Sheet

- * Deal with things early (^{phys.}_{mental})
- * Mostly focused on the statistics of a single variable.

↳ But most of our questions are about at least two variables.

→ First, plot your data.

- * "When one variable increases, by how much does the other increase or decrease?"



"When I get one more year of education, by how much does my starting salary increase?"

→ The most common way to condense this information in an informative way is to fit a line through it.

- * Why a line?

- easy interpretation ($\frac{\text{rise}}{\text{run}}$)
- not a lot to be gained by making it more complicated.

(*) Today: the mechanics of doing a good job fitting a line through sampled data

Next Class: using that line to say things about the population.

(*) You will need:

- an independent variable (x) (years of education)
- a dependent variable (y) (income)
- a calculator (optional)

→ There are lots of ways to pick a line (once there's more than 2 points)
 $y = a + bx$

→ Definitely want to minimize some concept of "distance" between the line & each point.

→ The one we usually use is the sum of the squared distance between the predicted y -value (\hat{y}_i) and the actual y -value (y_i) for each x -observation (x_i).

the y -val
on the
line

→ each dot is a observation
 $i = \{1, \dots, n\}$

④ We're looking for a line

$$\hat{Y}_i = a + bX_i$$

that minimizes the sum of the squared differences between the Y_i 's & \hat{Y}_i 's for each X_i .

That is

$$\min_{\{a, b\}} \sum_{x_i} (Y_i - \hat{Y}_i)^2$$

→ We're going to actually do this.
but let's rewrite something first.

→ there are going to be a lot of $Y - \bar{Y}$ & $X - \bar{X}$, so

$$y = Y - \bar{Y}$$

$$x = X - \bar{X}$$

$$\hat{Y}_i = a + b\tilde{x}_i + b\bar{x} - b\bar{x}$$

$$\begin{aligned} &= (a + b\bar{x}) + b\tilde{x}_i - b\bar{x} \\ &= \tilde{A} + b(\tilde{x}_i - \bar{x}) = \tilde{A} + bX_i \end{aligned}$$

doesn't depend on i

③

→ so now the task is

$$\min_{\{a, b\}} \sum_i (Y_i - \tilde{A} - bx_i)^2$$

* Calculus Reminder: To minimize or maximize something, you take the derivative & set it = 0,

→ Don't be intimidated by the \sum_i .
(You can always write it out)

$$S = (Y_1 - \tilde{A} - bx_1)^2 + (Y_2 - \tilde{A} - bx_2)^2 + \dots + (Y_n - \tilde{A} - bx_n)^2$$

To find chain rule!

$$\frac{\partial S}{\partial \tilde{A}} = \sum_i \cancel{2(Y_i - \tilde{A} - bx_i)}^1 (-1) = 0$$

$$\sum_i (Y_i - \tilde{A} - bx_i) = 0$$

$$\sum_i Y_i - \sum_i \tilde{A} - b \sum_i x_i = 0$$

(1) (2)

Show that
 $\sum_i (x_i - \bar{x}) = 0$

$$\textcircled{1} \quad \sum_{i=1}^n \tilde{A} = \underbrace{\tilde{A} + \tilde{A} + \dots + \tilde{A}}_n = n \tilde{A}$$

$$\textcircled{2} \quad b \sum_{i=1}^n x_i = b \left[\sum_{i=1}^n (x_i - \bar{x}) \right] = b [n\bar{x} - n\bar{x}] = 0$$

$$\sum_i Y_i - n\tilde{A} - 0 = 0$$

$$\sum_i Y_i = n\tilde{A}$$

$$\frac{\sum Y_i}{n} = \tilde{A}$$

$$\bar{Y} = \tilde{A} = a + b\bar{x}$$

$$\bar{Y} - b\bar{x} = a$$

need this
guy now.

To find b :

$$\frac{\partial S}{\partial b} : \sum_i 2(Y_i - \tilde{A} - bx_i)(+x_i) = 0 \quad (\text{div. by } 2)$$

$$\sum_i (x_i Y_i - x_i \tilde{A} - bx_i^2) = 0$$

$$\cancel{\sum_i x_i Y_i - \tilde{A} \sum_i x_i - b \sum_i x_i^2} = 0$$

$\cancel{=} 0$

$$\sum_i x_i Y_i - b \sum_i x_i^2 = 0$$

$$\sum_i x_i Y_i = b \sum_i x_i^2$$

$$\boxed{\frac{\sum x_i Y_i}{\sum x_i^2} = b}$$

no you can't
combine the \sum 's.

$$b = \frac{\sum x_i Y_i}{\sum x_i^2}$$

$$a = \bar{Y} - \left(\frac{\sum x_i Y_i}{\sum x_i^2} \right) \bar{x}$$

→ these are them! These are the linear regression coefficients.

- ↳ we have a lot to do to show that they're good for anything (they are).
- ↳ but for now we have them

(*) What the slope coefficient does and doesn't mean:

→ in an ideal world we would randomly assign X's (we would randomly assign levels of education & then measure income)

↳ but most of the time (especially in economics) that is not at all possible

- if by some miracle you have randomly assigned data

$b =$ the increase in Y caused by a 1 unit increase in X .

- if on the other hand if you do not have randomly assigned dream data

$b =$ the increase in Y associated with a one unit increase in X .

→ it includes the confounding factors associated with X (family background, grades, study habits, etc.).

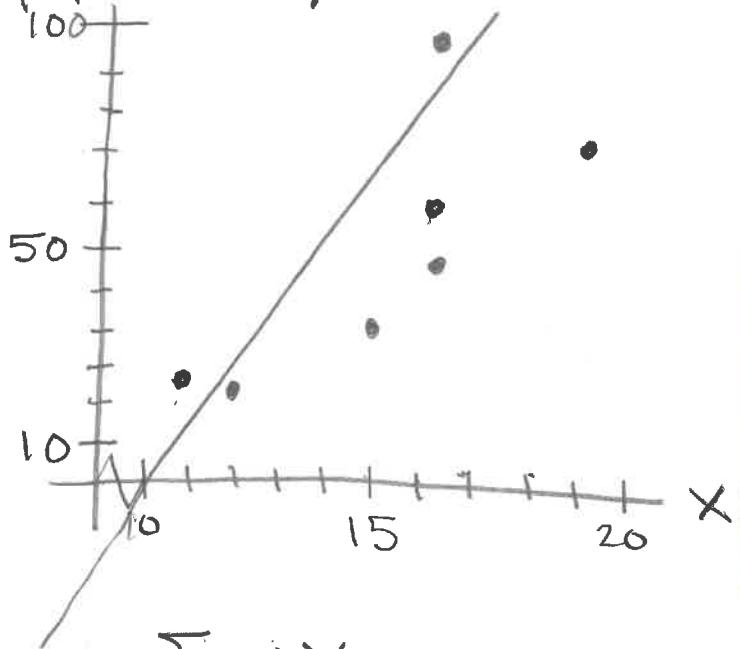
* Example: ⁽¹⁾ Calculate the regression line for the following data. ⁽²⁾ Interpret the slope coefficient. ⁽³⁾ What are some potentially confounding variables?

x_i (yrs. of ed.)	y_i (in thous.)	$x_i - \bar{x}$	x_i^2	$x_i y_i$
	(income)		$(x - \bar{x})^2$	$(x - \bar{x})Y$
11	24	-3.86	14.90	-100.36
12	25	-2.86	8.18	-71.5
15	32	0.14	0.02	4.48
16	65	1.14	1.30	74.1
14	45	1.14	1.30	51.3
16	100	1.14	1.30	114
18	75	3.14	9.86	235.5
$\Sigma =$	104	368	36.86	307.52

$$\bar{x} = 14.86$$

$$\bar{Y} = 52.57$$

① Y (income)



Minimizing the squared errors used to seem like magic. How does this eq do that ??

(Not magic. Calculus.)

$$b = \frac{\sum x_i Y_i}{\sum (x_i^2)} = \frac{307.52}{36.86} = 8.34$$

Always write the formula!

$$\begin{aligned} a = \bar{Y} - b \bar{X} &= 52.57 - (8.34) 14.84 \\ &= 52.57 - 123.93 \\ &= -71.36 \end{aligned}$$

$$\hat{Y} = -71.36 + 8.34 X$$

$$\begin{aligned} 0 &= -71.36 + 8.34 x \\ 71.36 &= 8.34 x \\ 8.56 &= x \end{aligned}$$

② There is a \$8,340 increase in income associated with one extra year of schooling.

③ • Whether you finish a degree
• School rank
• mother/father income

* Estimating Y's.

→ once you have an equation for Y in terms of X, you can use it to estimate Y for different X's.

Ex: Using the regression line we estimated, calculate your expected income for the number of years of education you have.

Mine is :

$$\hat{Y}_{\text{Prof.Lned+Re}} = -71.36 + 8.34(21) \\ = 103.78$$

The Simple Regression Model

→ Last time we focused on the mechanics of creating a regression line from sample data

↳ today we're going to write down a mathematical model of the economy that allows us to use that to make statistical statements about the real world.

* Assumptions

→ all models of the world come with assumptions (otherwise you'd just be like "sometimes stuff happens" & go get a coffee)

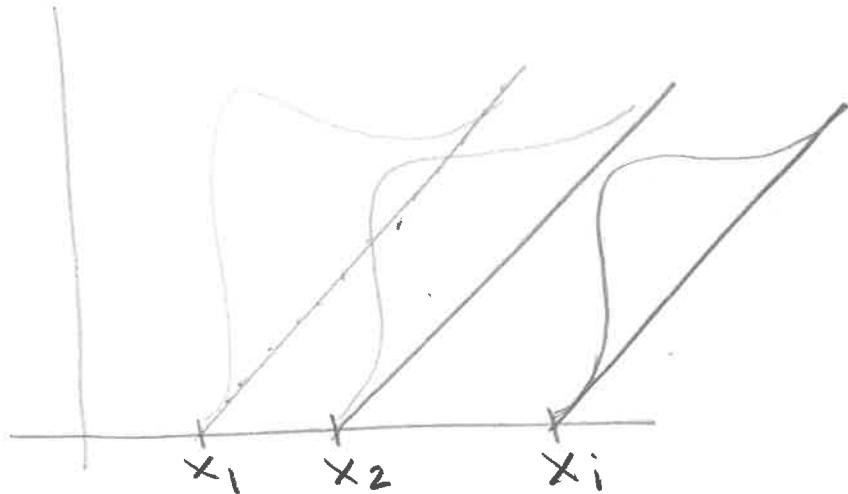
→ my advisor used to tell us: "All models are wrong. Some models are useful."

↳ Economists (or anyone doing mathematical modeling) always have to balance

- too many assumptions → not realistic
- too few assumptions → can't say anything useful

Simple Regression Model Assumptions

- * The distributions of $Y_i | x_i$ all have the same variance ^{event}



→ for every x , there's a distribution of Y_i 's possible (due to the randomness of life)

↳ these assumptions describe how we think about those distrib. & that randomness.

1. Homogeneous Variance:

→ all of those distributions have the same variance.

$$V[Y_1 | x_1] = V[Y_2 | x_2] = \dots = V[Y_n | x_n] = \sigma^2$$

$$V[Y_i | x_i] = \sigma^2 \quad \forall i$$

2. Linearity

→ the expected value/mean of each distribution is a linear function of X .

↳ it's, in fact, the same linear function of X .

$$E[Y_1 | X_1] = \underline{\alpha} + \underline{\beta} X_1$$

$$E[Y_2 | X_2] = \underline{\alpha} + \underline{\beta} X_2$$

(review
subscripts)

$$E[Y_n | X_n] = \underline{\alpha} + \underline{\beta} X_n$$

$$E[Y_i | X_i] = \underline{\alpha} + \underline{\beta} X_i \quad \forall i$$

3. Independence

→ Each Y_i is statistically independent from every other Y_j .

↳ Knowing Y_i gives you no information about Y_j .

→ these can all be stated in one sentence:

The random variables, Y_1, Y_2, \dots, Y_n , are independently distributed with mean $\alpha + \beta X_i$ and variance σ^2 .

→ most of the time we talk about the distribution of the error term rather than Y .

$$E[Y_i | x_i] = \alpha + \beta x_i \quad \text{then}$$

$$Y_i = \alpha + \beta x_i + u_i$$

$$E[u_i | x_i] = 0 \quad \text{show that if}$$

$$\sqrt{u_i | x_i} = \sigma^2$$

→ these are equivalent.

Population Regression Line:

$$E[Y_i | x_i] = \alpha + \beta x_i$$

the true
data fall
around the line -

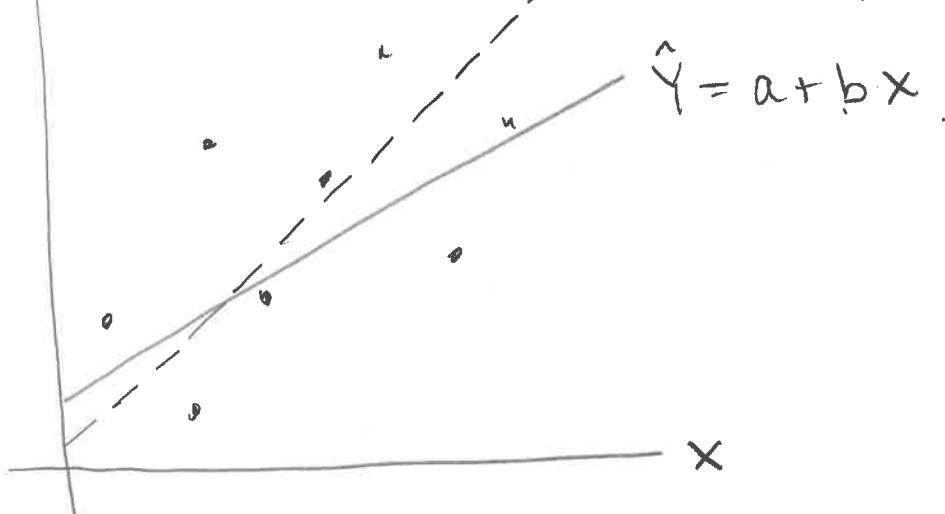
$$Y_i = \alpha + \beta x_i + u_i$$

this part is the line

Sample Regression Line:

$$\hat{Y} = a + b x \quad E[Y|x] =$$

$$\sim Y = \alpha + \beta x$$



} the assumpt.
are about
building
the
populatio
line

(*) Most of our interesting economic questions boil down to questions about β (not b).

→ is it > 0

→ is it < 0

→ is it just $\neq 0$ (there is some relationship)

Bad News: We really don't know what β is.

Good News: We definitely know what b is.

↳ We don't just know what it is, we know its distribution.

The slope estimate, b , is normally distributed with mean β and variance $\frac{\sigma^2}{\sum x^2}$.

→ this is it. This is why there is a rest of the course and we don't just pack up & go home.

* Estimating σ^2

→ since we almost never know σ^2

↳ we have to do what we always do: approximate the distribution using a sample estimate for s^2 and the t -distribution.

$$s^2 = \frac{1}{n-2} \sum_i (y_i - \hat{y}_i)^2$$

↑
because
we did a
and b

$$SE = \frac{s}{\sqrt{\sum x^2}}$$

→ so the distribution we actually use to do anything is

t_{n-2} with mean β and variance

$$\frac{s^2}{\sum x^2}.$$

→ now we can use all those statistical tools (hypothesis testing, confidence intervals, etc.)

if time

1. Perform a Hypothesis Test using the data from last class. Test the null hypothesis that $\beta = 0$.

2. Construct a 95% Confidence Interval for β .

For both of those it will take to construct s^2 & SE :

Y_i	\hat{Y}_i	$Y_i - \hat{Y}_i$	$(Y_i - \hat{Y}_i)^2$
26	20.38	5.62	31.58
25	28.72	-3.72	13.84
32	53.74	-21.74	472.6
65	62.08	2.92	8.53
45	62.08	-17.08	291.73
100	62.08	37.92	1437.93
75	78.76	-3.76	14.14
			2270.38

$$\hat{Y} = -71.36 + 8.34 X$$

$$\sum x^2 = 36.86$$

$$n = 7$$

$$n-2 = 5$$

$$S^2 = \frac{1}{5} (2270.38) \\ = 454.08$$

$$SE = \sqrt{\frac{21.31}{36.86}} = \frac{21.31}{6.07} = 3.51$$

* Next time, part of class will be for a review session.
→ your job is to bring questions

* Read Stata Tutorial

2.

Confidence Interval:

$$\begin{aligned}\beta &= b \pm t_{0.025}^{n-2} SE \\ &= 8.34 \pm (2.57)(3.51) \\ &= 8.34 \pm 9.02\end{aligned}$$

$$[-0.68, 17.36]$$

→ Since 0 is in the confidence interval, we fail to reject the null hypothesis that $\beta = 0$.

* You should be comfortable performing hypothesis tests using all methods (t-score, p-value, CI) for β .

1. Hypothesis Test

P

β = the slope coeff. of the true regression line.

b = the slope coeff. of the sample regression line

$$n = 7$$

$$\alpha = 0.05$$

H

$$H_0: \beta = 0 \quad H_a: \beta \neq 0$$

(8)

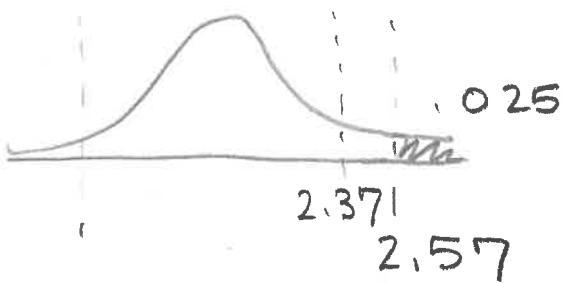
(A) $b \sim N(\beta, \frac{\sigma^2}{\sum x_i^2})$

↳ use t-distribution
 $df = n - 2 = 5$

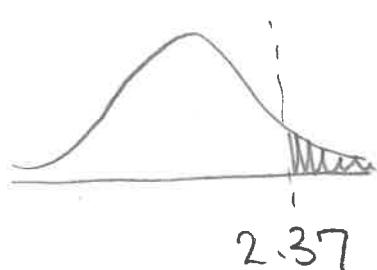
(N) Two-Sided t-test

(T) $t = \frac{b - 0}{SE} = \frac{8.34}{3.51} = 2.37$

t-score:



p-value



between 2.02 & 2.57

↳ $P \in (.05, .025)$

(M) Fail to Reject null
 that $\beta = 0$.

Goodness-of-Fit

* So far, we don't have any way of saying if the regression line is any good at explaining Y.

↳ We could regress height on favorite color and get a number.

↳ But how do we know whether your favorite color actually explains any of the variation in height? (It doesn't.)

→ to fix this, we're going to calculate R^2 .

* R^2 tells us the fraction of sample variation in Y explained by X.

→ No more.

→ No less.

(make sure you understand what things like R^2 do and do not say)

→ To calculate R^2 , then, we need some measures of variation.

1. Total Sum of Squares (SST)

$$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

(note: different from sample variance b/c didn't divide by $n-1$)

→ This measures the total sample variation in Y . (It's our denominator.)

→ how spread out is Y ?

If it's super spread out, all else eq., x will do a worse job, low R^2

→ there are two things that affect Y : x & ^{other} stuff

2. Explained Sum of Squares (SSE)

$$SSE = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

→ This is the total variation of the estimated \hat{Y} 's.

↳ These are the estimates of Y based entirely on x .

→ If these are also very spread out, we're doing ok. (2)

→ a lot of statistical programs will just tell you SSR.

↳ So it's useful to write R^2 in terms of SSR.

$$R^2 \stackrel{\text{also}}{=} 1 - \frac{SSR}{SST} .$$

We know this because

$$R^2 = \frac{SSE}{SST}$$

$$SST = SSE + SSR$$

$$SST - SSR = SSE$$

$$\begin{aligned} R^2 &= \frac{(SST - SSR)}{SST} = \frac{SST}{SST} - \frac{SSR}{SST} \\ &= 1 - \frac{SSR}{SST} . \end{aligned}$$

* R^2 measures exactly what we said: the fraction of total variation in Y explained by x.

→ Some very informative regressions have low R^2

→ Some useless regressions have very high R^2

3. Residual Sum of Squares (SSR)

$$SSR = \sum_{i=1}^n \hat{u}_i^2 = \left(\sum_{i=1}^n (\hat{u}_i - \bar{u}_i)^2 \right) \\ = \sum_{i=1}^n (\hat{u}_i - o)^2$$

→ This is the sample variation in the OLS residuals ($y_i - \hat{y}_i$)
 $= y_i - (a + b x_i)$

$$SST = SSE + SSR$$

(I encourage you to prove this.)

④ Now we can define R^2

$$R^2 = \frac{SSE}{SST}$$

variation in Y expl. by X
total var. in Y

→ sometimes it's way more of a pain in the buns to calculate

$$SSE = \sum_{i=1}^n (\hat{y}_i - \bar{y}_i)^2$$

↳ you have to calculate all the \hat{y}_i 's & subtract, & square.

→ doesn't mean it isn't a useful measure.

↳ It measures a very specific thing perfectly.

Example: Calculate the R^2 for the regression we did of income on years of schooling.

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

→ so let's be strategic about this: Which formula should we use?

↳ either way, we need SST (sadly we didn't calculate the variance of Y yet)

→ But last time, we calculated $(Y_i - \hat{Y}_i)^2$! That's \hat{u}_i^2 !

↳ so we have $SSR = 2270.38$

Y_i	$Y_i - \bar{Y}$	$(Y_i - \bar{Y})^2$	$\bar{Y} = 52.57$
26	-26.57	705.94	
25	-27.57	760.10	
32	-20.57	423.12	
65	12.43	154.50	
45	-7.57	57.30	
100	47.43	2,249.60	
75	22.43	503.10	
		4853.68	

good way to
check.
 $SST > SSR$
 SSE

$$R^2 = \frac{1 - \frac{SSR}{SST}}{\text{don't forget the } \underline{\underline{1}}}$$

$$= 1 - \frac{2270.38}{4853.68}$$

$$= 1 - 0.47 = \underline{\underline{0.53}}$$

* 53% of the variation in income
is explained by years of schooling.

Stata

→ depends on computer situation

cps78-85.dta

reg lwage educ
↑
log of
wage

① Write down the regression line

$$\hat{lwage} = 0.98 + 0.069 \text{ educ}$$

~~$$lwage = 0.98 + 0.069 \text{ educ}$$~~

② Is the coeff. on educ stat. sig.?

3 Ways:

1. Critical Value $12.08 > 1.96$

2. P-Value .000

3. Confidence Interval

$[0.058, 0.080]$

③ R^2 ? 0.12

12%

Multiple Regression

Simple Regression Eq: $Y = \alpha + \beta X + u$

→ probably the sketchiest assumption we made in the Simple Regression Model was that everything contained in u is uncorrelated with X .

$$E[u|x] = 0.$$

- * If we're looking at a regression of wage on years of educ.

$$\text{Wage} = \alpha + \beta \text{educ} + u$$

→ It's not a problem that there are other things that affect wage

↳ we may just not be interested in them

→ What is a problem is if those things are correlated with education

- ↓ What types of things might be in that u term?

→ things that affect wage that aren't education?

$u \longrightarrow \{$ experience
gender
family income
race
marital status
:
 $\}$

→ Do you find it believable that any of those things aren't correlated with years of education?

↳ I can tell you there almost no married women of color from poor backgrounds in PhD programs.

(*) One of the assumptions of the Simple Regression Model fails.

↳ need a different model.

Multiple Regression Model

→ Just like the SRM but instead of just the one X , there are several X 's.

New Terminology/Conventions:

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + u$$

↑ → there are k explanatory variables. (x_1, x_2, \dots, x_k) ↑ same
same

→ α becomes β_0 , just so we don't go crazy w/ our greek letters

$\beta_0, \beta_1, \beta_2, \dots, \beta_k$ looks nice.

→ there are now k coefficients:

$\beta_1, \beta_2, \dots, \beta_k$

↳ (k+1) total parameters:

$\beta_0 + (\beta_1, \beta_2, \dots, \beta_k)$

→ Each of these k explanatory vars has n observations (for each observation person)

↳ Helpful to think of it like a matrix or an spreadsheet.

	educ	exper.	gender	marital status	
person 1	10	4	m	...	x_k
person 2	•	•	•	...	S
⋮	⋮	⋮	⋮	⋮	⋮
person n	•	•	•	...	•

n × k matrix

(*) You need to be very comfortable with the regression equation & the structure of the data

↳ What's k? What's n? What does a given number in the matrix tell you?

(*) So how do we actually calculate the estimates for those β 's?

↳ What we started out interested in is β_1 .

↳ We may or may not be interested in β 's $2 \rightarrow K$.

↳ Need to include them to get (*) the right β_1 .

(*) We're still going to use Ordinary Least Squares (OLS)

↳ minimize the sum of squared errors.

$$\min \left\{ \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right\}$$

what makes up
 \hat{y}_i is different.

just the same

$$\min_{\beta_0, \beta_1, \dots, \beta_k} \left\{ \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_1 - \hat{\beta}_2 x_2 - \dots - \hat{\beta}_k x_k)^2 \right\}$$

→ the process is the same:

You don't need to actually do all this to tell me how you would. We'd take the derivatives and set them to 0

do algebra until our eyes bleed.

→ But now we have $k+1$ equations.

↳ that sounds like a job for a computer!

↳ Stata!

Command: (email link?)

```
use http://fmwww.bc.edu/ec-p/data/
woolridge/wage1
```

regress wage educ exper female
nonwhite married
(use same order as me)

① What is the regression equation?

$\hat{wage} = -1.78 + 0.58 \text{educ} + 0.06 \text{exper}$
 $- 2.07 \text{female} - 0.03 \text{nonwhite}$
 $+ 0.16 \text{married}$

* Interpreting the Coefficients

→ The β estimates ($\hat{\beta}$'s) give us the *ceteris paribus* effects of each x ,
 (other things) (equal)

$\Delta x_i = \text{the change in } x_i = x_i^a - x_i^b$

$$\Delta \hat{Y} = \hat{Y}^a - \hat{Y}^b$$

$$= (\hat{\beta}_0 + \hat{\beta}_1 x_1^a + \hat{\beta}_2 x_2^a + \dots + \hat{\beta}_k x_k^a)$$

$$- (\hat{\beta}_0 + \hat{\beta}_1 x_1^b + \hat{\beta}_2 x_2^b + \dots + \hat{\beta}_k x_k^b)$$

$$= \hat{\beta}_1 (x_1^a - x_1^b) + \hat{\beta}_2 (x_2^a - x_2^b) + \dots + \hat{\beta}_k (x_k^a - x_k^b)$$

$$= \hat{\beta}_1 \Delta x_1 + \hat{\beta}_2 \Delta x_2 + \dots + \hat{\beta}_k \Delta x_k$$

"What is the change in \hat{Y} associated with a one unit change in x_1 , holding other factors fixed?"

→ Other factors fixed?

one unit
change in
 x_1

$$\Delta x_2 = 0$$

$$\Delta x_3 = 0$$

$$\Delta x_1 = 1$$

→ plug those in.

$$\Delta \hat{Y} = \hat{\beta}_1(1) + \hat{\beta}_2(0) + \dots + \hat{\beta}_k(0)$$
$$= \hat{\beta}_1$$

→ The calculus version is much easier.

$$\frac{\partial \hat{Y}}{\partial x_1} = \frac{\partial (\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_k x_k)}{\partial x_1}$$
$$= \hat{\beta}_1$$

* The point of MR is that even when we don't have "other factors fixed" data, we can still make "other factors fixed" conclusions

Stata

* What is the change in wage associated with one more year of education?

$$\hat{\beta}_1 = 0.58 \quad \text{\$.58 an hour}$$

④ Run a simple regression of wage on educ.

↳ What is the effect here?

$$\hat{\beta} = 0.54$$

⑤ Which estimate do you believe more, why?

$\hat{\beta}$, because all those other factors are left in u in the simple regression

↳ assumptions fail

⑥ You need to be comfortable using generic variables and specific, contextual variables.

Properties of Multiple Regression

β 's are true parameters, b 's are

* What makes a good estimator? estimators

- unbiased $E[b_i] = \beta_i$

↳ kind of a bare minimum

- simple straightforward function of the y 's.

- low variance

↳ not too spread out

→ Let's discuss the assumptions of the model & see how the estimator do.

Assumption #1: Y is a linear function of the β 's.

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + u$$

(Linear here means the power on each β_i is 1, as opposed to $\beta_i^2 X_i$)

Assumption #2: The data

x_1, x_2, \dots, x_k for each i
are the result of a random
sample.

Assumption #3: There are no
perfect linear relationships between
any combinations of x 's.

↳ no perfect collinearity.

Assumption #4: $E[u | x_1, x_2, \dots, x_k] = 0$
↑ true error

* Keep these last two straight.

#3: what relationships can &
can't exist between the x 's &
each other.

#4: what relationships can &
can't exist between the x 's and u .

* Theorem: Under the assumptions #1-4, the OLS estimators are unbiased. That is,

$$E[b_0] = \beta_0$$

$$E[b_1] = \beta_1$$

$$\vdots$$

$$E[b_k] = \beta_k$$

* What happens if the x 's you use to estimate Y are different from the "true" determinants of Y ?

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

~~~~~  
true relationship

Case 1: You estimate

$$\hat{Y} = \hat{b}_0 + \hat{b}_1 x_1 + \hat{b}_2 x_2 + \hat{b}_3 x_3 + \hat{u}$$

(you include too many variables)

Case 2: You estimate

$$\tilde{Y} = \tilde{b}_0 + \tilde{b}_1 x_1 + \tilde{u}$$

(you include too few)

## Case 1: Irrelevant Variables

We can write the true model as:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$$

but  $\beta_3 = 0$

→ all the  $\beta_i$ 's are variable specifying numbers about the real world.

→ we could write down an equation for  $Y$  that included all the irrelevant  $X$ 's & their  $\beta$ 's that = 0 but we don't because it would be silly (& impossible)

Let's check our assumptions.

#1: Linear in  $\beta$ 's? ✓

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$$

#2: Random Sample? ✓

→ including  $X_3$  has no effect on this

#3: No perfect collinearity? ✓

→ no effect

#4:  $E[u | X_1, X_2, X_3] = 0$ ? ✓

→ Nothing about including  $X_3$  necessitates a violation of these assumptions.

↳ the theorem holds!

↳ all the  $b$ 's are unbiased!

Specifically  $E[\hat{b}_3] = R_3 = 0$

## Case 2: Omitted Variables

The real equation is

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + u$$

but you estimate

$$\tilde{Y} = \tilde{\beta}_0 + \tilde{\beta}_1 X_1 + \tilde{u}$$

→ now  $\tilde{u} = \beta_2 X_2 + u$

→ If  $X_2$  is not correlated with  $X_1$ , then we're fine

$$\tilde{Y} = \tilde{\beta}_0 + \tilde{\beta}_1 X_1 + \tilde{u}$$

satisfies the simple regression model assumptions & we can go about our business,

→ But we noted last time, that's usually not true (think education and family background).

$$E[x_2|x_1] \neq 0.$$

\* How does this violate the MR assumptions?

$$E[\tilde{u}|x_1] = E[\beta_2 x_2 + u|x_1]$$

$$= E[\beta_2 x_2|x_1] + E[u|x_1]$$

$$\underbrace{\hspace{1cm}}_{\neq 0}$$

→ So we can't conclude  $\hat{\beta}_1$  is unbiased.

\* We need one more assumption (about the variance of  $u$ ) to conclude these estimators for  $\beta_i$ 's are the best ones to use.

Assumption #5: The error term,  $u$ , has the same variance for any values of the  $x$ 's.

↳ homoskedasticity

→ So whether  $\text{educ} = 10$  or  $18$ ,  $u$  has the same variance

$$V[u|18, f, \text{mar.}] = \sigma^2$$

$$V[u|10, m, \text{single}] = \sigma^2$$

### Gauss-Markov Theorem:

Under Assumptions #1-5, the estimators  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ , computed using OLS, for  $\beta_0, \beta_1, \dots, \beta_k$  are the Best Linear Unbiased Estimators.

Best: Lowest Variance

Linear: Can be written as a linear function of  $y$ 's

Unbiased:  $E[\text{estimators}] = \text{true values}$

\* What the thm says: if we compute estimates for  $\beta_i$ 's (call them  $\hat{\beta}_i$ ) using OLS, they will be "linear" estimator that are unbiased. Among "linear" & unbiased estimators for  $\beta_i$ 's, these have the smallest variance.

\* What it's good for: the method for computing our estimates for the  $\beta$  coefficients is not only practical (it's "linear"), and reasonable (unbiased), but it's the least noisy (lowest variance) one we might reasonably use.

---

### Avoiding Symbol Overload

\* A lot of moving parts:

- variables ( $Y$ )
- data ( $Y_i$ )
- parameters ( $\beta$ )
- estimators ( $b_i$ )

\* The "True" / Theoretical Model  
(variables, parameters)

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

no hats.  
↳ hats  
mean  
estimators

(variables, estimators)  
(\*) The equation we estimate

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \hat{u}$$

hats mean  
estimate  
or estimator  
(squiggles too)

(\*) The estimation method we use:  
(estimators, data)

OLS

→ produces specific estimators  
using data → an actual  
numbers

$$\hat{\beta} = \frac{\sum x_i y_0}{\sum x_i^2} \quad (\text{simple reg. version})$$

→ Write down the wage, educ, etc.  
model out in this way.

- the true model
- the equation we estimated
- the actual estimates of the  $\hat{\beta}_i$ 's.

## Multiple Regression: Inference

→ "Inference" just refers to formal hypothesis testing

↳ what can we actually infer from what we have?

\* Typically questions about  $\beta$ 's

- Is  $\beta_i = 0$ ? (today)

- Are  $\beta_1 = \beta_2 = \beta_3 = 0$  together?  
(next time)

→ Need probability distribution of  $\beta$ 's to do this

↳ the distribution of the  $\beta$ 's depends on the distribution of  $Y|X$  ("Linear")

↳ which depends on the distribution of  $u$

Assumption #6: The population error,  $u$ , is independent of the explanatory var  $x_1, \dots, x_k$ , and is distributed  $N(0, \sigma_u^2)$ .

Theorem: Under assumptions #1-6, conditional on  $x$ 's,

$$b_j \sim N(\beta_j, \sigma_{\beta_j}^2)$$

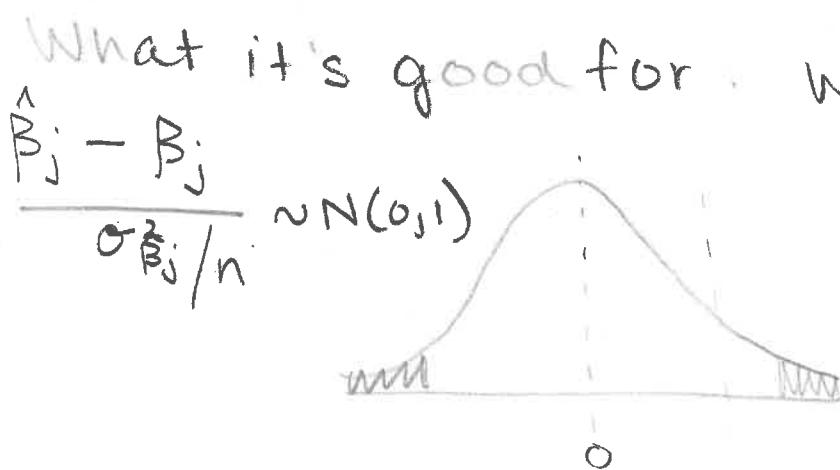
→ the variance of  $\hat{\beta}_j$  is a complicated formula (described in the book) almost always use a  $\hat{\sigma}_{\beta_j}^2$

$$\sigma_{\beta_j}^2 = \frac{\sigma_u^2}{SST_j(1-R_j^2)}$$

total variance in  $x_j$        $R^2$  of a regression of  $x_j$  on all the other  $x$ 's.

↳ so all of that is insane which is why we ask Stata to do it.

What the theorem says: under the assumptions we've outlined, the  $\hat{\beta}$ 's we estimate have this specific distribution.



We can use our <sup>normal</sup> <sub>distr</sub> hyp. testing tools to say things about the  $\beta$ 's!

(2)

→ just one more step to take this from a math problem to an implementable process.

\* As (almost) always: we don't actually know  $\sigma_{\hat{\beta}_j}^2$ .

↳ We don't know  $\sigma_u^2$

↳ we assumed it into existence (it's probably a reasonable assumption that  $u \sim N(0, \sigma_u^2)$  but we certainly don't know what  $\sigma_u^2$  is)

$$\sqrt{[u]} = E[(u - E[u])^2]$$

definition of var.

$$= E[(u - \bar{o})^2] = E[u^2]$$

→ closer! but we don't observe  $u$  either

$$\frac{1}{n} \sum_{i=1}^n \hat{u}_i^2$$

$n - (k + 1)$   
↑ num parameters  
num obs to estimate

$$\hat{\sigma}_{\hat{\beta}_j}^2 = \frac{\frac{1}{n-(k+1)} \sum_{i=1}^n \hat{u}_i^2}{SST_j (1 - R_j^2)}$$

→ So if you were on a desert island and someone asked you for an estimator for the variance of  $\hat{\beta}_j$ , that's what you would do.

→ For now: Stata.

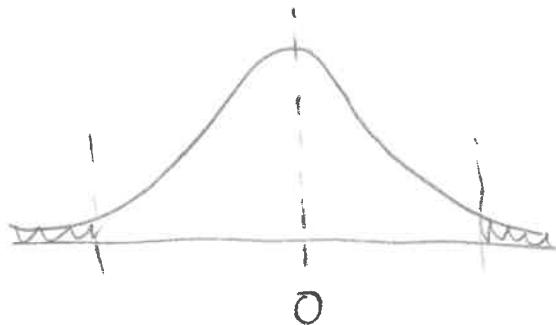
→ But now we can actually do some hypothesis testing

↳ all the ingredients for a t-test.

\* We've come a long way from our  $\bar{x}$ , std. dev. hyp. test days

→ the picture & the process are the same.

$$\frac{\hat{b}_j - \beta_j}{\sqrt{\hat{\sigma}_{\hat{b}_j}^2}} \sim t_{n-(k+1)}$$



(\*) Today: test one  $\beta$  at a time

↳ next time: relationships between  $\beta$ 's.

(\*) Hypothesis Testing of a single  $\beta$ :

Question: Does variable  $x_i$  have a statistically significant effect on the dependent variable,  $Y$ ?

Test: Is  $\beta_i$  (statistically sign.)  $\neq 0$ ?

→ can do one-sided tests but usually just do two-sided

3 Different Ways:

(Typically)  
use  $\alpha = 0.05$

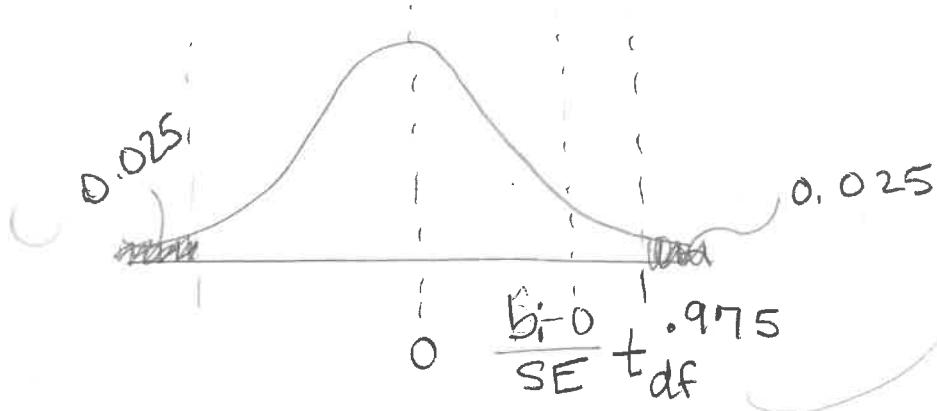
1. Critical Value

$$\frac{\beta_i - 0}{SE(\hat{\beta}_i)} \quad H_0: \beta_i = 0 \quad H_a: \beta_i \neq 0$$

2. P-Values

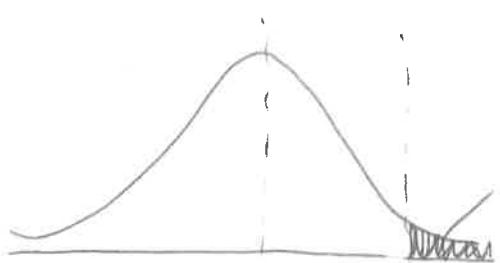
3. Confidence Intervals

1. Critical Value

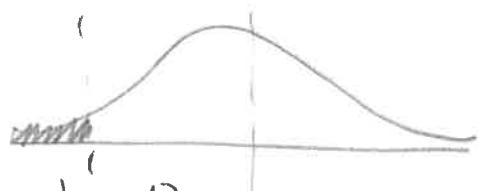


look up  
the t<sub>df</sub>  
critical  
value  
if  $\frac{\beta_i - 0}{SE} < t$  fail to  
reject  
 $> t$  reject (5)

## 2. P-Values



p-value



$$-\frac{b_i - 0}{SE(b_i)}$$

Calculate the  
P-value

↳ area in  
the tail  
(whichever tail)

if p-value > 0.02

↳ fail to  
reject

if  
< .025  
reject

"there is a p% chance we  
would see this data if the  
null ( $\beta_i = 0$ ) is true."

## 3. Confidence Interval

→ calculate the 95% confidence  
interval

↳ "I am 95% confident that the  
true value of  $\beta_i$  is contained  
in this interval."

$$[b_i^{LB}, b_i^{UB}]$$

if the null value ( $\beta_i = 0$ ) in [ ] → fail to  
reject  
not in [ ] → reject

(\*) Keeping all this null, reject, fail to reject stuff straight.

- The null hyp is that  $\hat{\beta}_i = 0$ .
    - ↳ that var.  $x_i$  is not a statistically sign. determinant of  $Y$ .
  - If we fail to reject the null.
    - ↳ can't conclude  $\hat{\beta}_i \neq 0$ .
    - ↳ can't conclude  $x_i$  a sign. det. of  $Y$
  - If we reject the null
    - ↳ reject  $\hat{\beta}_i = 0$
    - ↳ can conclude  $x_i$  a sign. det. of  $Y$
- Stata!
- load wage data  
→ run regression from last time

(\*) Is education a statistically significant determinant of wage?

$$H_0: \hat{\beta}_1 = 0$$

$$H_a: \hat{\beta}_1 \neq 0$$

## 1. Critical Value

df?

$$df = n - (k+1) = 526 - (6) = 520$$

$$t_{.520}^{.975} = 1.96$$

$$\text{test value} = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} = \frac{0.58}{0.052} = \underline{\underline{11.15}}$$

test value  $> t_{520}^{.975}$  → reject the null  
(difference from rounding)

## 2. P-Value

$$P\text{-value}_1 = 0.000$$

P-value  $< 0.025$  → reject the null

## 3. Confidence Interval

$$[0.48, 0.68]$$

null value (0) not in interval  
↳ reject the null

Conclusion: Reject the null  $\beta_1 = 0$ ,  
conclude education is a statist.  
significant determ. of wage.

\* Are there any variables that we included that aren't sign. at the (two-sided) 5% (2.5%) level?

non-white  
married

\* Are there any that are close?  
married  
the constant

→ The rule I use is: if you have to mess with the  $\alpha$  or the test to get significance, you probably can't conclude it's significant.

## Mult. Reg.: Inference (Cont'd.)

- Last time: individual coefficients statistically signif.
- This time: relationships between and among coefficients.

Question: Which plays a bigger determ. of your wage: level of education or level of experience?

| Variable | $\hat{\beta}$ |
|----------|---------------|
| const    | -1.78         |
| educ     | 0.58          |
| exper.   | 0.056         |
| female   | -2.07         |
| nonwhite | -0.025        |
| married  | 0.066         |

$\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \hat{\beta}_4, \hat{\beta}_5$

→ The coeff. on education is larger than experience

↳ but just like we can't conclude married is sign. even though the coeff.  $> 0$ , we can't conclude this w/o formally testing it.

\* What are we asking?

"Is  $\hat{\beta}_1 \neq \hat{\beta}_2$ , statistically significant"

~~don't assume the direction~~

$$H_0: \hat{\beta}_1 = \hat{\beta}_2 \rightarrow \hat{\beta}_1 - \hat{\beta}_2 = 0$$

$$H_a: \hat{\beta}_1 \neq \hat{\beta}_2 \rightarrow \hat{\beta}_1 - \hat{\beta}_2 \neq 0$$

↑                                   ↑  
know                                   know  
the distr.                           the distr.

→ Need distribution of the random variable  $\hat{\beta}_1 - \hat{\beta}_2$ .

$$\hat{\beta}_1 - \hat{\beta}_2 \sim N(\beta_1 - \beta_2, ?)$$

$$\sqrt{[\hat{\beta}_1 - \hat{\beta}_2]} ?$$

$$= \sqrt{[\hat{\beta}_1]} + \sqrt{[\hat{\beta}_2]} - \text{Cov}[\hat{\beta}_1, \hat{\beta}_2]$$

\_\_\_\_\_       \_\_\_\_\_       \_\_\_\_\_  
got           got this      a huge mess  
this

Options:

- calculate it analytically (calculus)
- do some algebraic gymnastics w/ the regression (in the book)
- \* • ask the nice computer

→ Cue up our regression from last time

\* the command for these super complicated fancy tests is: test.

↳ uses whatever regression you ran last.

↳ name the  $\beta$ 's using their associated variables.

Command: test educ  $(=)$  exper

$\begin{matrix} > \\ = \\ = \end{matrix}$

→ it lists the null:

$$\text{educ} - \text{exper} = 0$$

We'll talk about the F part in a minute.

→ gives the p-value (& the test statistic)

$$\text{Prob} > F = 0.000$$

→ Since p-value less than 0.025 reject null!

Conclusion:  $\hat{\beta}_1$  is statistically significantly greater than  $\hat{\beta}_2$   
↳ Years of education is a stronger determinant of your wage than years of experience.

\* You can do even crazier equations than  $\hat{\beta}_1 - \hat{\beta}_2 = 0$ .

↳ If you wanted to know "Is  $2\hat{\beta}_1 + 3\hat{\beta}_2 + \hat{\beta}_4 = \hat{\beta}_3$ ?"

test  $2*\text{educ} + 3*\text{exper} + \text{nonwhite}$   
= = female

---

\* We know how to test if  $\hat{\beta}_1$  is signif. and we know how to test if  $\hat{\beta}_2$  is signif. but how do we test if they're jointly significant?

Question: If we control for education & experience, does your race, gender, and marital status play a (statist. sign.) role in determining your wage?

Or: Are  $\hat{\beta}_3$ ,  $\hat{\beta}_4$ , and  $\hat{\beta}_5$  jointly equal to zero?

\* How is this diff. from testing each one at a time?

→ When we test each individually  
We're testing  $\hat{\beta}_3 = 0$  or  $\hat{\beta}_4 = 0$   
or  $\hat{\beta}_5 = 0$ .

→ This is testing  $\hat{\beta}_3 = 0$  and  $\hat{\beta}_4 = 0$   
and  $\hat{\beta}_5 = 0$

\* We're asking if this estimated equation (Unrestricted Model)

$$Y = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X_3 + \hat{\beta}_4 X_4 + \hat{\beta}_5 X_5 + \hat{u}$$

is statistically different from this estimated equation

$$Y = \tilde{\beta}_0 + \tilde{\beta}_1 X_1 + \tilde{\beta}_2 X_2 + \tilde{u} \quad (\text{Restricted Model})$$

$$(\tilde{\beta}_3 + \tilde{\beta}_4 + \tilde{\beta}_5 = 0)$$

where  $\tilde{\beta}_3$ ,  $\tilde{\beta}_4$  and  $\tilde{\beta}_5$  are restricted to 0

The null hypothesis is

$$H_0: \hat{\beta}_3 = 0 \text{ and } \hat{\beta}_4 = 0 \text{ and } \hat{\beta}_5 = 0$$
$$\hat{\beta}_3 = \hat{\beta}_4 = \hat{\beta}_5 = 0$$

$H_a: H_0$  is not true

At least one of  $(\hat{\beta}_3, \hat{\beta}_4, \hat{\beta}_5)$  are not  $= 0$ .

\* The Restricted Model is the one where we impose the restrictions from the null hypothesis.

\* The Unrestricted Model has no such restrictions.

\* Each variable set to 0 is a restriction  
 $q \equiv$  the number of restrictions

= 3 in this case  
(also called numerator df)

→ We can actually run this test using the  $R^2$ 's from the two model estimations.

$$\frac{(R_{UR}^2 - R_R^2) / q}{(1 - R_{UR}^2) / (n - (k+1))} \sim F_{q, n-(k+1)}$$

④ Using Stata:

regress wage ... married

Unrestricted  $R^2 = 0.3158$

$$n - (k+1) = 520$$

$$F_{3, 520} = 2.60$$

regress wage educ exper

Restricted  $R^2 = 0.2252$

$$q = 3$$

$$\begin{aligned} F \text{ statistic} &= \frac{(0.3158 - 0.2252) / 3}{(1 - 0.3158) / 520} \\ &= \frac{0.0906 / 3}{0.6842 / 520} \\ &= \frac{0.0906}{3} \cdot \frac{520}{0.6842} = \frac{47.112}{2.0526} \\ &= 22.95 \end{aligned}$$

$$22.95 > 2.60$$

reject null

⑦

\* Can also use "test"

Command: test (female = 0)

(nonwhite = 0) (married = 0)

→ don't forget that it uses the most recently estimated model

↳ reestimate LR model

P-value = 0.000

Reject null

Conclusion:  $\hat{\beta}_3$ ,  $\hat{\beta}_4$ , and  $\hat{\beta}_5$  are not all jointly 0.

↳ education and experience aren't the only thing that determine wage.

## Functional Form

\* We're pretty much in doing linear regression variables.

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$$

→ the linear part refers to the  $\beta$ 's.

→ the  $x$ 's are just doing pretty much anything as they satisfy the

$$Y = \beta_0 + \beta_1 (7w_1) + \beta_2 (1 - x_1)$$

→ from a practical perspective just mult.  $w_1$  by 7 in  $x_1$  and reg  $Y | x_1$  as

\* There are pretty much things we can do to a variable

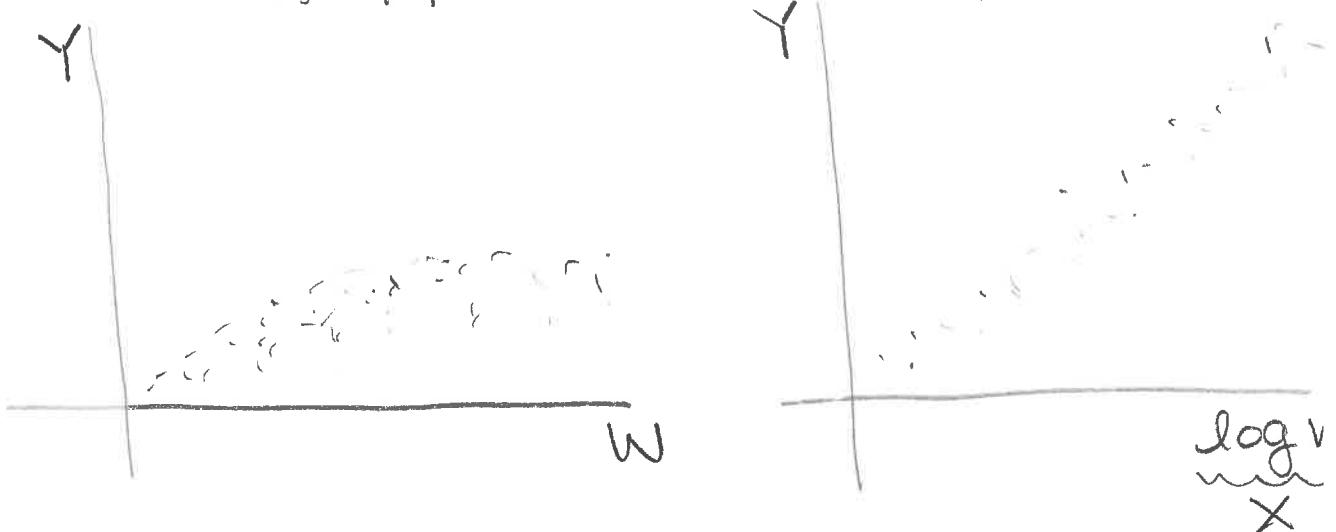
1. Scale it:

→ add a constant

→ mult by a constant

→ this is pretty cosmetic  
↳ do it to make things easier to work with (fewer 0's, more standard interpretation, etc.)

2. Take the natural log
  - $\log(x)$
  - We do this when  $y$  doesn't appear to be linearly related to  $X$ , appears to be exponential



3. Raise it to a power
  - usually squared:  $x^2$
  - can be other ones:  $x^3, \sqrt{x}$ , etc
  - square variables when we want to make big numbers bigger
    - age, experience

17 vs. 21

2 vs. 6

→ first one: doesn't change much,  
straightforward

→ 2nd & 3rd: a little trickier

## \* Scaling

$Y = \text{college GPA}$

$W_1 = \text{SAT score}$

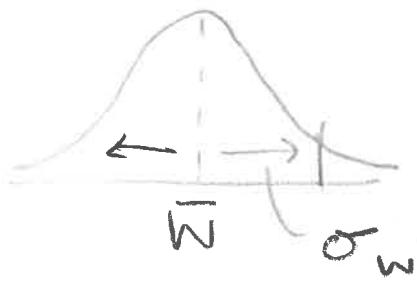
$W_2 = \text{ACT score}$

$$Y = \beta_0 + \beta_1 W_1 + \beta_2 W_2 + u$$

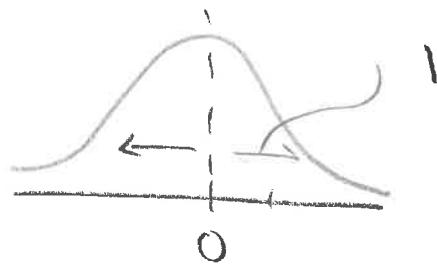
→ standardize everything!

Standardize: subtract the mean,  
divide by std. deviation.

↳ construct a z-score



Standardize  
====>



$$\frac{Y - \bar{Y}}{\sigma_Y}, \quad \frac{W_0 - \bar{W}_0}{\sigma_0}, \quad \frac{W_1 - \bar{W}_1}{\sigma_1}, \quad \frac{W_2 - \bar{W}_2}{\sigma_2}$$
$$z_y, \quad z_0, \quad z_1, \quad z_2$$

$$\frac{u - \bar{u}}{\sigma_u}, \quad \frac{u - 0}{\sigma}$$

③

→ take our equation

$$Y = \beta_0 + \beta_1 W_1 + \beta_2 W_2 + u$$

→ and average it

$$\sum_{i=1}^n Y = \sum_{i=1}^n (\beta_0 + \underbrace{\dots}_{\beta_1} + u)$$

$$\bar{Y} = \beta_0 \cdot 1 + \beta_1 \bar{W}_1 + \beta_2 \bar{W}_2 + \underbrace{u}_{=0}$$

→ subtract  $\bar{Y}$  from  $Y$

\* always remember your Algebra I.

→ got to do the same thing to both sides of the equation.

$$Y - \bar{Y} = \cancel{\beta_0 \cdot 1} + \beta_1 W_1 - \beta_1 \bar{W}_1 + \beta_2 W_2 - \beta_2 \bar{W}_2 + u - \cancel{u}$$

$$\frac{Y - \bar{Y}}{\sigma_Y} = \underbrace{\beta_1 (W_1 - \bar{W}_1)}_{\sigma_1} + \underbrace{\beta_2 (W_2 - \bar{W}_2)}_{\sigma_2} + u$$

→ classic algebra trick:

$$\text{mult. by } 1 = \frac{\sigma_1}{\sigma_1} = \frac{\sigma_2}{\sigma_2}$$

$$\frac{Y - \bar{Y}}{\sigma_Y} = \underbrace{\beta_1 (W_1 - \bar{W}_1) \frac{\sigma_1}{\sigma_1}}_{\sigma_Y} + \underbrace{\beta_2 (W_2 - \bar{W}_2) \frac{\sigma_2}{\sigma_2}}_{\sigma_Y} + l$$

$$\frac{Y - \bar{Y}}{\sigma_Y} = \beta_1 \left( \frac{W_1 - \bar{W}_1}{\sigma_1} \right) \frac{\sigma_1}{\sigma_Y} + \beta_2 \left( \frac{W_2 - \bar{W}_2}{\sigma_2} \right) \frac{\sigma_2}{\sigma_Y} + \frac{u}{\sigma_Y}$$

$$z_y = \beta_1 \frac{\sigma_1}{\sigma_y} z_1 + \beta_2 \frac{\sigma_2}{\sigma_y} z_2 + \frac{u}{\sigma_y}$$

$\downarrow$        $\downarrow$   
 $b_1$        $b_2$

$$\hat{Y} = b_1 x_1 + b_2 x_2 + v$$

Before:  $\beta_1$  told us the effect of hard SAT score on GPA level

↳ one more SAT point

Now:  $b_1$  tells us the effect of 1 std. deviation impr. in SAT score on your GPA relative to everyone else

↳ one SAT std. dev from the mean associated with a  $b_1$  std. dev. impr. in GPA.

-----  
Stata!

load gpa.dta

hmm. no ACT

↳ but high school size

↳ use that for  $W_2$  like we meant to all along!

1. Create the Standardized Vars.

| Var    | Mean   | Std. Dev |
|--------|--------|----------|
| cumgpa | 2.08   | 0.9896   |
| sat    | 898.91 | 168.19   |
| hssize | 312.77 | 198.78   |

(sum)

"summarize cumgpa sat hssize"

- ④ to create a new variable use "generate" (gen)

$$\text{gen z-gpa} = (\text{cumgpa} - 2.08) / 0.9896$$

$$\text{gen z-sat} = (\text{sat} - 898.91) / 168.19$$

$$\text{gen z-hs} = (\text{hssize} - 312.77) / 198.78$$

→ now we have three new vars we can just use as usual.

↳ although we should specify no constant

reg z-gpa z-sat z-hs, noconstant  
one word

$$\hat{z\text{-gpa}} = 0.1521 z\text{-sat} + -0.00005 z\text{-hs}$$

- ④ A 1 std. dev. in SAT score is assoc. with a 0.15 std. dev. increase in college GPA.

## \* Log Transformations

→ When we use log transformation to make the variables work better

↳ sometimes they make the data better satisfy the G-M Assumptions

→ They can have really nice Economic interpretations.

Y - depend.

X - ind.

→ If we regress

$$\log(Y) = \beta_0 + \underline{\beta_1} \log(x) + u$$

$\beta_1$  is the elasticity of Y

with respect to X

Elasticity: the  $\% \Delta$  of one variable associated with a  $1\% \Delta$  in another

Ex: the price elasticity of demand is the % $\Delta$  in quantity demanded associated with a 1% $\Delta$  in price.

$$E_D = \frac{\% \Delta Q}{\% \Delta P}$$

Ex: We have data

$Y$  = number of houses sold in US counties

$X$  = avg. annual salary in the counties

regress

$$\log(\hat{Y}) = \hat{\beta}_0 + \hat{\beta}_1 \log(X)$$

$\hat{\beta}_1$ , estimator for the income elasticity of demand

→ A 1% increase in income is associated with a  $\hat{\beta}_1$ % increase in houses sold.

\* If we also included a second independent variable, like county population that isn't logged (b/c it's already linearly related with  $Y$ )

$$\hat{\log(Y)} = \hat{\beta}_0 + \hat{\beta}_1 \log(x_1) + \hat{\beta}_2 x_2$$

→  $\beta_2$  is not an elasticity

↳ it's called a semi-elasticity

"A 1 unit increase in  $x_2$  is associated with a  $(100\beta_2)$  % increase/decrease in  $Y$ ."

"Y then X"

| Model       | Dep. Var  | Ind. Var  | Interpret.                                    |
|-------------|-----------|-----------|-----------------------------------------------|
| level level | $Y$       | $X$       | $\Delta Y = \beta_1 \Delta X$                 |
| level log   | $Y$       | $\log(X)$ | $\Delta Y = (\frac{\beta_1}{100})\% \Delta X$ |
| log level   | $\log(Y)$ | $X$       | $\% \Delta Y = 100 \beta_1 \Delta X$          |
| log log     | $\log(Y)$ | $\log(X)$ | $\% \Delta Y = \beta_1 \% \Delta X$           |

What we  
do

→ Every Metrics textbook should have a version of this table to refer to.

↳ yours is on p. 39

✳ to predict  $Y$  when you regressed  $\log(Y)$ , adjust  $e^{\log Y}$  by

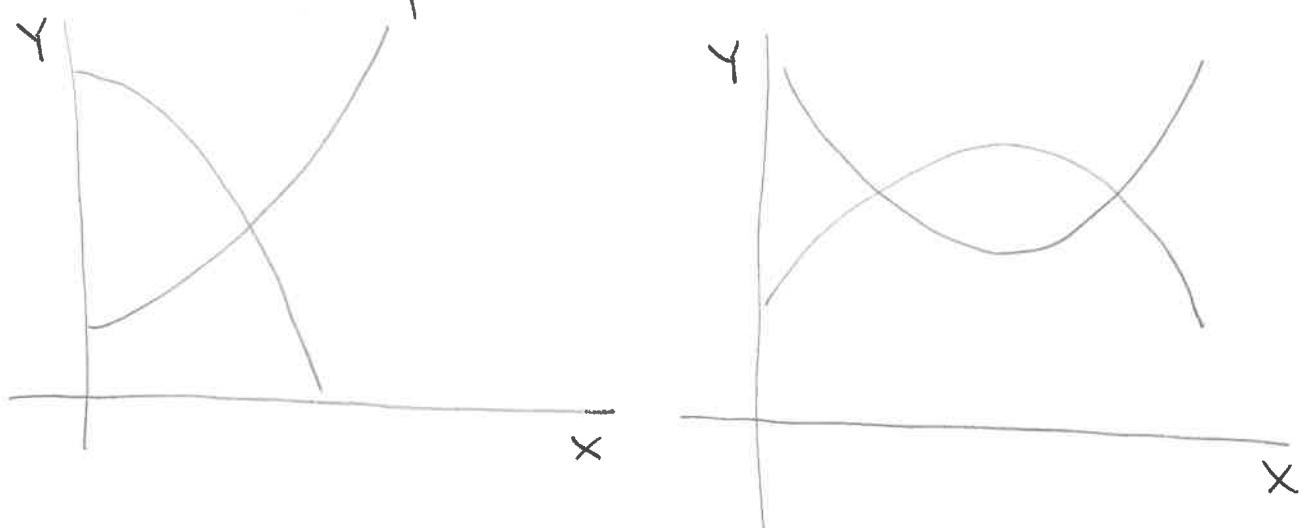
$$a = \frac{1}{n} \sum_{i=1}^n e^{\hat{u}} = \bar{e}^{\hat{u}} \rightarrow \hat{Y} = a \cdot e^{\log Y}$$

## Functional Forms: Quadratics & Interactions

- last time: first two things we can do to variables (scale them, take the log)
- today: the third thing
  - ↳ multiply variables together
  - ↳ sometimes it's the same variable:  $x^2, x^3$ , etc.
  - ↳ sometimes it's different variables:  $x_1 \cdot x_2$

### \* Quadratics:

- When the variable appears non-linear, quadratic



$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \dots + u$$

①

→ squaring things exaggerates differences

- big numbers get bigger
- small numbers get smaller
- that's why we square the errors in OLS

→ usually the squared variable is also there as a level

- not a problem to implement just call  $x_2 = x_1^2$
- doesn't violate perfect collinearity
  - BUT: makes the interpretation of the effect of  $x$  on  $Y$  a little trickier
  - If  $\Delta x_1 = 1$ ,  $\Delta Y = \beta_1$  or  $= \beta_2$  or  $= \beta_1 + \beta_2$ , etc.

→ I'm afraid we'll need a smidgen of calculus today

Calculus [Reminder?]

$$f(x) = \sqrt{x^0} + 4\sqrt{x^2} + \sqrt{x^3}$$

$$\frac{df(x)}{dx} = \frac{d(x^0)}{dx} + \frac{d(4x^2)}{dx} + \frac{d(x^3)}{dx}$$

$$= 1x^0 + 4(2x^1) 3x^2$$

$$= 1 + 8x + 3x^2$$

partial der.  $\rightarrow$  just treat other vars as cons.

$$\begin{aligned}\frac{\partial Y}{\partial x_1} &= \frac{\partial (\beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \dots)}{\partial x_1} \\ &= \frac{\partial \beta_0}{\partial x_1} + \frac{\partial \beta_1 x_1}{\partial x_1} + \frac{\partial \beta_2 x_1^2}{\partial x_1} + \dots \\ &= 0 + \beta_1 \cdot 1 + \beta_2 (2x_1) + 0 \\ &= \underbrace{\beta_1 + 2\beta_2 x_1}_{\text{"Partial Effect"}} \leftarrow \text{"Partial Effect"}\end{aligned}$$

\* the effect of a change in the variable  $x_1$  depends on the particular value of  $x_1$ .

Example:  $\hat{Y} = 3 + 2x_1 + 4x_1^2$

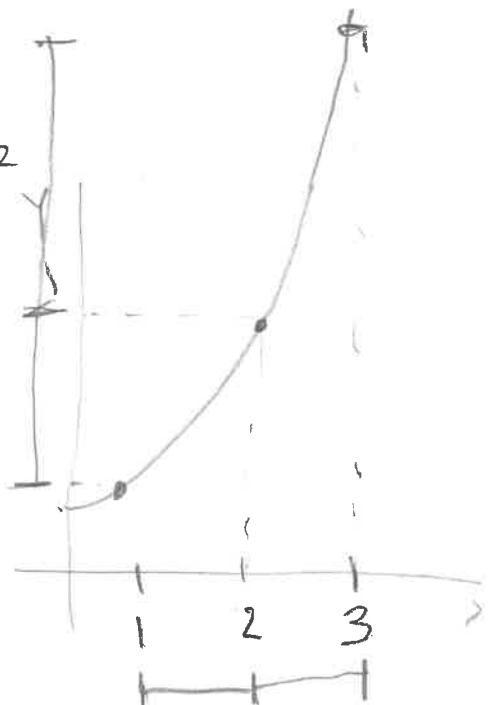
$$\frac{\partial \hat{Y}}{\partial x_1} = 2 + 8x_1$$

When  $x_1 = 1$

$$\Delta Y = 2 + 8(1) = 10$$

When  $x_1 = 2$

$$\Delta Y = 2 + 8(2) = 18$$



\* What about when  $\beta_1 > 0$  and  $\beta_2 < 0$ ?  
 ↳ What the heck does that mean?

Example: Wage =  $3.73 + \underline{.298} \exp - \underline{.0061} \exp^2$

Effect of  $x$  on  $\hat{Y}$ ?

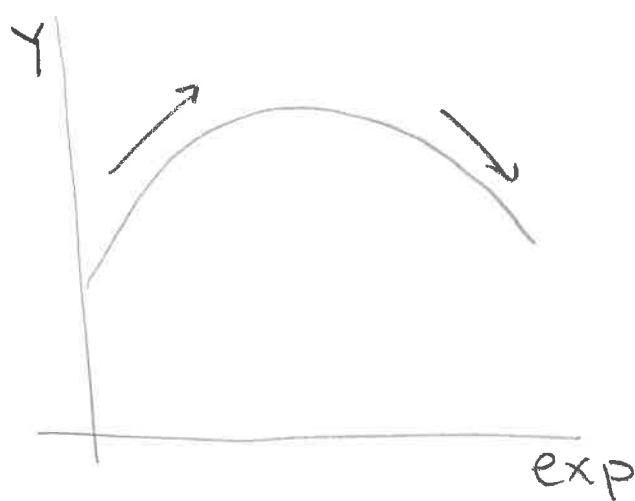
$$\frac{\partial \hat{Y}}{\partial \exp} = 0.298 - 0.0122 \exp$$

CHANGE  
in  $\hat{Y}$

→ diminishing returns

make \*  
you keep what  
you're graphing  
straight!

→ the change to  $\hat{Y}$  is positive for small values of  $\exp$  but eventually it's negative



(economically relevant)

→ A natural question is: what's the turning point?

↳ thankfully our high school algebra classes gave us a formula for this!

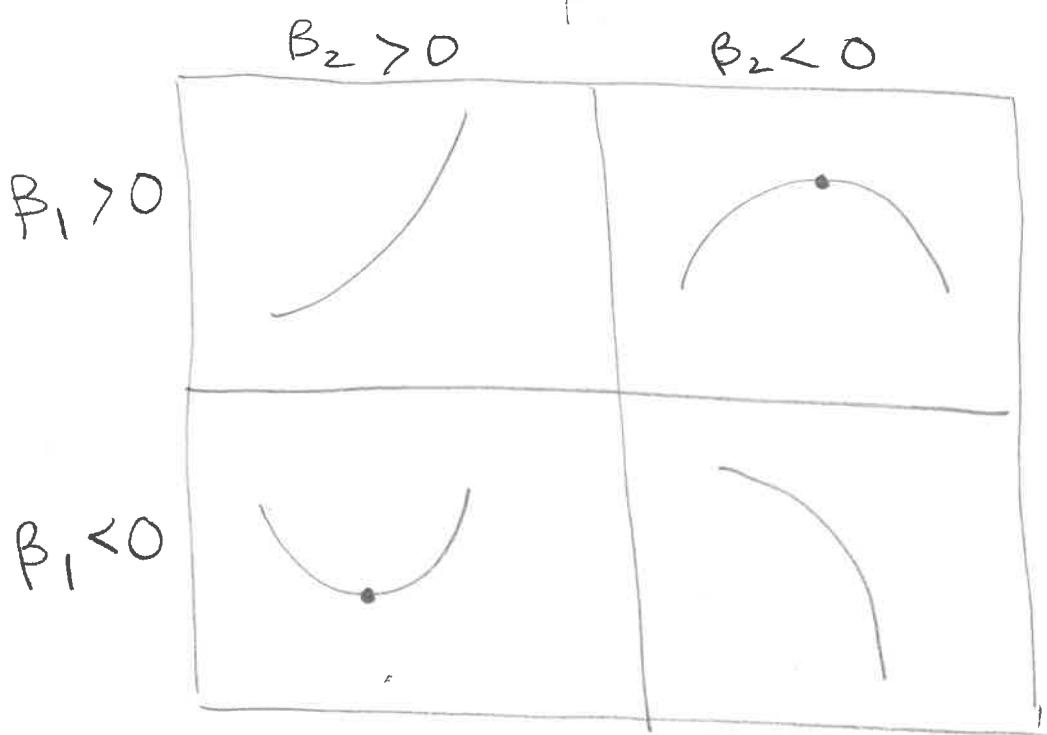
$$x^* = \left| \frac{\hat{\beta}_1}{2\hat{\beta}_2} \right|$$

↑  
turning  
point  $x^*$

the turning point tells us when it stops being good & starts being bad (or vice versa)

Turning Point Experience Level =

$$\left| \frac{0.298}{-0.0122} \right| = |-24.4| = \underline{\underline{24.4}}$$



## \* Interaction Terms:

- Sometimes the effect of a variable depends on the value of another variable
- If you've ever watched Flip or Flop (or really any HGTV show) you know that the number of bedrooms has diminishing returns ( $\times^2$ ) but also depends on the square footage.

$$\text{price} = \hat{\beta}_0 + \hat{\beta}_1 \text{sqft} + \hat{\beta}_2 \text{rooms} + \hat{\beta}_3 \text{rooms}^2 \\ + \hat{\beta}_4 \text{sqft} \cdot \text{rooms}$$

- Now we need to be a little more specific about partial derivatives.
- When you have a function with multiple variables:

$$f(x, w, z) = x^2 + zx + 2wz + 4w^2x$$

and you just want the derivative (change with respect to one variable)

$$\frac{\partial f(x, w, z)}{\partial w} \quad \text{"the partial derivative of } f \text{ with respect to } w$$

- hard to say, easy to do

→ just treat all the variables that aren't that variable like constants

$$\frac{\partial f(x, w, z)}{\partial w} = 0 + 0 + 2(1)z + 4(2w)x \\ = 2z + 8wx$$

1. What is the effect on price associated with sqft?

$$\frac{\partial \hat{p}_{rice}}{\partial \text{sqft}} = 0 + \hat{\beta}_1 + 0 + 0 + \hat{\beta}_4 \text{rooms} \\ = \hat{\beta}_1 + \hat{\beta}_4 \text{rooms}$$

→ the effect of square feet depends on how many bedrooms there are.

2. What is the effect on price associated with rooms?

$$\frac{\partial \hat{p}_{rice}}{\partial \text{rooms}} = 0 + 0 + \hat{\beta}_2 + 2\hat{\beta}_3 \text{rooms} \\ + \hat{\beta}_4 \text{sqft} \\ = \hat{\beta}_2 + 2\hat{\beta}_3 \text{rooms} + \hat{\beta}_4 \text{sqft}.$$

→ the effect of the number of bedrooms depends on both the number of rooms and the square feet.

→ it's all well and good to have an equation describing the effect associated with a variable

↳ But we would really prefer to have a single number.

↳ "dollars per sq ft"

## (\*) Average Partial Effect

$$\text{Partial Effect (Rooms)} = \hat{\beta}_2 + 2\hat{\beta}_3 \text{rooms} + \hat{\beta}_4 \text{sqft}$$

→ average across all values

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \text{PE(rms)} &= \frac{1}{n} \sum_{i=1}^n (\hat{\beta}_2 + 2\hat{\beta}_3 \text{rooms}_i + \hat{\beta}_4 \text{sqft}_i) \\ &= \hat{\beta}_2 + 2\hat{\beta}_3 \frac{1}{n} \sum_{i=1}^n \text{rooms}_i + \hat{\beta}_4 \frac{1}{n} \sum_{i=1}^n \text{sqft}_i \\ &= \hat{\beta}_2 + 2\hat{\beta}_3 \bar{\text{rooms}} + \hat{\beta}_4 \bar{\text{sqft}} \end{aligned}$$

→ definition: APE =  $\bar{\text{PE}}$

computed by plugging in  $\bar{x}$ 's.

→ "on average, an increase of one bedroom is associated with an increased (APE) dollars."

# Stata!

- load hprice data
- sum price area rooms

| Var   | Mean      | Sta. Dev. |
|-------|-----------|-----------|
| Price | 946100.66 | 43223.73  |
| area  | 2106.73   | 694.96    |
| rooms | 6.59      | .9012     |



- generate rooms-sq = rooms \* rooms
- generate rooms-area = rooms \* area
- regress price area rooms rooms-sq  
rooms-area

$$\begin{aligned} \text{APE(rooms)} &= \hat{\beta}_2 + \hat{\beta}_3 \overline{\text{rooms}} + \hat{\beta}_4 \overline{\text{area}} \\ &= 54,013.65 + (-4209.14)(6.5 \\ &\quad + 4.62(2106.73)) \\ &= 54013.65 - 27,738.23 \\ &\quad + 9,825.49 \\ &= 36,100.91 \end{aligned}$$

"On average, one more bedroom is assoc. with \$36,100.91 in price increase."

# Binary Variables

\* So far: quantitative data  
(measure a quantity)

- mostly continuous:
  - wage
  - sq. feet
  - police per-capita
- some discrete
  - education
  - bedrooms
  - num. children

\* A lot of interesting & important questions involve qualitative data

binary { - gender

- marital status

categorical - region of the US (N, S, E, W)

→ this <sup>section</sup> ~~chapter~~

use this type of information correctly.

## (\*) Binary Variables

→ take on the value 0 or 1

→ what's a 0 and what's a 1  
is up to you.

### Words of Caution:

#### 1. Be consistent

- with what we do in class
- across your work

#### 2. Name your variable according to what's = 1

- female = 1 if female  
0 if male  
(not "gender")
- married = 1 if married  
0 if not married
- always tell people  
what's 1 and what's 0.

#### 3. Never include both (for example) male and female

- perfect collinearity.

| female |   |
|--------|---|
| wave   | 0 |
| male   | 1 |
| educ   | 0 |
| ...    | 0 |
| wave   | 0 |

Example:

Model:

$$\text{Wage} = \beta_0 + \delta_1 \text{male} + \beta_1 \text{educ} + u$$

$$\begin{aligned} \text{male} &= 1 \quad \text{if male} \\ &= 0 \quad \text{if female} \end{aligned}$$

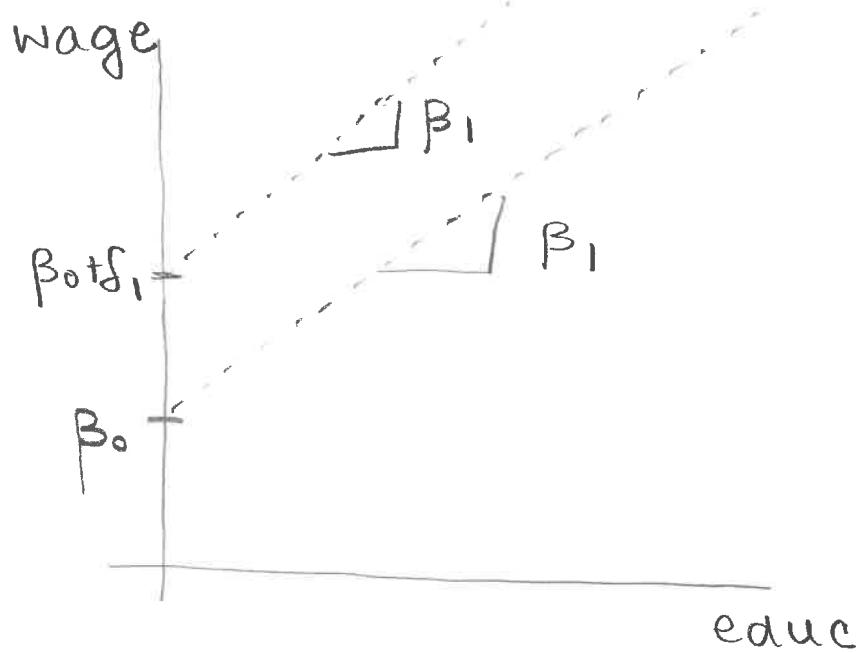
$$E[\text{wage} | \text{male}, \text{educ}] = \beta_0 + \delta_1 \text{male} + \beta_1 \text{educ}$$

↑  
Plug in  
values

$$\begin{aligned} E[\text{wage} | 1, \text{educ}] &= \beta_0 + \delta_1(1) + \beta_1 \text{educ} \\ &= \underline{\beta_0 + \delta_1} + \beta_1 \text{educ} \end{aligned}$$

$$\begin{aligned} E[\text{wage} | 0, \text{educ}] &= \beta_0 + \delta_1(0) + \beta_1 \text{educ} \\ &= \underline{\beta_0} + \beta_1 \text{educ} \end{aligned}$$

→ graph them!



assumes

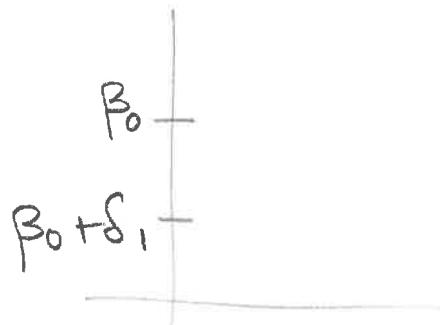
→ returns to educ don't depend on gender

→  $\delta_1$  is the difference in wage due to being male ( $\delta_1 > 0$ )

→ make sure you keep your definition straight

↳ if we had used "female"

$$\delta_1 < 0$$



→ the variable that is a 0 is called the base group or benchmark group.

↳ what you're making comparison to.

### \* Categorical & Ordinal Variables

→ it's straight forward to extend the binary variable process to cover multiple categories.

Example: Regions of the Country

$$\text{region} = \begin{cases} 1 = \text{North} \\ 2 = \text{South} \\ 3 = \text{East} \\ 4 = \text{West} \end{cases}$$

→ comes in two flavors:  
categorical & ordinal

ordinal: the order matters

$$\text{credit rating} = \begin{cases} 1 = D \\ 2 = C \\ 3 = B \\ 4 = A \end{cases}$$

→ categorical: order doesn't matter

can switch

$$\text{region} = \begin{cases} 1 = N \\ 2 = E \\ 3 = S \\ 4 = W \end{cases}$$

→ ordinal: can say something like "higher is better"

↳ direction means something.

④ To actually implement this it's best to basically make a binary variable for every (-1) category.

→ Could make credit rating one variable, but that would assume constant partial effects: going from D → C is the same as B → A

→ making region its own variable makes no sense.

↳ because order doesn't matter.

$$\text{north} = \begin{cases} 1 & \text{north} \\ 0 & \text{other} \end{cases} \quad \text{east} = \begin{cases} 1 & \text{east} \\ 0 & \text{other} \end{cases}$$

$$\text{south} = \begin{cases} 1 & \text{south} \\ 0 & \text{other} \end{cases}$$

④ Always leave out one category  
    ↳ perfect collinearity

$$CR-D = \begin{cases} 1 & \text{if D} \\ 0 & \text{if other} \end{cases} \quad CR-C = \begin{cases} 1 & \text{if C} \\ 0 & \text{if other} \end{cases}$$

$$CR-B = \begin{cases} 1 & \text{if B} \\ 0 & \text{if other} \end{cases}$$

(leave out A)

→ the one you leave out is the base case

→ Note: sometimes it's not practical to create separate dummies

Example: school rank

→ can just treat as qualitative

→ break into percentiles, etc.

— — — — — — — — — —

Stata!

- load data from spreadsheet
- regress  
    price male - author

## \* Interpret Coefficient

- make genre dummies

$$\text{sci-fi}_i = \begin{cases} 1 & \text{scifi} \\ 0 & \text{other} \end{cases}$$

⋮  
⋮  
⋮

regress price male - author

$$\text{sci-fi} \quad - \quad - \quad - \quad -$$

## \* Interpret coefficients

# Heteroskedasticity

\* The assumption that took our OLS estimates from unbiased to BLUE was homoskedasticity

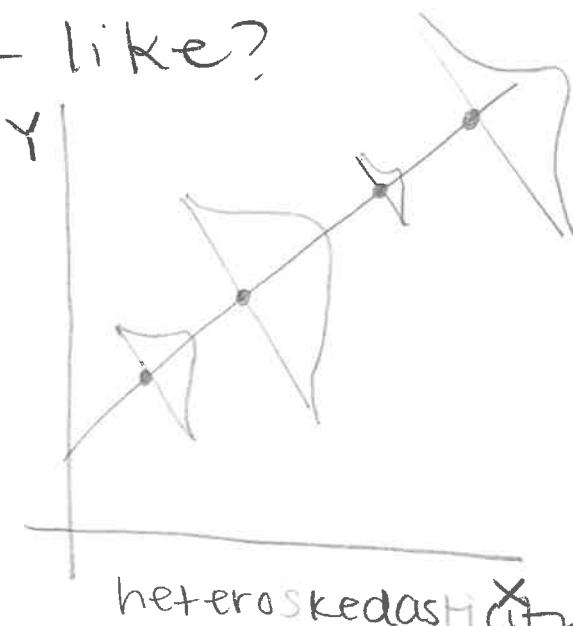
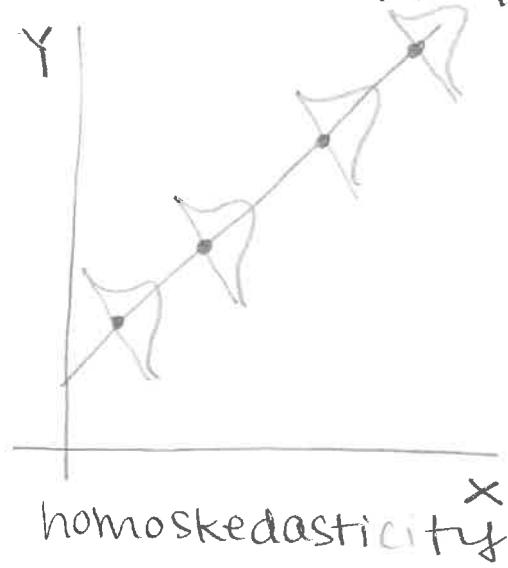
$$V[u | x = x_i] = \sigma^2 \quad \forall x_i$$

→ the variance of  $u$  doesn't change with values of  $x$ .

→ what happens when that assumption doesn't hold?

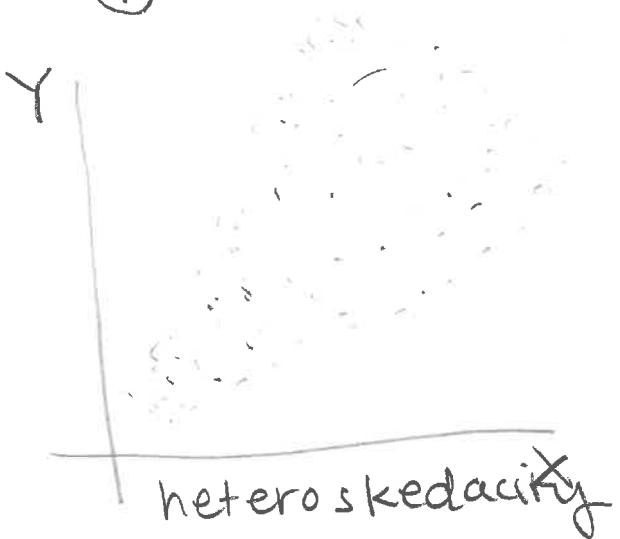
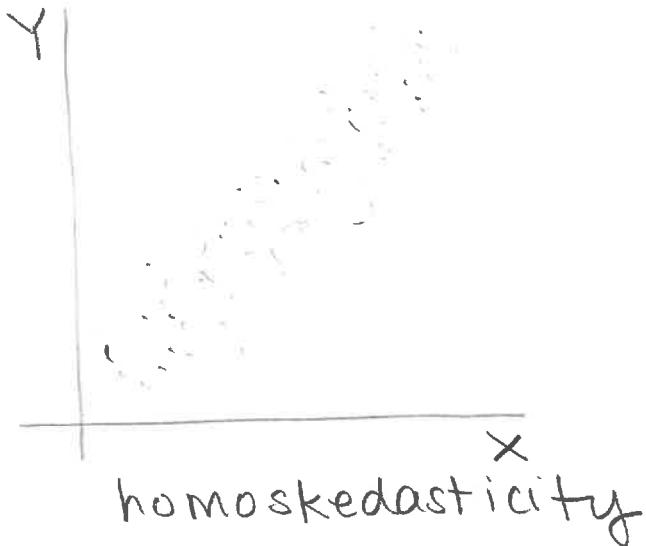
- What does that look like?
- What problems does it cause?
- how do we test for it?
- how do we fix the problems?

\* What does it look like?



①

$$V[u | X = x_i] = \sigma^2$$



① What problems does it create?

Reminder: Homoskedasticity was not an assumption needed for unbiasedness  
 ↳ heteroskedasticity can't lead to bias

The Problem:

→ need an estimate for the variance of the  $\beta$ 's to test for significance

↳ no  $V[\hat{\beta}]$  → no distribution  
 ↳ no t-tests

- the estimate of  $V[\hat{\beta}]$  requires an estimate for  $V[u]$

↳ can't assume  $V[u]$

↳ can't assume  $V[\hat{\beta}]$

↳ can't test for  $\hat{\beta}$  significance

→ so you'll have the right  $\beta$  (it's unbiased) but you can't say whether it's statistically  $\neq 0$ .

### \* How do we test for it?

Null Hyp: Everything is homoskedastic and grand.

$$V[u | x = x_i] = \sigma^2 \leftarrow \text{a constant}$$

→ We'll do some statistics and either reject the null (heteroskedastic) or fail to reject.

What is the definition of variance?

$$V[u] = E[(u - E[u])^2]$$

$\uparrow$   
 $0$

$$= E[(u - 0)^2]$$

(3)

$$V[u|x] = E[u^2|x]$$

$$V[u|x] = \sigma^2 \text{ (under null)}$$

$$\underline{E[u^2|x] = \sigma^2 \text{ (under null)}}$$

→ under the null hypothesis, there's no relationship between  $u^2$  & the  $x$ 's.

→ Let's say we do our normal regression:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \hat{u}$$

if we take these & square them

$$\hat{u}^2 = \hat{\delta}_0 + \hat{\delta}_1 x_1 + \hat{\delta}_2 x_2 + \hat{v}$$

→ if our null hyp is true,  $\hat{\delta}_1 = \hat{\delta}_2 = 0$

→ we know how to test that!

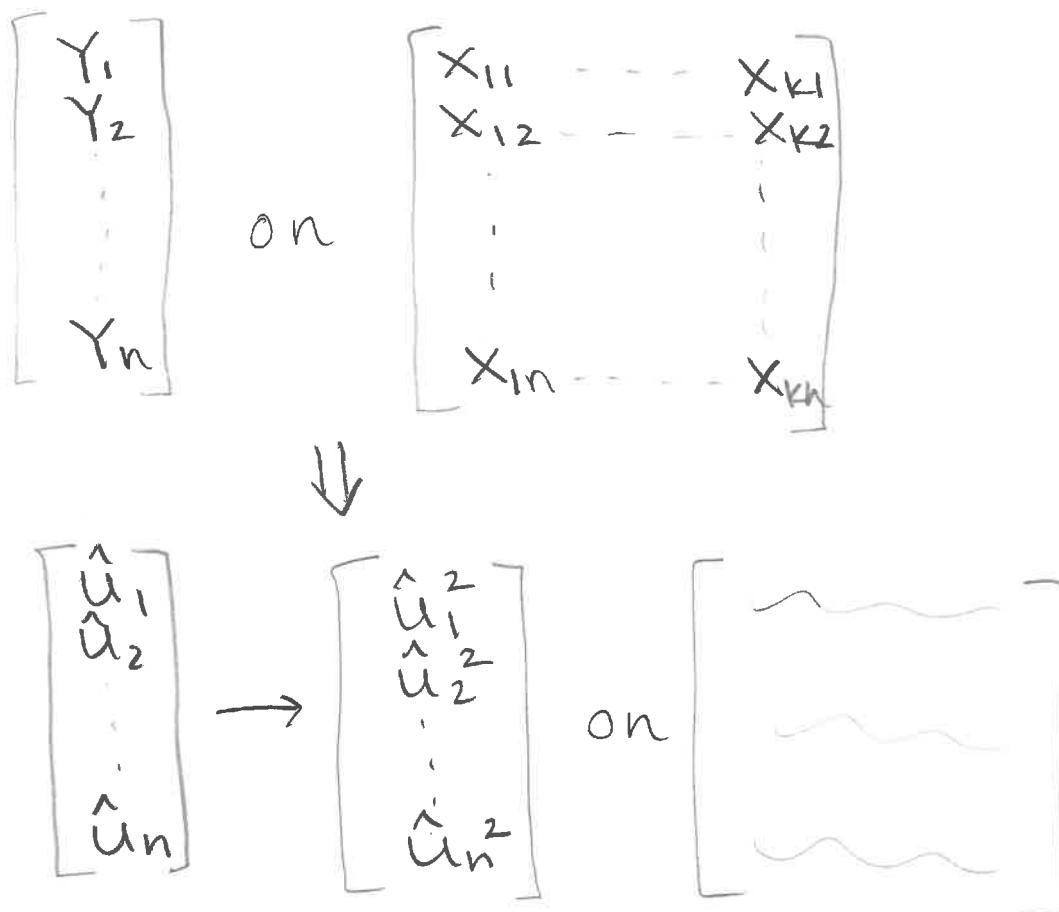
↳ just an F-test

$$H_0: \hat{\delta}_1 = \hat{\delta}_2 = 0$$

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

# Bruesch - Pagan Test Outline:

1. Regress  $Y$  on  $X$ 's.
2. Calculate the residuals,  $\hat{u}$
3. Square them
4. Regress  $\hat{u}^2$ 's on  $X$ 's
5. Jointly test  $X$  coefficients = 0
  - if reject null: heterosk.
  - if fail to reject: homosk.



→ it's so cool!

↳ we made this data! But  
it can tell us something  
meaningful about the world.

\* How do we fix the problems it causes?

→ the problem is that we have no estimate for the variance / std error of the  $\hat{\beta}$ 's.

In the simple regression case

$$Y = \alpha + \beta X + u$$

$$V[\hat{\beta}] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{\left[ \sum_{i=1}^n (x_i - \bar{x})^2 \right]^2}$$

(we can do a multiple regression version but it's notationally crazy)

→ if there's homoskedasticity

$$\sigma_i^2 = \sigma^2 \forall i$$

↳ doesn't vary with  $i$

$$V[\hat{\beta}] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma^2}{\left[ \sum_{i=1}^n (x_i - \bar{x})^2 \right]^2} = \frac{\sigma^2 \sum_{i=1}^n (x_i - \bar{x})^2}{\left[ \sum_{i=1}^n (x_i - \bar{x})^2 \right]^2}$$

$$= \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \begin{matrix} \leftarrow \text{estimate} \\ \leftarrow \text{calculate done!} \end{matrix}$$

→ but if heteroskedastic  
can't pull that  $\sigma^2$  out

$$V[\hat{\beta}] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{\left[ \sum_{i=1}^n (x_i - \bar{x})^2 \right]^2}$$

→ but we can estimate this in  
pretty good way

↳ use  $\hat{u}_i^2$  as an estimator for  $\sigma_i^2$

→ so we can just use:

$$V[\hat{\beta}] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{u}_i^2}{\left[ \sum_{i=1}^n (x_i - \bar{x})^2 \right]^2}$$

(\*) "Pretty Good" in what way?

- this method is valid "as the sample size tends to infinity"  
→ for large samples this is good.

- if no heteroskedasticity, the regular t-tests are valid for any sample size.

- heteroskedasticity is not a huge problem
- ↳ check for it.
  - ↳ if present, use  $\hat{U}$  ("robust standard errors")
  - ↳ if not, don't
- it's better if it isn't there but if it is, hopefully you have large enough data sets that it isn't a problem.
- 

Stata!

load cps data

scatter lwage exper

→ eh hard to say

scatter lwage educ

→ yep. heterosk.

reg lwage educ exper expersq

hettest

→ reject null

reg lwage educ exper expersq,  
robust

# Binary Variables: Interaction Terms

→ Before: used interaction terms to allow for differential effects

$$\text{sqft} \cdot \text{bedrooms}$$

→ Binary Variables can also have differential effects

Example: the effect of being married on your wage may be different if you're male or female

Example: the effect of years of education may be different if you're male or female.

## 2 Types:

1. Binary - Binary interactions  
female - married

2. Binary - Qualitative interactions  
female - education

→ using two binary variables allow for four different intercepts

female = 1 female  
              0 male

married = 1 married  
              0 not married

What's the base case?

male and not married

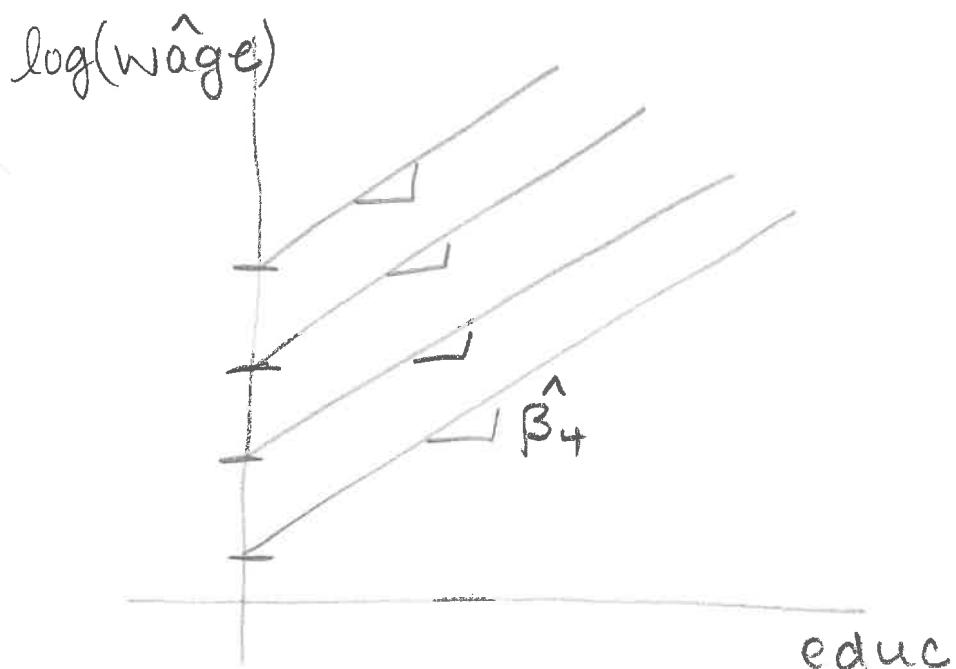
Four Cases:

|         |   | female        |             |
|---------|---|---------------|-------------|
|         |   | 1             | 0           |
| married | 1 | married women | married men |
|         | 0 | single women  | single men  |

base case

→ each of these have a different intercept

$$\log(\hat{wage}) = \hat{\beta}_0 + \hat{\beta}_1 \text{female} + \hat{\beta}_2 \text{married} + \hat{\beta}_3 \text{female} \cdot \text{married} + \hat{\beta}_4 \text{educ}$$



→ We'll calculate these things in a moment.

→ but this still assumes the returns to education are the same for everyone ( $\hat{\beta}_4$ )

↳ to allow for differential returns (different slopes)

female · educ

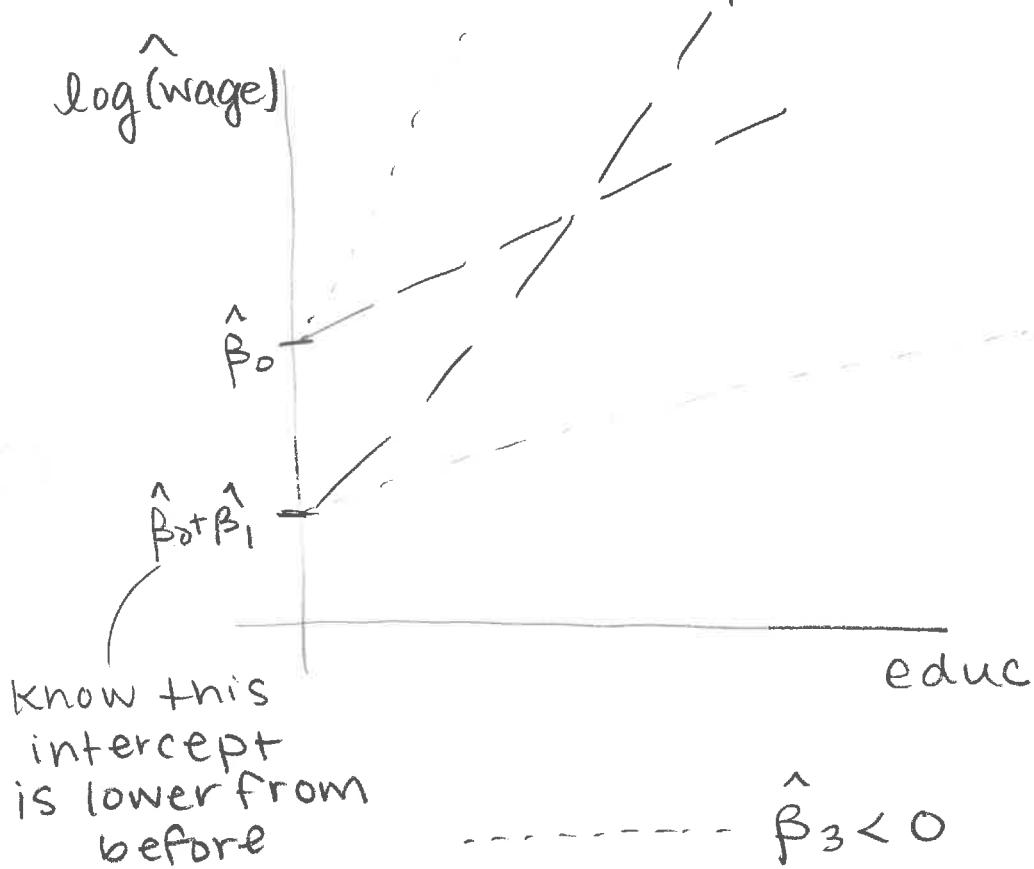
$$\log(\hat{\text{wage}}) = \hat{\beta}_0 + \hat{\beta}_1 \text{female} + \hat{\beta}_2 \text{educ} + \hat{\beta}_3 \text{female} \cdot \text{educ}$$

$$\begin{aligned} E[\log(\hat{\text{wage}}) | \text{female} = 0] &= \hat{\beta}_0 + 0 + \hat{\beta}_2 \text{educ} \\ &\quad + \hat{\beta}_3 \cdot 0 \cdot \text{educ} \\ &= \hat{\beta}_0 + \hat{\beta}_2 \text{educ} \end{aligned}$$

base case relationship  
③

$$E[\log(\text{wage}) | \text{female} = 1] = \hat{\beta}_0 + \hat{\beta}_1(1) \\ + \hat{\beta}_2 \text{educ} + \hat{\beta}_3(1) \text{educ}$$

$$= (\underbrace{\hat{\beta}_0 + \hat{\beta}_1}_{\text{female intercept}}) + (\underbrace{\hat{\beta}_2 + \hat{\beta}_3}_{\text{female slope}}) \text{educ}$$



$$\cdots \cdots \cdots \hat{\beta}_3 < 0$$

→ female has flatter slope

→ returns  $(\hat{\beta}_2 + \hat{\beta}_3)$  are smaller

$$\cdots \cdots \cdots \hat{\beta}_3 > 0$$

→ female has steeper slope

→ returns  $(\hat{\beta}_2 + \hat{\beta}_3)$  larger

## (\*) Joint Testing

Question: Is there any effect on wage from being a female?

→ need to test  $\hat{\beta}_1$  &  $\hat{\beta}_3$  jointly

$$H_0: \hat{\beta}_1 = \hat{\beta}_3 = 0$$

$H_a: H_0$  is not true  
(at least 1 is not = 0)

(could throw married int too)

→ unrestricted model:

$$\begin{aligned} \text{log}(\hat{\text{wage}}) = & \hat{\beta}_0 + \hat{\beta}_1 \text{female} + \hat{\beta}_2 \text{educ} \\ & + \hat{\beta}_3 \text{female} \cdot \text{educ} \end{aligned}$$

→ restricted model:

$$\text{log}(\hat{\text{wage}}) = \hat{\beta}_0 + \hat{\beta}_2 \text{educ}$$

→ 2 restrictions ( $q = 2$ )

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

# Stata!

- load cps 75-85 data

## ① Female & Married

- create female\_married
- regress lwage female married  
female\_married educ

$$\hat{lwage} = 0.87 - 0.15 \text{female} + 0.25 \text{married}$$
$$- 0.16 \text{female\_married} + 0.073 \text{educ}$$

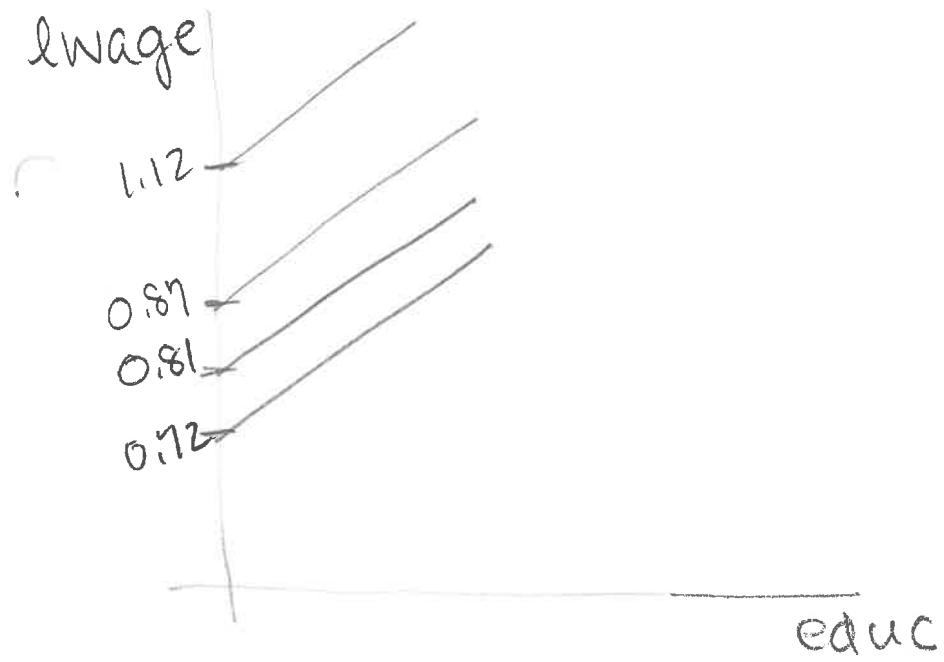
- graph

$$E[\hat{lwage} | \text{female} = 0, \text{married} = 0] \\ = 0.87 + 0.073 \text{educ}$$

$$E[\hat{lwage} | 0, 1] = \underbrace{0.87}_{1.12} + 0.25(1) \\ + 0.073 \text{educ}$$

$$E[\hat{lwage} | 1, 0] = \underbrace{0.87}_{0.72} - 0.15(1) + 0.073 \text{educ}$$

$$E[\hat{lwage} | 1, 1] = \underbrace{0.87}_{0.81} - 0.15(1) \\ + 0.25(1) - 0.16(1) + 0.073 \text{educ}$$



## ② Female & Education

- gen female • educ

- reg lwage female educ female-educ

$$\begin{aligned} \text{lwage} = 1.14 - 0.51 \text{female} + 0.065 \text{educ} \\ + 0.018 \text{female-educ} \end{aligned}$$

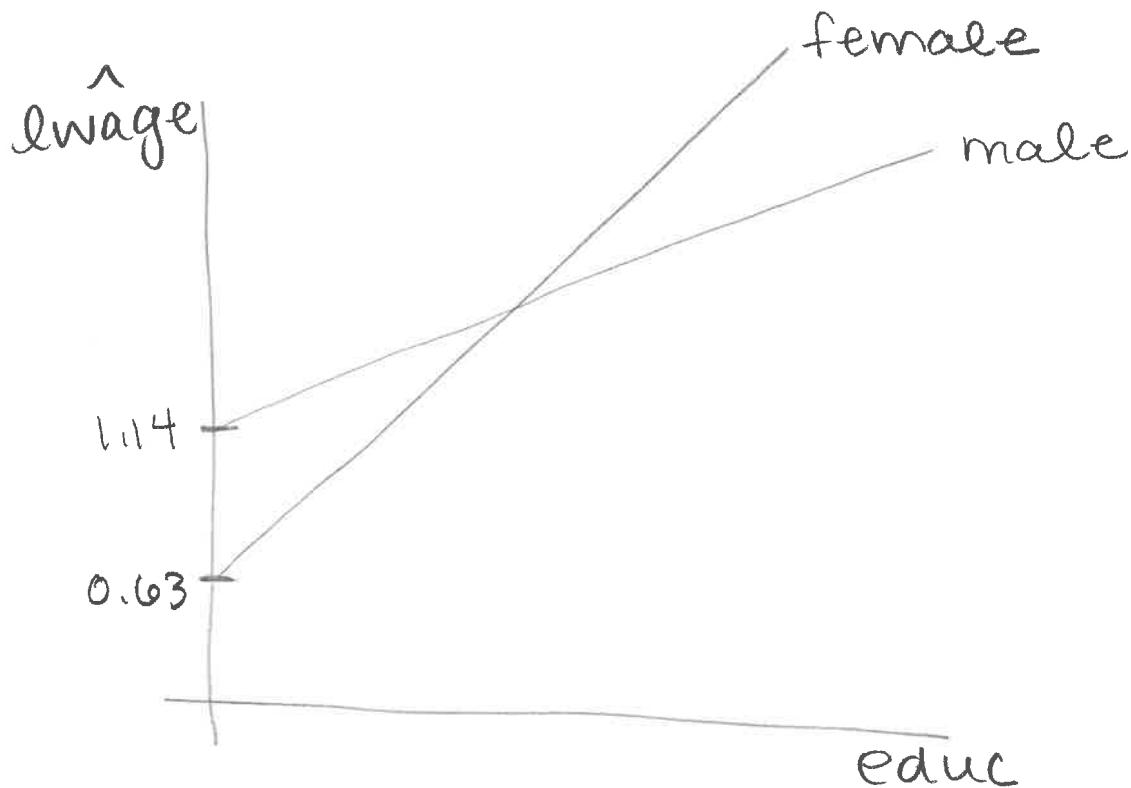
- graph

$$E[\text{lwage} | \text{female}=0] = 1.14 + 0.065 \text{educ}$$

$$E[\text{lwage} | \text{female}=1] = 1.14 - 0.51$$

$$+ 0.065 \text{educ} + 0.018 \frac{\text{female}}{\text{edu}}^{(1)}$$

$$= 0.63 + 0.083 \text{educ}$$



- men have a lower intercept  
(w/ no education, men's wages are higher)
- women have a steeper slope  
(an increase of education of one year earns women a larger increase in wage)

- test

Does being a woman affect your wage?

$$H_0: \hat{\beta}_1 = \hat{\beta}_3 = 0$$

$H_a$ :  $H_0$  is not true

test  $(\text{female} = 0)(\text{female} - \text{educ} = 0)$

(8)

$$F(2, 1080) = 41.44$$

$$P\text{-value} = 0.000$$

reject the null that both  
are 0.

# Time Series Econometrics

✳ This last section of the class is meant to expose you to some of the different genres of econometrics  
(techniques)

- time series (today)
- panel data
- differences in differences
- big data

→ each of these could be their own course

↳ our goal is to get an introduction to each

↳ if you have to use any of these in the future...

- 1) you know you need to do something different from the standard
- 2) you know where to start

→ Up First: Time Series Econometrics

\* Time Series Econometrics is just econometrics applied to time series data

↳ data that occurs across time.

Examples:

- annual GDP
- quarterly sales
- quarterly unemployment
- weekly bond rates

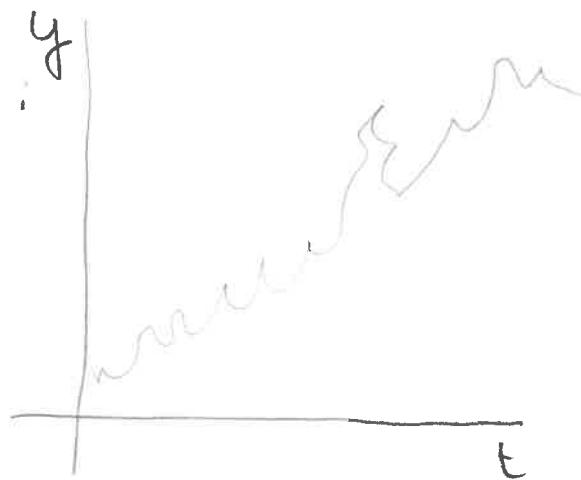
→ annual, quarterly, weekly, daily  
→ you should always specify the periodicity of your data

→ typically, time series data are used a lot in macro

↳ plenty of exceptions: home prices, sales, etc.

→ graphing time series:

- x: time
- y: data



↳ comes with an automatic x-var built in.

## \* Notation

→ the value of  $y$  at time  $t$  may depend on the value of  $x$  at time  $t$

$$y_t = \alpha + \beta x_t$$

- house prices today depend on interest rates today
- inflation this year depends on unemployment this year

→ the value of  $y$  at time  $t$  may depend on the value of  $x$  at time  $t-1$ .

$$y_t = \alpha + \beta x_{t-1}$$

- stock prices today depend on consumer sentiment yesterday
- tax revenue this year depends on employment last year.

→ the value of  $y$  at time  $t$  may depend on the value of  $y$  at time  $t-1$

$$y_t = \alpha + \beta y_{t-1}$$

- aggregate output this year depends on aggregate output this year
- stock prices today depend on stock prices yesterday. (3)

→ When we use a previous time period of a variable, we call it a lag of that variable.

- $y_{t-1}x_{t-1}$ : first lag of  $y, x$
- $x_{t-2}$ : 2nd lag
- $x_{t-h}$ :  $h^{th}$  lag

("lag" looks a lot like "log"; make sure you read it right)

| time period | $y_t$ | $y_{t-1}$ | $x_t$ | $x_{t-1}$ |
|-------------|-------|-----------|-------|-----------|
| 0           | 5     | —         | 20    | —         |
| 1           | 9     | 5         | 17    | 20        |
| 2           | 7     | 9         | 35    | 17        |
| 3           | 2     | 7         | 42    | 35        |
| 4           | 4     | 2         | 57    | 42        |
| 5           | 9     | 4         | 19    | 57        |

"the value of  $y_{t-1}$  at  $t=2$  is 9" but it isn't useful because it doesn't match with anything → need a value for all vars to include in regression

→ Note: by including 1 lag, we lose one row of data → num obs decreases by 1 (4)

- ✳ When you include a lag of  $y$ , it's likely the errors are correlated
  - ↳ this is called serial correlation and it has an entire chapter to itself.
  - ↳ we're going to focus on the lags of other variables

$$y_t = \beta_0 + \delta_0 x_t + \delta_1 x_{t-1} + \dots + \delta_h x_{t-h} + u$$

"Finite Distributed Lag (FDL)  
Model of Order  $h$ "

### ✳ Interpretation

$\delta_0$ : the immediate change in  $y$  from a 1 unit change in  $x$

$\delta_1$ : the change in  $y$  from a 1 unit change in  $x_{t-1}$

↳ the effect on  $y$  now from a change yesterday

$\delta_h$ : the effect on  $y$  now from a 1 unit change in  $x$   $h$  periods ago (5)

## \* Trends & Seasonality

→ many variables contain a time trend

↳ they're growing (or shrinking) over time



→ if we just regress one on the other, we may find false relationships because they're both moving together in time.

→ include t (omitted var bias otherwise)

$$Y_t = \beta_0 + \delta_0 x_t + \delta_1 x_{t-1} + \beta_1 t + u$$

→ the coefficient on t describes the general time trend

↳  $\beta_1$ : the change in  $Y$  associated with a unit increase in time

↳ time moving forward by 1 period (6)

→ Some data vary with certain types of times.

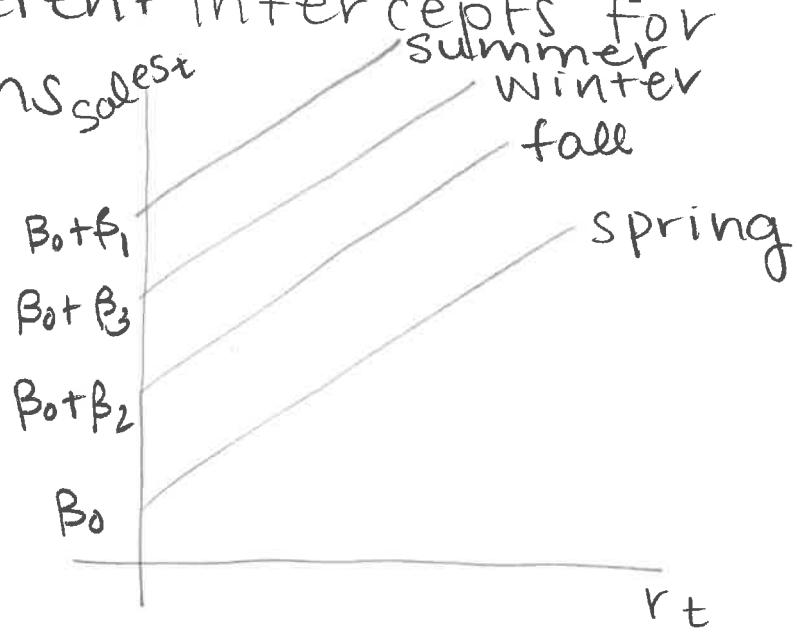
- house sales (& rentals as many of us know) are much higher in the summer than in other seasons.
- gym memberships in January
- output growth during wartime years

→ We deal with this by including dummies for (each - 1) "season" (month, year, etc.)

$$\text{house\_sales}_t = \beta_0 + \beta_1 r_t + \beta_1 \text{summer} \\ + \beta_2 \text{fall} + \beta_3 \text{winter} + u$$

(Always remember to omit one category.)

→ creates different intercepts for different seasons



# Stata!

\* What is the effect of tax exemptions for children on the US fertility rate?

fertility rate: births per 1000 women

→ When you get more money back on your taxes per child, how does it affect the number of children you have?

Load fertil3.dta  
\* tell stata this is a time series  
sum gfr pe year

tsset year

① FDL Order 2

$$gfr_t = \beta_0 + \delta_0 p_{e_t} + \delta_1 p_{e_{t-1}} + \delta_2 p_{e_{t-2}} + u_t$$

→ You can generate the lag variables directly with gen.

↳ but you can just type ↳ before a variable to lag it

reg gfr pe L<sub>-</sub>pe L<sub>2-</sub>pe

→ slightly different output

L<sub>-</sub>pe: the effect on fertility today of the pe increasing last year

→ When the pe increases by 1 a year ago, the number of children decreases by 0.02

test joint sign. of all pe's  
↳ not jointly significant

\* Trends & Seasonality

- year
- WWII

- the advent of the pill

$$gfr_t = \beta_0 + \beta_0 \cdot pe_t + \beta_1 WW2_t + \beta_2 pill + \beta_3 t + u_t$$

(if time, drop year → pill is sig)

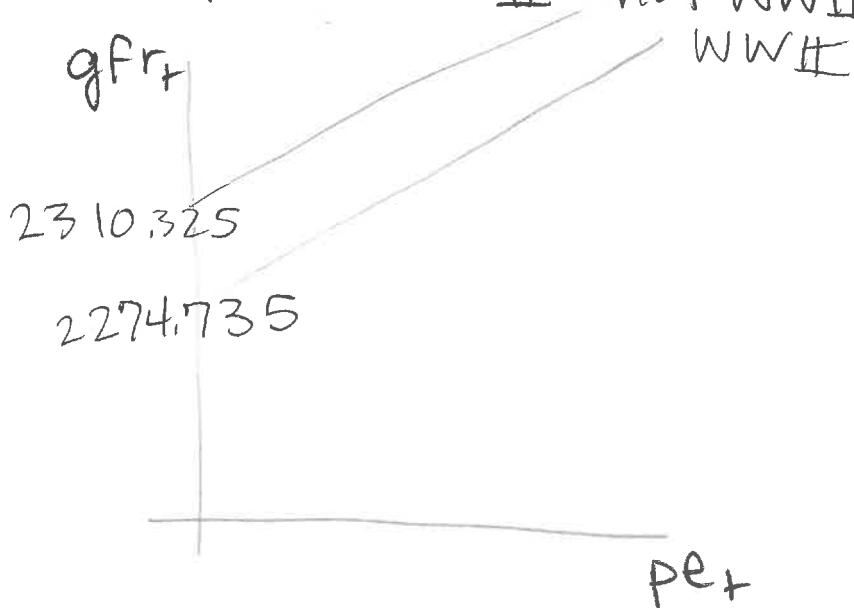
↳ time trend might have been the real driver

(\*) What kind of time trend?

→ negative

→ on average, every year  
the number of  
births per 1000 women decreases  
by 1.14 children

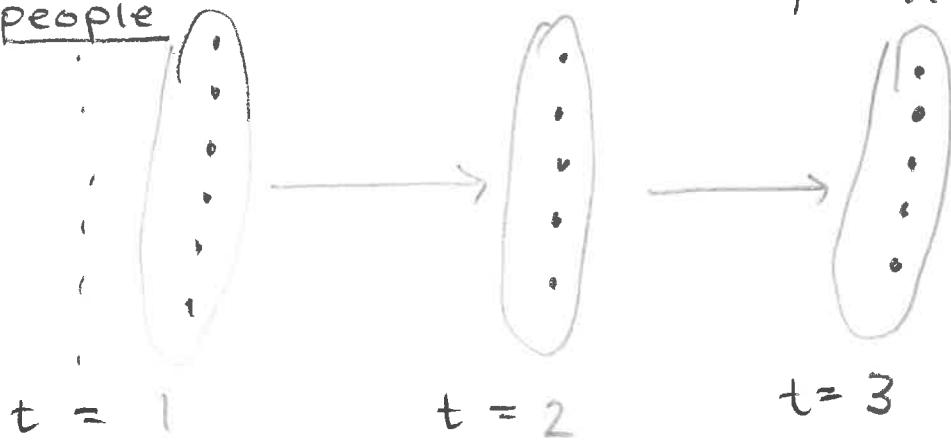
(\*) Graph the effect of  $p_{et}$  on  $gfr_t$ .  
Compare WW II not WW II.



## Difference in Differences ("Diff-in-Diff")

\* The next two classes deal with data on a cross section, collected over time

people



→ allows us to track changes over time

→ particularly useful for policy evaluation.

### Examples:

- Did a work training program lead to a decrease in unemployment?
- Does subsidizing health insurance improve health?
- Did the introduction of the birth control pill decrease the fertility rate?

\* There are 2 ways to collect this data:

1. Sample the same people from a population ~~across~~ multiple time periods

→ "Panel Data" (makes a panel)

| Person | t=1 | t=2 |
|--------|-----|-----|
| 1      | 7   | 9   |
| 2      | 2   | 9   |
| 3      | 10  | 10  |
| 4      | n   | n   |
| 5      | n   | n   |

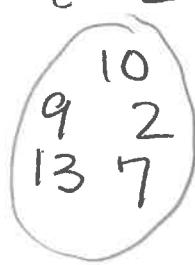
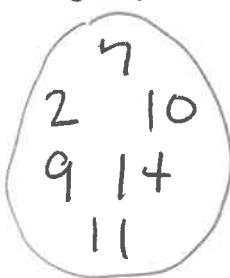
(Next Class)

→ [class example]

2. Sample a different random sample of people from a population in each time period

→ "Pooled Cross Section" (You pool the data)

t=1                    t=2



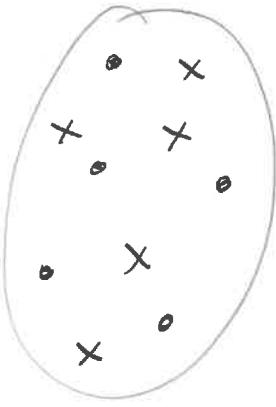
(Today)

→ [class example]

Difference-in-Differences uses a pooled cross section to analyze a policy or change.

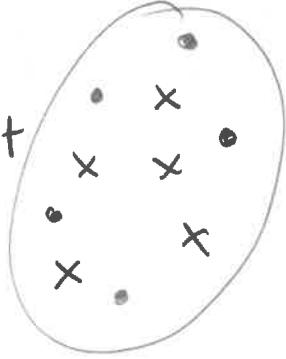
\* Ingredients for a Diff-in-Diff Estimation

- 2 Time Periods
  - before & after the policy applies
  - Why needed: there are other things that change over time besides the policy and don't want to erroneously attribute that to the policy.
- 2 Different Groups of [People]
  - one to which the policy applies ("Treatment Group")
  - one to which the policy does not apply ("Control Group")
  - Why needed: there may be things the people in the treatment (or control) have in common
    - ↳ don't want to erroneously capture that either.



$t=1$   
(Before)

- control
- ✗ treatment



$t=2$   
(After)

## \* How to Implement a Diff-in-Diff Estimation

→ create a dummy variable for the treatment group

$$\text{treat} = \begin{cases} 0 & \text{control} \\ 1 & \text{treatment} \end{cases}$$

→ create a dummy variable for one of the time periods

$$t-2 = \begin{cases} 0 & t=1 \\ 1 & t=2 \end{cases}$$

→ interact them

$$\text{treat} \cdot t-2 = \text{treat} \times t-2$$

(\*) Regress!

Because the time is capturing the policy.  
 Treat is allowing for similarities in the groups that aren't policy related.

$$Y = \beta_0 + \delta_0 t - 2 + \beta_1 \text{treat} + \delta_1 \text{treat}_{-t-2}$$

$\beta_0$  = average outcome for the control in 1st time period.

$\delta_0$  = change for all [people] between the two time periods

$\beta_1$  = difference between the two groups of people not due to the policy

$\delta_1$  = the diff in diff estimator  
 = the change in outcome before and after the policy due to the policy

|                     | $t=1$<br>Before     | $t=2$<br>After                            | After-Before          |
|---------------------|---------------------|-------------------------------------------|-----------------------|
| Control             | $\beta_0$           | $\beta_0 + \delta_0$                      | $\delta_0$            |
| Treatment           | $\beta_0 + \beta_1$ | $\beta_0 + \delta_0 + \beta_1 + \delta_1$ | $\delta_0 + \delta_1$ |
| Treatment - Control | $\beta_1$           | $\beta_1 + \delta_1$                      | $\delta_1$            |

\* Example: The effect of a new garage incinerator on nearby housing prices

$$\text{near-inc} = \begin{cases} 0 & \text{if a house is } > 3 \text{ mi. from the location of the incinerator} \\ 1 & \text{if } < 3 \text{ mi} \end{cases}$$

$$y_{-81} = \begin{cases} \text{Year} = 1978 \\ \text{Year} = 1981 \end{cases}$$

$$\text{price} = \beta_0 + \delta_0 y_{-81} + \beta_1 \text{near-inc} + \underline{\delta_1 y_{-81} \cdot \text{near-inc}} + u$$

$\delta_0$  = the effect on all housing prices over time (inflation)

$\beta_1$  = the location effect not due to the incinerator (farther from most jobs)

$\delta_1$  = the effect on housing prices near the incinerator due to the incinerator

Britney Example

## \* Stata!

Load INJURY.dta

- policy: the maximum amount of weekly earnings that were covered by workers comp. increased.
- if you have a low wage, you weren't hitting the old cap, this won't affect you

↳ Control

→ high wage → treatment

→ data before and after the change ("afchange")

\* Does this policy increase the duration of workers comp claims?

gen after\_high = afchange \* highearn  
reg log(durat) afchange highearn after\_high

$$S_i = 0.188$$

→ duration increased by about 19% due to the policy

→ coeff on afchangeinsig.  
↳ didn't affect low income  
workers as expected.

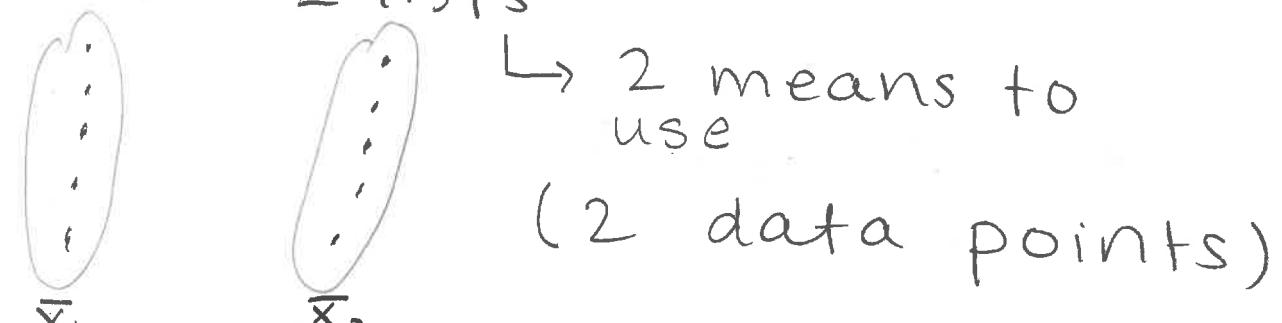
## Panel Data: Fixed Effects Model

→ rather than pooled data across people, cities, schools, etc.

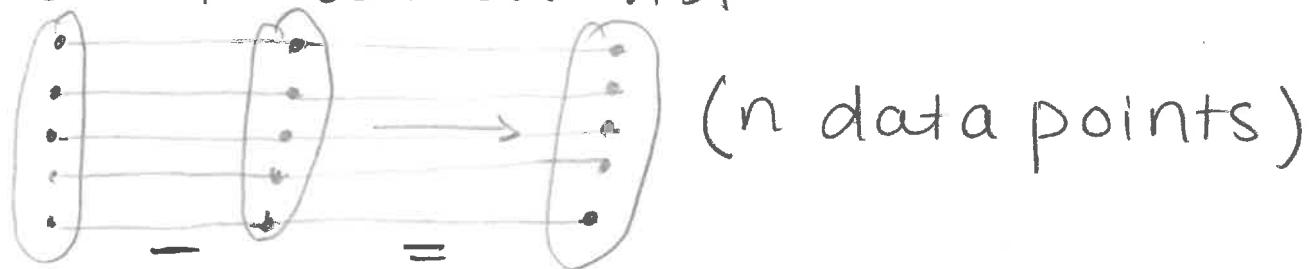
↳ track the same people, cities, schools, etc. in multiple time periods

(\*) Why this is useful: contains more information.

→ Intuition: it's like when you have a paired 2-sample t-test, & you're interested if the diff between them is  $> 0$ .  
↳ if it's unpaired you just have 2 lists



↳ but if it's paired, you can create a new list



## Organizing Panel Data:

| individual | Variable X | $x_{it}$ |          |
|------------|------------|----------|----------|
|            | $t=1$      | $t=2$    | $t=3$    |
| 1          | $x_{11}$   |          |          |
| 2          | $x_{21}$   |          |          |
| 3          | $x_{31}$   | $x_{32}$ | $x_{33}$ |
| 4          | $x_{41}$   |          |          |
| 5          | $x_{51}$   |          |          |

→ this is just one variable, & look at how much space it takes up

→ normally we're working with multiple variables

| indiv | t | Y        | $x^1$      | $x^2$ |
|-------|---|----------|------------|-------|
| 1     | 1 | $y_{11}$ | $x_{11}^1$ |       |
|       | 2 | $y_{12}$ | $x_{12}^1$ |       |
|       | 3 | $y_{13}$ | $x_{13}^1$ |       |
| 2     | 1 |          |            |       |
|       | 2 |          |            |       |
|       | 3 |          |            |       |
| 3     | 1 |          | $x_{31}^1$ |       |
|       | 2 |          | $x_{32}^1$ |       |
|       | 3 |          | $x_{33}^1$ |       |

\* We're using this data to answer questions about how something has changed over time

↳ often before & after a policy is implemented.

→ interested in the effect of some policy ( $\text{policy} = \begin{cases} 0 & \text{otherwise} \\ 1 & \text{if the policy applies to the city} \end{cases}$ ) on some outcome variable

→ why can't we use:

$$y_{it} = \beta_0 + \beta_1 x_{it} + u_{it}$$

First Problem: there are probably

things that changed over time that affect the outcome that aren't the policy

- ↳ don't want to erroneously attribute those to the policy
- ↳ need to take time into account somehow

Second Problem: there are probably people, city, school, etc. specific things that affect the outcome

↳ don't want to attribute those to the policy either.

↳ many of these are unobservable or at least unobserved in your panel

→ new equation:

$$y_{it} = \beta_0 + \beta_1 t - 2_t + \beta_2 x_{it} + a_i + u_{it}$$

outcome

time Period dummy

explanatory var

Individual effect  
→ no + subscript.  
→ only varies by individual  
→ unobserved

④ Dealing with the Individual Effects  
→ also called the Fixed Effects  
(b/c they're fixed over time)

→ there are 3 ways we use to deal with  $a_i$

1. First Differences Model
- \* 2. Fixed Effects Model (14-1)
3. Random Effects Model

most commonly used

→ Basically, we want  $a_i$  to go away but still be able to answer the question. & take advantage of all the power of the panel data

$$y_{it} = \beta_0 + \beta_1 x_{it} + a_i + u_{it}$$

→ for each  $i$ , average the equation over time

$$\begin{aligned} \sum_{t=1}^T y_{it} &= \sum_{t=1}^T (\beta_0 + \beta_1 x_{it} + a_i + u_{it}) \\ &= \sum_{t=1}^T \beta_0 + \sum_{t=1}^T \beta_1 x_{it} + \sum_{t=1}^T a_i + \sum_{t=1}^T u_{it} \end{aligned}$$

$$\bar{y}_i = \beta_0 + \beta_1 \bar{x}_i + a_i + \bar{u}_i$$

→ Subtract  $\bar{y}_i$  from  $y_{it}$  (on both sides of the equation)

→ "time demeaning"

$$\begin{aligned} y_{it} - \bar{y}_i &= \beta_0 + \beta_1 x_{it} + \alpha_i + u_{it} \\ y_{it} - (\beta_0 + \beta_1 \bar{x}_i + \alpha_i + \bar{u}_i) &= \beta_1 (x_{it} - \bar{x}_i) + u_{it} - \bar{u}_i \end{aligned}$$

$$\ddot{y}_{it} = \beta_1 \ddot{x}_{it} + \ddot{u}_{it}$$

→ gets rid of  $\alpha_i$

↳ anything about the individual constant over time, is averaged out

→ but keeps all the power of the panel data

\* Example: Does the use of a grant for worker training decrease the number of scrapped products on an assembly line?

| factory | t | grant | t-2 | grant | t-2 |
|---------|---|-------|-----|-------|-----|
| 1       | 1 | 0     | 0   | 0.5   | 0.5 |
|         | 2 | 1     | 1   |       |     |
| 2       | 1 | 0     | 0   | 0     | 0.5 |
|         | 2 | 0     | 1   |       |     |
| 3       | 1 | 0     | 0   | 0.5   | 0.5 |
|         | 2 | 1     | 1   |       |     |

| factory | t | grant | t-2  |
|---------|---|-------|------|
| 1       | 1 | -0.5  | -0.5 |
|         | 2 | 0.5   | 0.5  |
| 2       | 1 | 0     | -0.5 |
|         | 2 | 0     | 0.5  |
| 3       | 1 | -0.5  | -0.5 |
|         | 2 | 0.5   | 0.5  |

→ all the information

↳ just shifted around the time mean

Stata!

→ open JTRAIN.DTA

→ open data browser  
fcode — firm  
year

→ need to tell stata what the individual variable is (fcode) and the time variable (year)

[xtset fcode year  
state year

[xtreg lscrap d88 d89 grant grant-  
Y  $x_1$   $x_2$   $x_3$   $x_4$   
fe,

→ obtaining a grant in 1988 decreased the scrap rate in 1989  
 $(e^{-0.422} - 1 = -0.344) \ 34.4\%$

# Introduction to Big Data

→ Big Data Econometrics is it's own course, this is merely an introduction

↳ goal: should you need to use Big Data tools you'll have an idea of where to start

→ courses in Computer Science & Statistics

↳ new Data Science major

## \* Big Data

→ it's called big data because the files/databases containing the data are prohibitively big.

↳ standard data-handling program (STATA, Excel, etc) would crash.

→ need special tools & software to:

- store
- access ("query")
- mess with ("computation")

## \* Storage & Accessing

→ usually done across multiple servers

→ Google, Amazon, etc. have made it possible to rent data storage space.

→ for most statistical processes you need to access many if not all values.

↳ just calculating  $\bar{x}$ !

- find each  $x_i$
- query the value of  $x_i$
- add each one
- divide by  $n$  (which is huge)

## \* Most of the time when we refer to Big Data Techniques we're referring to Machine Learning.

↳ usually used to try to predict some  $y$  as a function of some predictors ( $x$ 's).

↳ in the presence of computational challenges. (2)

→ in Econ we're usually trying to detect relationships.

↳ Machine Learning has tools that can help us do that job

\* Two Important Differences

1. Testing & Training

→ an important difference from "Normal" Econometrics: When we have this much data we can use only some of the data to build the model and some of the data to test the model.

→ always viewed as a good idea

↳ "we didn't have enough data"

↳ now we do!

→ "cross validation"

↳ use some subset of the data to build the model (e.g. estimate the regression coefficients)

↳ use the model to guess the corresponding values for another subset (e.g. plug in the  $x$ 's to get  $\hat{y}$ )

↳ see how far off you are  
 $(y - \hat{y})$ ?

## 2. Controlling Complexity

→ We intuitively know that while always including more & more explanatory variables can make the errors smaller

↳ but having too many is cumbersome / overly complex

↳ very non-rigorous way limit the  $x$ 's.

→ Machine Learning goes one step further: mathematically penalizes too many  $x$ 's / complexity

→ For example: Instead of minimizing the sum of squared errors

$$\min (y - \beta_0 - \beta_1 x_1 - \dots - \beta_k x_k)^2$$

minimize that plus some count of the number of nonzero  $\beta$ 's

$$\min (y - \beta_0 + \beta_1 x_1 - \dots - \beta_k x_k + \lambda(\beta_i \neq 0))^2 \quad (4)$$

## \* Two Particular Tools:

1. Classification Trees
2. Variable Selection

→ both concerned with determining what are the most important characteristics in determining some outcome.

### 1. Classification Trees

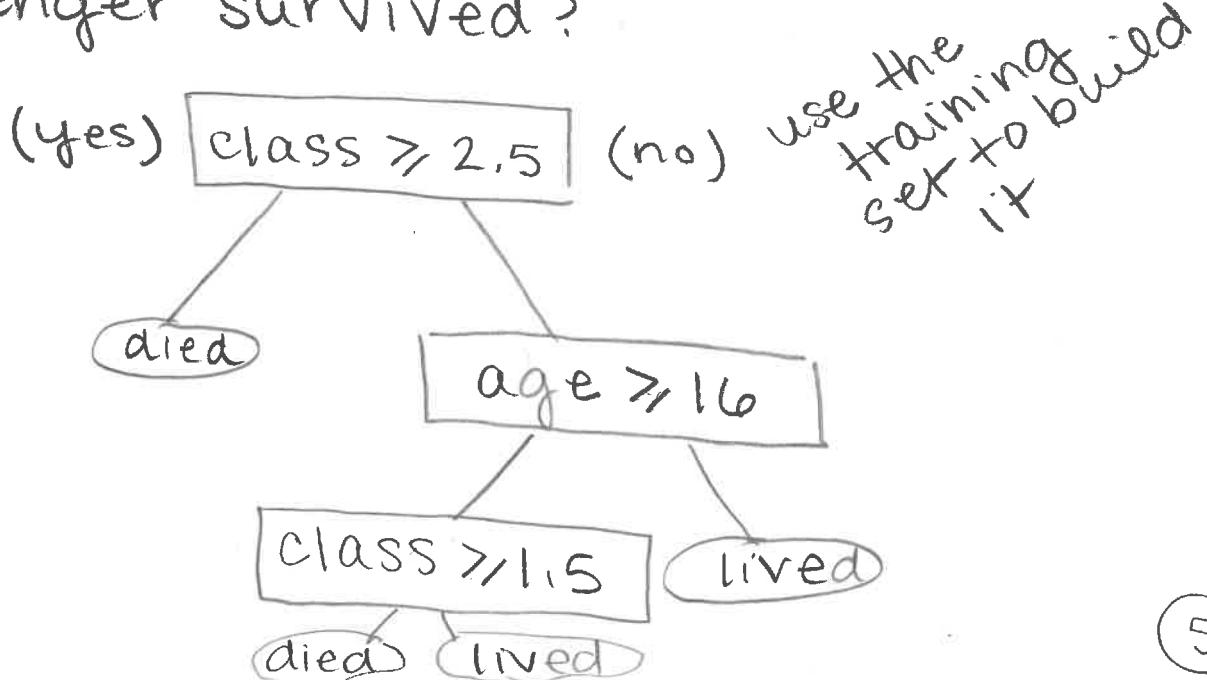
→ binary outcomes

→ trees are a particularly efficient data structure

↳ can handle lots of data

(Dark)

Example: What role did class & age play in whether a titanic passenger survived?



→ use the training set of data to make the tree.

→ then use each person in the testing set and ask the tree for it's prediction (like plinko)

↳ compare the prediction to the real life outcome.

Data

| Class | Age | Prediction | Real Outcome |
|-------|-----|------------|--------------|
| 1     | 17  | L          | L            |
| 3     | 17  | D          | L            |
| 2     | 34  | D          | D            |
| 1     | 80  | L          | L            |
| 2     | 27  | D          | D            |

→ this tree classified 723 of 1046 passengers correctly

## 2. Variable Selection

→ regression but where too many variables are penalized

→ you have so many explanatory variables to choose from

$$\min \left[ (y - \beta_0 - \beta_1 x_1 - \dots - \beta_k x_k)^2 + \lambda \sum_{i=1}^k |\beta_i|^2 \right]$$

~~~~~  
punishes non-zero β 's

Example: What factors are important in determining GDP?

- lots of options
- Sala-i-Martin data has 42 explanatory variables
 - ↳ he enumerated all the combinations
 - ↳ could just use Machine learning regression
 - more comp. efficient
 - does a good job

Most Important Factor:

1. Equipment Investment

* Model Uncertainty

→ we've spent a lot of time throughout this course talking about the uncertainty created by sampling from a population

↳ we almost never talk about the uncertainty / error created by picking one model over another (using age & age^2 but not age, age^2 & age^3)

→ Big Data (both the data itself & the tools) make it possible to average across the predictions of many models!

Example: "Forests" of Trees.

→ create a bunch of different classification trees (a "forest")

↳ ask each for a prediction

↳ average across predictions

→ If 7 of 10 trees say a person would die on the titanic, the prediction is they would die. (8)

→ Econometricians, Computer Scientists and Statisticians alike agree that this is better than just one model.

(*) Big Data Takeaways:

→ Challenges

- needs separate software for storage, access, computation

→ Benefits

- more information
- allows us to do stuff we said we should do but needs lots of data
 - cross validation
 - complexity control
 - using multiple models