

# FALL 2021 EC 282 LECTURE NOTES

09/09/2021

## LESSON PLAN

- SELF-INTRODUCTION
- COURSE INFO + SYLLABUS + EXPECTATIONS
- COURSE LOGISTICS : R-STUDIO

&

STACKOVERFLOW  
GITHUB

## EC282 : INTRODUCTION TO ECONOMETRICS

ECONOMETRICS



ECON

THEORY

LIS MEASUREMENT  
(QUANTIFY)

→ DATA SCIENCE  
STATISTICAL  
LEARNING  
REGRESSION ANALYSIS  
MACHINE LEARNING

AI

→ SMART PHONES

COLLECTING MORE  
STORING DIGITAL  
MANAGING INFORMATION.

4

→ SQL, GRAPH DATA  
BASE

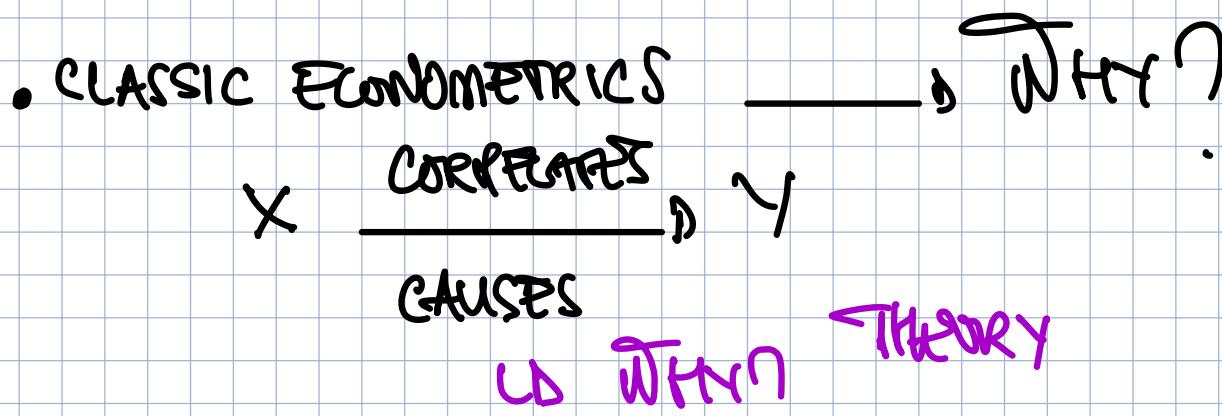
### Two important changes

- BIG DATA - PERSPECTIVE:

WE BECOME CAPABLE OF  
SERVING,  
CLOUDS

- PREDICTION MODELLING CHANGED:  
MACHINE LEARNING (SUPERVISED OR)  
UNSUPERVISED

FY: इयां प्रेडिक्शन अस्टिनेट्स  
(NEVER PERFECT)



Covid-19 JACCINES

CAUSAL RESEARCH  
DESIGN

FDA APPROVAL AFTER LARGE SCALE RCT.

THIS COURSE: INTRODUCTION TO क्षमा:

- LEARN BASIC TOOLS OF PREDICTION:
- DIFFERENTIATE THE CORRELATION & CAUSATION

- SPREAD SHEETS  
ACCESS, EXCEL

↪ मुझे नहीं बहुत धूम्रपान  
IN THE FUTURE FOR  
DATA MANAGEMENT &  
ANALYSIS

MORE COMPLEX DATA SYSTEMS : PYTHON, R → INDUSTRY STANDARD

→ NOT SPSS / SPSS

## CLASS LOGISTICS

WEB PAGE : PDF / HTML VERSIONS

OFFICE HOURS : ONLINE

TEXTBOOK : GET THE GREATEST VERSION

R & R-STUDIO

(SHOW IF WE HAVE TIME)

} STACKOVERFLOW

ECONOMETRICS USING R  
POST ON BIMUB

## GITHUB ACCOUNT

PROGRAMMING : BEST OF YOUR EFFORT IS FINE  
(MORE ON THIS LATER)

GRADING : HIGH STAKES ASSESSMENTS

2 MIDTERMS + FINAL → 70 %

LOW STAKES ASSESSMENTS

1 HOMEWORK ASSIGNMENTS → SUBMITTED THROUGH  
→ 0% PASS / FAIL BLACKBOARD

- PASS / FAIL
- GROUP WORK (MAX 2)
- SUBMIT THROUGH BB (EACH STUDENT)
- DIFFERENT DATASETS

NO EXAMS REQUIRE **PROGRAMMING SKILLS**.

## LD RUNNING CLUB

12/09/2021

LESSON PLAN: REVIEW OF PROBABILITY

- RANDOM VARIABLES
- PROBABILITY DISTRIBUTIONS
- $E[X]$ ,  $JAR[X]$ , STD. DEV

$$M_x \quad J_x^2 \quad J_x$$

- STANDARDIZED RANDOM VAR  $\rightarrow Z$
- 2 RANDOM VARIABLES  $X, Y$
- JOINT & MARGINAL DISTRIBUTIONS
- INDEPENDENCE
- COVARIANCE & CORRELATION

**RANDOM VARIABLES:**  $X, Y, Z$

NUMERICAL SUMMARY OF A RANDOM OUTCOME

- HEADS / TAILS → PURELY RANDOM

• COVID INFECTION → RANDOM COMPONENT +  
 YES / NO DETERMINISTIC  
 COMPONENT

$$\hookrightarrow \{0, 1\} \rightarrow Y = \{0, 1\}$$

if YES  
NO UNLIKELY

YOUR GRADE FROM ECONOMETRICS:

$$Y = \{0, 0.7, 1, 1.3, 1.7, 2, 2.3, 2.7, 3, 3.3, 3.7, 4\}$$

UNLIKELY

DISCRETE RANDOM VARIABLES: → IT IS A COUNTABLE  
 NB OF OUTCOMES  
 WE WILL MOSTLY WORK WITH THESE

CONTINUOUS RANDOM VARIABLES: ANY NUMERICAL  
 VALUE IN AN INTERVAL OR COLLECTION OF  
 INTERVALS.

EXPECTED VALUE: LONG-RUN AVERAGE

$$E[Y], \mu_Y$$

$$Y = \{Y_1, Y_2, \dots, Y_k\} \text{ OUTCOMES}$$

$$Pr(Y=Y_1), Pr(Y=Y_2), \dots, Pr(Y=Y_k)$$

DISCRETE PROBABILITIES

$$E[Y] = \sum_i Y_i P_i = Y_1 P_1 + Y_2 P_2 + \dots + Y_k P_k$$

↳ WEIGHTED AVERAGE OF OUTCOMES

↳ WEIGHTS → DISCRETE PROBABILITIES

$$\sum_i P_i = 1$$

HANDOUT QUESTION:

$Y = \{0, 1\} \rightarrow$  BERNULLI / BINARY

$$PR(Y=1) \rightarrow P$$

$$PR(Y=0) \rightarrow (1-P)$$

$$E[Y] = 0 \times (1-P) + 1 \times P$$

$$= P$$

$$Y = \{0, 1\} \rightarrow \text{COVID}$$

$$PR(Y=1) = 0.01 \quad E[Y] = 0 \times 0.99 +$$

$$PR(Y=0) = 0.99 \quad 1 \times 0.01$$

$$= 0.01 \rightarrow P$$

$$\text{VARIANCE: } \text{VAR}(Y) \rightarrow \sigma_Y^2$$

MEASURES THE WEIGHTED SPREAD OF OUTCOMES AROUND THE LONG RUN AVERAGE ( $\mu_Y$ )

$$\begin{aligned}\sigma_Y^2 &= (Y_1 - \mu_Y)^2 p_1 + (Y_2 - \mu_Y)^2 \cdot p_2 + \dots \\ &\quad + (Y_K - \mu_Y)^2 \cdot p_K \\ &= \sum_{i=1}^K (Y_i - \mu_Y)^2 p_i\end{aligned}$$

) SAME THING!

$$\sigma_Y^2 = \sum (Y_i - \mu_Y)^2 \rightarrow \sigma_Y = \sqrt{\sigma_Y^2}$$

$$F(Y = 0, 1) =$$

$$\begin{array}{ll} Y_1 = 0 & Y_2 = 1 \\ (1-p) & p \end{array}$$

$\mu_Y = p$

(STANDARD DEVIATION)

WE LIKE THIS BETTER

(UNIT OF MEASURED IS THE SAME AS  $Y$ )

$$\begin{aligned}\sigma_Y^2 &= (0-p)^2 (1-p) + \\ &\quad (1-p)^2 p\end{aligned}$$

$$= p^2 (1-p) + (1-p)^2 p$$

$$= p(1-p)[p - (1-p)]$$

$$= p(1-p)$$

## HANDBOT QUESTION:

$$\sigma_y^2 = 0.01 \times 0.99$$

$$\sigma_y^2 = 0.0099 \quad \sigma_y = 0.0995$$

TWO RANDOM VARIABLES

X AND Y : BOTH DISCRETE RANDOM VARIABLES.

UPPER CASE

$\text{PR}(X=x)$  } PROB OF  $X \text{ AND } Y$  EQUAL TO  
 $\text{PR}(Y=y)$  } A SPECIFIC VALUE OF  $X, Y$

LOWER CASE

COLUMN MARGINAL PROBABILITY DISTRIBUTIONS

$\text{PR}(X=x, Y=y) \rightarrow \text{PROB. THAT } X=x \text{ AND }$

$\text{PR}(Y=y) = \sum_i \text{PR}(X=X_i, Y=y)$

HANDBOT

$Y = \{0, 1\}$

↓

WE SHORT

LONG

COMMUTE

COMMUTE

$X = \{0, 1\}$

↓

NO RAIN

RAIN

$$PR(Y=0) = PR(X=0, Y=0) + PR(X=1, Y=0)$$

$$PR(X=0, Y=0) = 0.15$$

$$PR(x=0, y=1) = 0.15$$

$$P(X=1, Y=0) = 0.07$$

$$PR(X=1, Y=1) = 0.63$$

MUTUALLY  
EXCLUSIVE

$$\sum_j \Pr(X=x_j, Y=y_j) = 1$$

$$\Pr(Y=1) = \Pr(X=0, Y=1) + \Pr(X=1, Y=1)$$

$$= 0.15 + 0.63$$

$$= 0.78$$

2015

# THEOREM

$$\Pr(Y=y \mid X=x) = \frac{\Pr(Y=y, X=x)}{\Pr(X=x)}$$

THE PROBABILITY THAT  $Y$  IS EQUAL TO  $y$   
VALUE CONDITIONAL ON  $X$  IS EQUAL TO  $\lambda$

PR (Short commute | RTW)

09/16/2021

- LESSON PLAN:
- BAYES THEOREM + CONDITIONAL
  - COND. EXPECTED VAL. PROBABILITY
  - LIE: LAW OF ITERATED EXPECTATIONS
  - INDEPENDENCE, COVARIANCE, CORRELATION

$$E[Y] = \sum_i y_i \Pr(Y=y_i)$$

$$E[Y|X=x] = \sum_i y_i \Pr(Y=y_i | X=x)$$

↳ Roll A DIE  $\rightarrow$  TWO RANDOM VARIABLES

$$Y = \{1, 2, 3, 4, 5, 6\} \quad \Pr(Y=y_i) = 1/6$$

$$X = \{0, 1\}$$

$\downarrow$  ODD  
EVEN

$$\Pr(X=0) = 1/2$$

$$\Pr(X=1) = 1/2$$

$$E[Y] = \sum_i y_i \Pr(Y=y_i)$$

$$= 1 \times 1/6 + 2 \times 1/6 + \dots + 6 \times 1/6 = 3.5$$

$$E[Y|X=1] = \sum_i y_i \Pr(Y=y_i | X=1)$$

$$= 1 \times \frac{\Pr(Y=1, X=1)}{\Pr(X=1)} + 2 \times \frac{\Pr(Y=2, X=1)}{\Pr(X=1)}$$

$$+ \dots + 6 \times \frac{PR(Y=6, X=1)}{PR(X=1)}$$

$$= 1 \times \frac{1}{2} + 0 + 3 \times \frac{1}{3} + 0 + 5 \frac{1}{2} + 0$$

$$\rightarrow E[Y|X=1] = 3$$

**LAW OF ITERATED EXPECTATIONS**

$$E[Y] = \sum_i E[Y|X=x_i] PR(X=x_i)$$

FINDING A PIECE:  $Y = \{0, 1, 2, 3, 4, 5, 6\}$

$$X = \{0, 1\}$$

$$E[Y] = E[Y|X=0] PR(X=0) + \\ E[Y|X=1] PR(X=1)$$

$$E[Y|X=1] = 3 \quad PR(X=1) = \frac{1}{2}$$

$$E[Y|X=0] = 4 \quad PR(X=0) = \frac{1}{2}$$

$$(1 \times 0 + 2 \times \frac{1}{3} + 3 \times 0 + 4 \times \frac{1}{3} + 5 \times 0 + 6 \times \frac{1}{3})$$

$$= 4$$

$$E[Y] = 3 \times \frac{1}{2} + 4 \times \frac{1}{3} = 3.5$$

SHORT NOTATION FOR LIE:

$$E[Y] = E[E[Y|X]]$$

INDEPENDENCE:  $PR(Y=y|X=x) = PR(Y=y)$

KNOWING X CARRIES NO  
INFO ON THE LIKELIHOOD

OF  $Y=y$

$$PR(Y=y|X=x) = \frac{PR(Y=y, X=x)}{PR(X=x)}$$

IF INDEPENDENT  
 $\downarrow PR(Y=y)$

$$PR(Y=y, X=x) = PR(Y=y) \times PR(X=x)$$

ONLY HOLDS IF X & Y ARE  
INDEPENDENT

CORRELATION:

$$\text{COR}(X, Y) = \overline{D_{XY}} = E[(X - \mu_x)(Y - \mu_y)]$$

$$= \sum_i I_j (X_i - \mu_x)(Y_i - \mu_y) PR(X=x_i, Y=y_i)$$

ONLY DIRECTION MATTERS NOT THE SCALE  $\angle \theta = 0, \pi/2$

$$\text{COR}(X,Y) = \frac{\sigma_{XY}}{\sigma_X \sigma_Y} \quad -1 \leq \text{COR}(X,Y) \leq 1$$

$\downarrow$   
PERFECT  
NEGATIVE
 $\downarrow$   
PERFECT  
POSITIVE

**HANDOUT**  
**QUES NO 2b**

$$A = \{0, 1\}$$

$$M = \{0, 1, 2, 3, 4\}$$

$$E[M | A=0]$$

$$E[M | A=1]$$

$$E[M]$$

USE BAYES RULE  
TO GET

- $PR(M=m | A=0)$
  - $PR(M=m | A=1)$
  - $PR(A=0)$
  - $PR(A=1)$
- } NEED

$$PR(A=0) = PR(M=0, A=0) + PR(M=1, A=0) + \dots + PR(M=4, A=0)$$

$$PR(A=1) = PR(M=0, A=1) + PR(M=1, A=1) + \dots + PR(M=4, A=1)$$

$$PR(A=0) = 0.5$$

$$PR(A=1) = 0.5$$

$$PR(M=0 | A=0) = \frac{PR(M=0, A=0)}{PR(A=0)}$$

$$\begin{aligned} E[M|A=0] &= 0 \times PR(M=0|A=0) + \\ &\quad 1 \times PR(M=1|A=0) + \\ &\quad \dots + \\ &\quad 4 \times PR(M=4|A=0) \\ &= 0.56 \end{aligned}$$

$$E[M|A=1] = 0.14$$

$$\begin{aligned} E[M] &= E[M|A=0] PR(A=0) + \\ &\quad E[M|A=1] PR(A=1) \\ &= 0.56 \times \frac{1}{2} + 0.14 \times \frac{1}{2} \\ &= 0.35 \end{aligned}$$

09/20/2021

## RANDOM SAMPLING & DISTRIBUTION OF THE SAMPLE AVERAGE

POPULATION: IN MOST CASES DATA DO NOT EXIST.

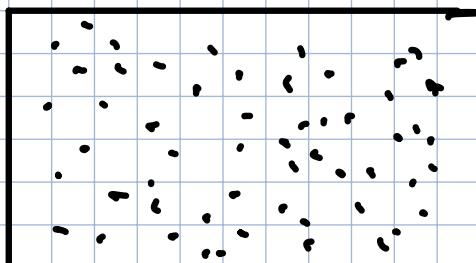
POPULATION PARAMETERS:

$E[Y]$ ,  $\mu_Y$ ,  $J(Y)$ ,  $\sigma_Y^2$ ,  $\sigma_Y$ ,  $Cov(Y, \cdot)$   
WR( $X_i \sim \eta$ ) ALL UNKNOWN

IN BUT WE WANT TO LEARN ABOUT THESE PARAMETERS.  
CONSUMPTION, PRICES, UNEMPLOYMENT, POLITICS,  
WELLBEING.

## RANDOM SAMPLING:

POPULATION FOR



Y  
RANDOM VARIABLE

n RANDOM OBS

$\{Y_1, Y_2, \dots, Y_n\}$

- EACH  $Y_i$  IS EQUALLY LIKELY TO BE DRAWN.
- EACH  $Y_i$  IS DRAWN FROM THE SAME RD.

iid → IDENTICALLY AND INDEPENDENTLY DISTRIBUTED

(2)

(1)

RANDOM SAMPLING INSURES THIS ↗

$$\bar{Y} = \sum_i Y_i / n$$

$\bar{Y}$  → RANDOM VARIABLE

DEF: SIMPLE AVERAGE OF n RANDOMLY DRAWN  $Y_i$

EVERY TIME, A RANDOM DRAWN n WILL GIVE YOU A DIFFERENT  $\bar{Y}$

$E[\bar{Y}]$ ,  $\text{Var}[\bar{Y}]$ , PROB. DISTRIBUTION

$$E[\bar{Y}] = E[Y_1 + Y_2 + \dots + Y_n] / n$$

$$\begin{aligned} E[\bar{Y}] &= 1/n (E[Y_1] + E[Y_2] + \dots + E[Y_n]) \\ &= 1/n (\mu_Y + \mu_Y + \dots + \mu_Y) \end{aligned}$$

$$E[\bar{Y}] = \mu_Y$$

$$\begin{aligned} \text{JAR}(\bar{Y}) &= \text{JAR}(1/n \sum_i Y_i) \\ &= \text{JAR}(1/n (Y_1 + Y_2 + \dots + Y_n)) \end{aligned}$$

BECAUSE  $Y_i, Y_j$  ARE INDEPENDENT

$$= \text{JAR}(1/n Y_1) + \text{JAR}(1/n Y_2) + \dots + \text{JAR}(1/n Y_n)$$

$$\text{JAR}(1/n Y_i) = 1/n^2 \text{JAR}(Y_i)$$

 PROOF IS IN THE APPENDIX

$$\sum_i (1/n Y_i - 1/n \mu_Y)^2 \rightarrow \text{follows THIS LOGIC}$$

$$\text{JAR}(\bar{Y}) = 1/n^2 \sum Y_i^2 \cdot n = \sum Y_i^2 / n = \sum \bar{Y}^2$$

$$\text{STD}\cdot\text{DEV}(\bar{Y}) = \sqrt{\sum Y_i^2 / n} = \sqrt{\sum \bar{Y}^2}$$

LARGE SAMPLE APPROXIMATIONS:

AS  $n \rightarrow \infty$

$$\begin{aligned} 2^{\circ} \text{ JAR}(Y_i) &= \text{Var}(Y_i) / n \\ 1. Y_1, Y_2, \dots, Y_n &\text{ iid} \\ \Rightarrow E[Y_i] &= \mu_Y \end{aligned}$$

LAW OF LARGE NUMBERS : UNDER GENERAL

CONDITIONS

Y IS A GOOD APPROXIMATION FOR

$\mu_Y$ .

$$\bar{Y} \xrightarrow{P} \mu_Y \text{ AS } n \rightarrow \infty$$

CENTRAL LIMIT THEOREM : AS  $n \rightarrow \infty$

$$\bar{Y} \sim N\left[\mathbb{E}[\bar{Y}], \frac{\sigma^2}{n}\right]$$

THE INITIAL PROBABILITY DISTRIBUTION  
OF Y DOES NOT MATTER.

-R SHOWCASE.

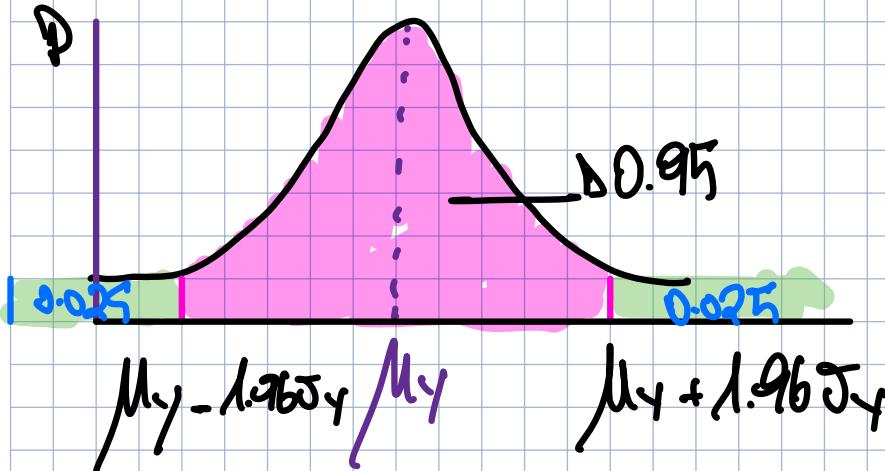
NORMAL DISTRIBUTION:

$$N(\mu_Y, \sigma_Y^2)$$



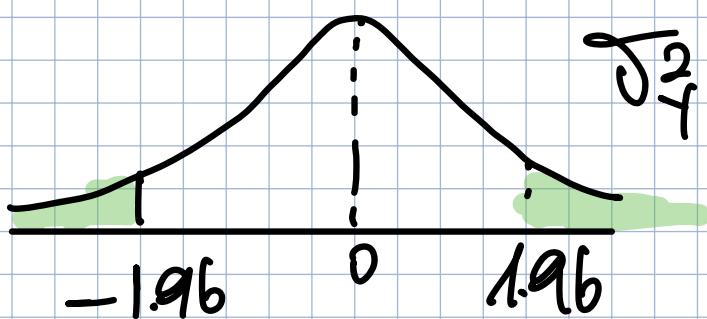
USUALLY BOTH

ARE UNKNOWN

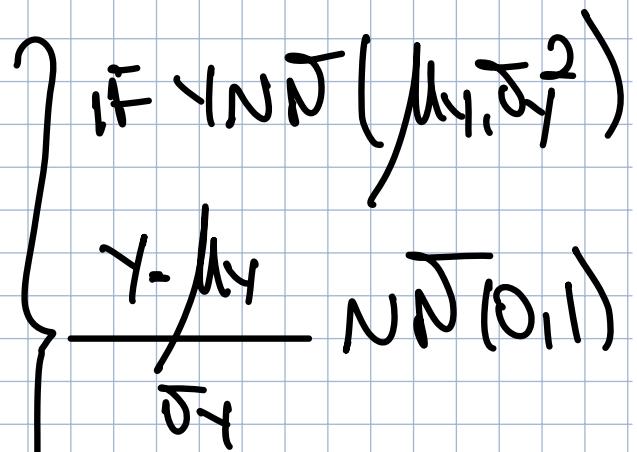


$$Z \sim N(0,1)$$

$$\mathbb{E}[Z] = 0$$



$$\sigma_Z^2 = 1 \quad \sigma_Y = 1$$



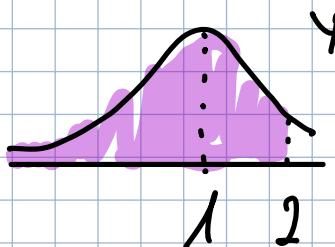
$$\Phi(c) = \Pr(Z \leq c)$$

$$\text{So } Y \sim N(1,4)$$

$$\ln \mu_Y$$

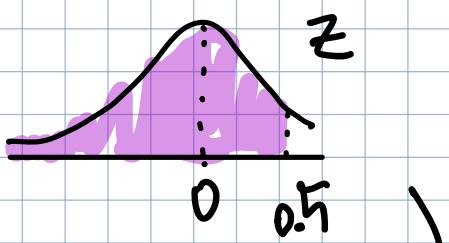
$$\frac{Y-1}{2} \sim N(0,1) \quad \text{OR} \quad Z \sim N(0,1)$$

$$\Pr(Y \leq 2) = \Pr\left(\frac{Y-1}{2} \leq \frac{2-1}{2}\right)$$



$$\Pr(Z \leq 1/2)$$

$$\Pr(Z \leq 0.5) = 0.691$$

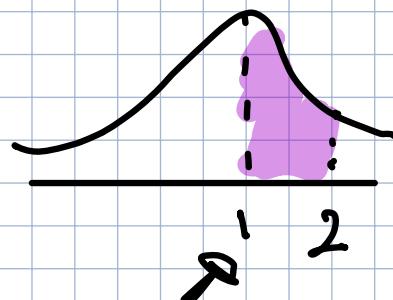


$$P(1 \leq Y \leq 2)$$

$$P(Y \leq 2) - P(Y \leq 1)$$

$$P(Z \leq 0.5) - P(Z \leq 0)$$

$$0.691 - 0.50 = 0.19$$



09/27/2021 REVIEW OF STATISTICS CH 3

ESTIMATOR: FUNCTIONS OF SAMPLE DATA DRAWN

↳ RANDOM VARIABLE  
(NUMERIC OUTCOME)

- MEAN
- MEDIAN

FROM AN UNKNOWN POPULATION

ESTIMATE: NUMERICAL VALUE OF THE ESTIMATOR

$\bar{x}$ : AVERAGE TRUCK EARNINGS OF COLLECTED GROUPS.

↳  $Y$  → RANDOM VARIABLE DEFINED AT THE POPULATION LEVEL

$y$  →  $\mu_y$  → POPULATION ATTRIBUTE IS WHAT WE WANT TO KNOW

RANDOM SAMPLE

$\{Y_1, Y_2, \dots, Y_n\}$

↳  $\bar{y}$  → MEAN  $\sum_i Y_i / n$  MEDIAN

IN GENERAL AN ESTIMATOR IS DENOTED AS  $\hat{M}_Y$

WHAT IS A GOOD ESTIMATOR?

1. UNBIASED  $E[\hat{M}_Y] = M_Y$

2<sup>o</sup> CONSISTENCY  $M_Y \xrightarrow{P} \hat{M}_Y$  AS  $n \rightarrow \infty$

3<sup>o</sup> JARVIANE  
CONSISTENCY

RANDOM

3<sup>o</sup> JARVIANE & EFFICIENCY

$\hat{M}_Y$  JS  $\tilde{M}_Y \rightarrow$  DIFFERENT ESTIMATORS

JAR( $\hat{M}_Y$ ) < JAR( $\tilde{M}_Y$ )

BEST  $\bar{Y}$  IS VALUE  $\rightarrow$  ESTIMATOR

↓ UNBIASED

LINEAR

$\sum_i Y_i / n$

$E[\bar{Y}] = M_Y$

$\bar{Y} \xrightarrow{P} M_Y$  AS  $n \rightarrow \infty$

BEST  $\rightarrow$  SMALLEST JARVIANE

$$\text{JAR}(\bar{Y}) = \text{JAR}\left(1/n \sum_i Y_i\right)$$

$$= 1/n^2 \text{JAR}(\sum_i Y_i)$$

$$\frac{1}{n^2} \left[ JAR(Y_1) + JAR(Y_2) + \dots + JAR(Y_N) \right]$$

$$\frac{1}{n^2} \cdot n JAR(Y_i) = \frac{1}{n} \bar{J}_Y^2 = JAR(\bar{Y})$$

$\sum_i (Y_i - m)^2 / n \rightarrow$  JAR CANCE TO MINIMIZE  
BY CHOOSING M

$$\frac{d \sum_i (Y_i - m)^2}{dm} = 0 \rightarrow \frac{d \sum_i (Y_i^2 - 2mY_i + m^2)}{dm}$$

$$\frac{d \sum_i Y_i^2}{dm} - \frac{d \sum_i 2mY_i}{dm} + \frac{dm^2}{dm} = 0$$

$$0 - 2 \sum_i Y_i + 2m = 0$$

PROOF:

$$-2 \sum_i (Y_i - m) = 0$$

$$\sum_i Y_i = nm \rightarrow m = \sum_i Y_i / n$$

$$m = \bar{Y}$$

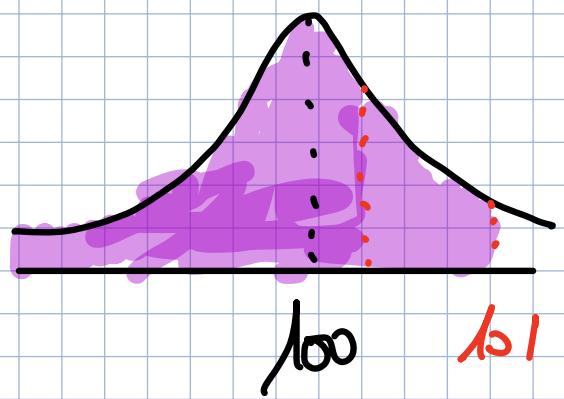
HANDOUT QUESTION 2<sup>o</sup>

$$CLT \bar{Y} \sim N(100, \bar{J}_Y^2 / n)$$

a.  $\mu_Y = 100 \quad \bar{J}_Y^2 = 43 \quad n = 100$

$$\bar{Y} \sim N\left(100, \frac{43}{100}\right)$$

$$\bar{Y} \sim N\left(100, 0.43\right)$$



$$P_R(\bar{Y} < 101)$$

$$\Leftrightarrow P_R\left(\frac{\bar{Y} - \mu_Y}{\sqrt{\sigma^2}} < \frac{101 - \mu_Y}{\sqrt{\sigma^2}}\right)$$

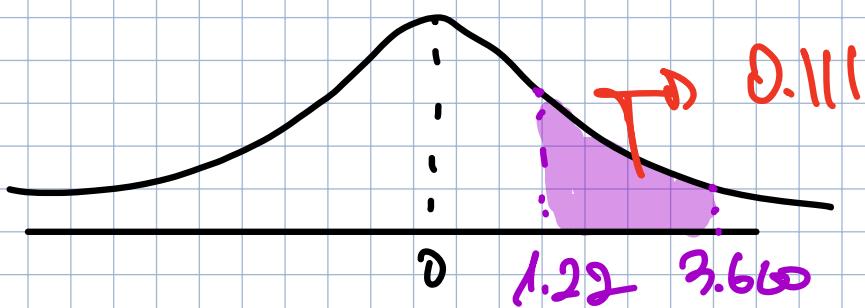
$$\Leftrightarrow P_R\left(Z < \frac{101 - 100}{\sqrt{0.43}}\right)$$

$$\Leftrightarrow P_R(Z < 1.525) = 0.9364$$

b.  $n = 64$   $\mu_Y = 100$   $\sigma^2_Y / n = 0.6719$

$$P_R(101 < \bar{Y} < 103) = P_R\left(\frac{101 - 100}{\sqrt{0.6719}} < Z < \frac{103 - 100}{\sqrt{0.6719}}\right)$$

$$\Leftrightarrow P_R(1.22 < Z < 3.660)$$



$$C. \quad n=165, \bar{Y} = \frac{\sum Y}{n} = \frac{43}{165} = 0.266$$

$$\Pr(\bar{Y} > 98) = 1 - \Pr(\bar{Y} \leq 98)$$

$$1 - \Pr\left(Z < \frac{98 - 100}{\sqrt{0.26}}\right)$$

$\approx 100\%$

## RANDOM SAMPLING + MEASUREMENT

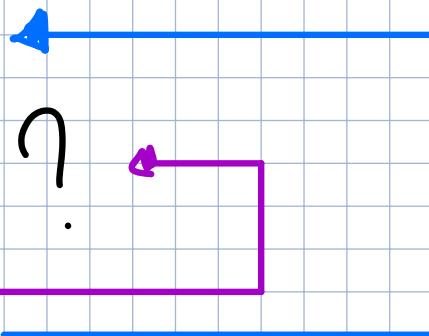
$y_1, y_2, \dots, y_n \rightarrow$  iid

NON-RANDOM SAMPLING:

UNEMPLOYMENT  $\rightarrow$  SUNDAY ?

CANCER RATE BY AGE ?

SAMPLE SELECTION BIASES:



## HYPOTHESIS TESTING CONCERNING THE POPULATION MEAN

1. SPECIFY THE HYPOTHESIS: WHAT DO YOU WANT

TO KNOW ABOUT THE POPULATION?

DO FACE MASKS WORK?

DO MASKS WORK?

IS THERE A GENDER / RACIAL GLASS CEILING?

Q<sup>0</sup> STATE THE NULL +<sub>0</sub> TRYING TO REJECT ALTERNATIVE +<sub>A</sub>

$$H_0: \mu_Y = \mu_{Y,0} \rightarrow \text{SPECIFIC VALUE}$$

↳ PRED MEAN

$$H_A: \mu_Y \neq \mu_{Y,0} \rightarrow \text{TWO-SIDED ALTERNATIVE}$$

IDEA: TAKE A RANDOM SAMPLE

$$\{Y_1, Y_2, \dots, Y_N\}$$

$$\mu_Y > \mu_{Y,0}$$

$$\mu_Y < \mu_{Y,0}$$

↳  $\sum_i Y_i / n = \bar{Y}$  → SAMPLE MEAN

NEVER BE EXACTLY EQUAL

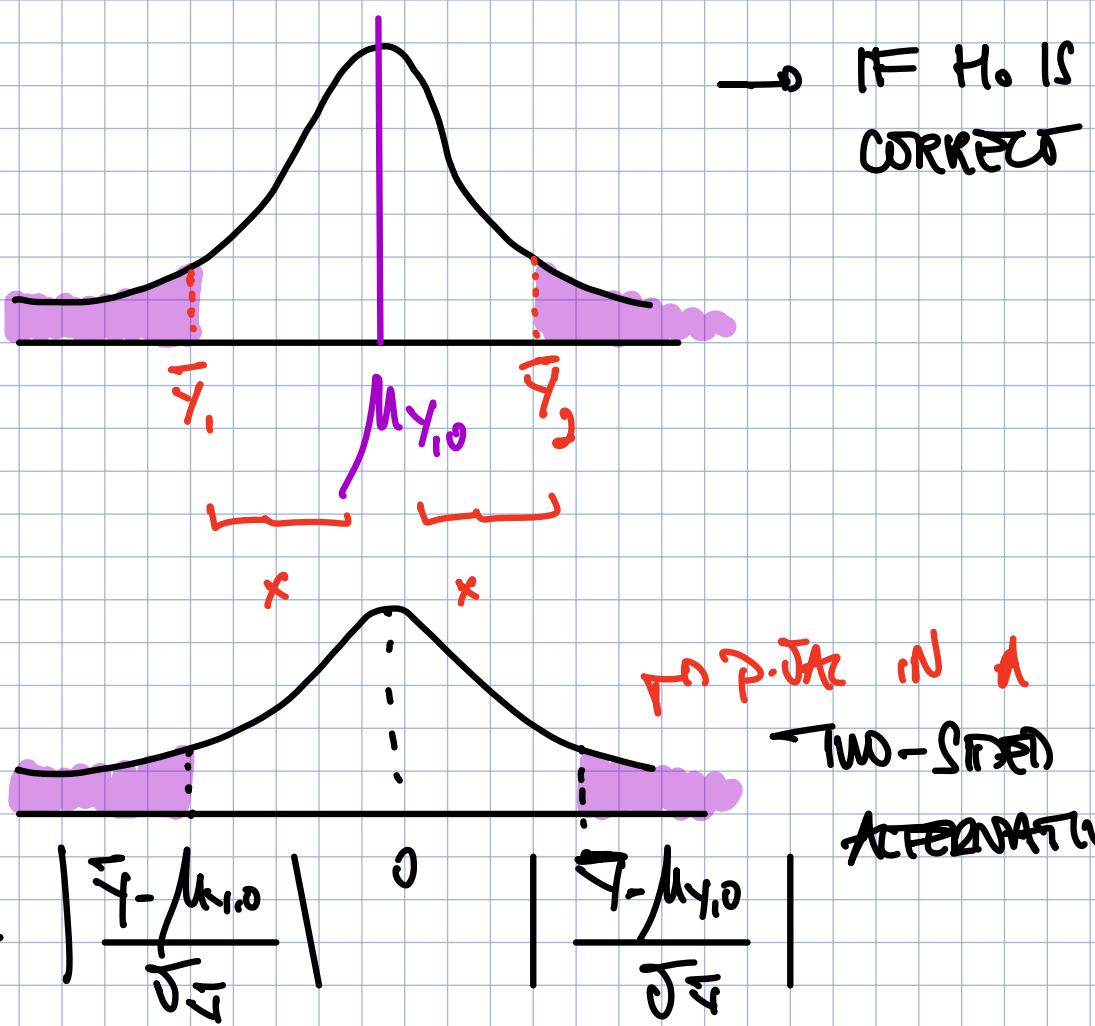
$$\bar{Y} = \mu_{Y,0}$$

↳ EVERY SAMPLE HAS A

1<sup>o</sup> ASSUME  $H_0$  IS TRUE: DIFF. MEAN

$$\bar{Y} \sim N(\mu_{Y,0}, \sigma^2_Y / n)$$

2<sup>o</sup> CRIT. P-VAL (PROB. THAT YOU DRAW A SAMPLE WHOSE MEAN IS AT LEAST AS EXTREME AS OBSERVED  $\bar{Y}$ )



FLIPPING A COIN →  $\mu_Y, \sigma_Y^2/n$   
 $\hookrightarrow \sigma_Y^2$

AVERAGE HOURLY EARNINGS       $\mu_Y, \sigma_Y^2 \rightarrow$  BOTH UNKNOWN  
 $n \rightarrow$  KNOWN

$$H_0: \mu_Y = \mu_{Y,0} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{You NEED}$$

$$H_A: \mu_Y \neq \mu_{Y,0} \quad \left. \begin{array}{l} \\ \end{array} \right\} \bar{Y} \times \sigma_Y^2/n$$

$\sigma_Y^2$  UNKNOWN

# THE SAMPLE VARIANCE, SAMPLE STANDARD DEV, AND SAMPLE STANDARD ERROR

$$S_y^2 = \frac{1}{n-1} \sum_i (y_i - \bar{y})^2 \xrightarrow{\text{REOXY}} S_y^2$$

$$S_y = \sqrt{S_y^2} \xrightarrow{\text{SAMPLE STD DEV}}$$

$$SE[\bar{y}] = \frac{S_y}{\sqrt{n}} \xrightarrow{\text{REOXY}} \frac{S_y}{\sqrt{n}}$$

$$P\text{-VAL} = 2 \Phi \left[ - \left| \frac{\bar{y} - \mu_{y,0}}{SE[\bar{y}]} \right| \right] \xrightarrow{\text{REOXY}} \frac{1}{\sqrt{n}}$$

$$t\text{-STAT} = \frac{\bar{y} - \mu_{y,0}}{SE[\bar{y}]} \xrightarrow{\text{REOXY}} P\text{-VAL} = 2 \Phi[-|t\text{-stat}|]$$

IF  $H_0$  IS TRUE  $\bar{y} \sim N(\mu_{y,0})$  THEN  
 $n$  IS LARGE

1. STATE  $H_0, H_A$   
2. CALCULATE TEST STATISTIC  $-t\text{-stat} \rightarrow P\text{-VAL}$

3<sup>rd</sup> COMPARE THE P-VAL AGAINST A PRE-SPECIFIED  
P THRESHOLD

(1) ECON: 1%, 5%, 10%

## 2 TYPES OF MISTAKES

- 1.  $H_0$  IS TRUE BUT YOU INCORRECTLY REJECT IT
  - 2.  $H_0$  IS FALSE BUT YOU INCORRECTLY ACCEPT IT OR FAIL TO REJECT IT. → HARD TO JUDGE  
DEPENDS ON THE DISTRIBUTION UNDER  $H_0$
- You CAN DETERMINE THE RISK THAT YOU WANT TO TAKE

### HANDOUT QUESTION:

a.  $Y = \{0, 1\}$   $E[Y] = \mu_Y = \sum_i Y_i P_i$

$$0 \times 0.5 + 1 \times 0.5 = 0.5 = E[Y]$$

b.  $V(Y) = \sum_i (Y_i - \mu_Y)^2 / 2$   
 $= [(1 - 0.5)^2 + (0 - 0.5)^2] / 2$   
 $\sigma^2_Y = 0.25 \quad \bar{\sigma}^2_Y = 0.25 / n$

c.  $\bar{\sigma}^2_Y = 0.25 / 100 = 0.0025$

$$\bar{Y} = 0.05$$

d.  $H_0: E[Y] = 0.5 \rightarrow$  COIN IS FAIR

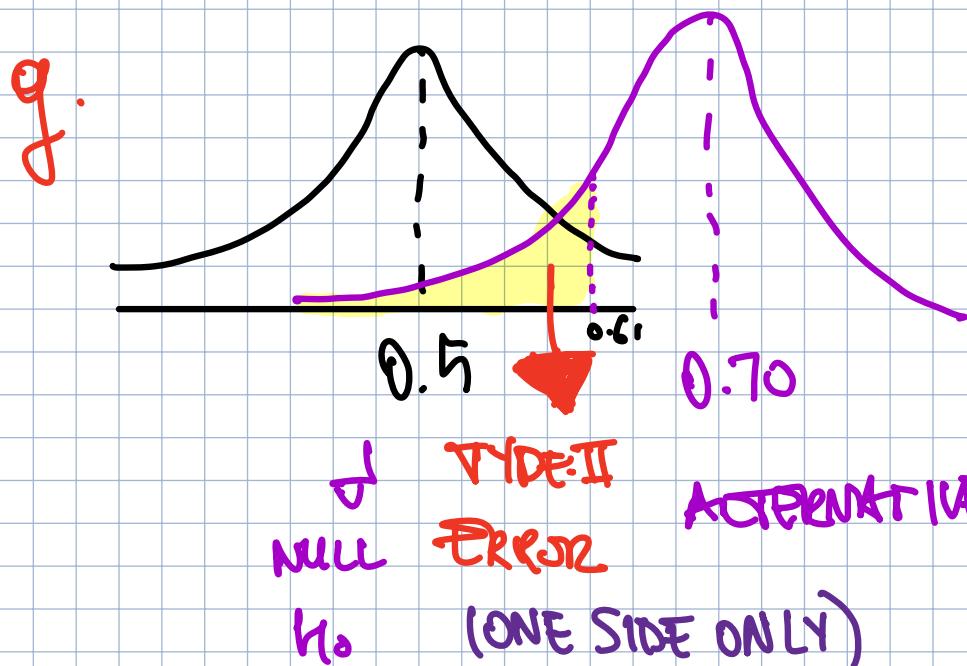
$H_1: E[Y] \neq 0.5 \rightarrow$  COIN IS BIASED

e.  $Z = \frac{\bar{Y} - \mu_{Y,0}}{\bar{\sigma}_Y} = \frac{0.61 - 0.5}{0.05} = 2.2$

$$P_{\text{JAL}} = 2 \Phi[-(Z\text{-score})]$$

$$P_{\text{JAL}} = 0.0278$$

f. 0.78% CHANCE TYPE-I ERROR



TYPE II ERROR

$$2 \times \Phi \left[ -\frac{0.61 - 0.70}{\sqrt{0.21/100}} \right] = 2 \times \Phi(-1.96)$$

UP TO 5%.

$$\begin{aligned} \text{JAR}(\bar{Y}) &= \frac{\sigma^2}{n} \\ &= 0.21/100 \end{aligned}$$

UP TO TYPE-II  
ERROR

$$\bar{Y} = 0.14 \quad \text{UP TO } 0.7$$

$$\begin{aligned} \text{JAR}((0-0.7)^2 \times 0.3 + (1-0.7)^2 \times 0.7) &= 0.21 \end{aligned}$$

04/11/2021

PETTER

LESSON PLAN: HYPOTHESIS TESTING CONCERNING  $\mu_y$

IF POP VARIANCE IS KNOWN  $\sigma_y^2 \rightarrow$

$$\text{TEST STATISTIC } Z\text{-SCORE} = \frac{\bar{Y} - \mu_{y,0}}{\sigma_y / \sqrt{n}}$$

$$\sigma_{\bar{Y}} = \sigma_y / \sqrt{n}$$

UNBIASED CASES

IF VARIANCE IS UNKNOWN  $\rightarrow$

$$t\text{-stat} = \frac{\bar{Y} - \mu_{y,0}}{s / \sqrt{n}}$$

$$s^2 = \frac{1}{n-1} \sum_i (Y_i - \bar{Y})^2$$

LOSAMPLE VALUES

$\rightarrow$  CONFIDENCE INTERVALS

$\rightarrow$  COMPARING MEANS FROM TWO

DIFFERENT POPULATIONS

? NEW MATERIAL

$$X, Y \quad H_0: \mu_x = \mu_y$$

$$\mu_x - \mu_y = d_0$$

$$H_A: \mu_x \neq \mu_y$$

$$\mu_x - \mu_y \neq d_0$$

$\rightarrow$  SCATTER PLOTS, SAMPLE COVARIANCE, SAMPLE CORRELATION

## TESTS ON SAMPLE MEAN

1. STATE  $H_0, H_A$        $H_0: \mu_y = 0.5$        $H_A: \mu_y \neq 0.5$

2. USE CCT TO PREDICT THE DIST. UNDER  $H_0$  IS TRUE

$$\bar{Y} \sim N(\mu_y, \sigma^2_y/n) \rightarrow \bar{Y} \sim N(0.5, 0.0025)$$

3. CALCULATE THE ACTUAL  $\bar{Y} = \sum Y_i/n = 0.61$

4. IF VARIANCE  $\sigma^2_y/n$  IS KNOWN:

$$Z = \frac{\bar{Y} - \mu_{y,0}}{\sqrt{\sigma^2_y}} = \frac{0.61 - 0.5}{\sqrt{0.0025}} = 2.2$$

5. CALCULATE P-VAL

$$2 \times \Phi[-|Z\text{-SCORE}|] = 2.78\%$$

6. COMPARE P-VAL TO SIGNIFICANCE LEVEL

1%, 5%, 10% IF P-VAL < SIGNIFICANCE LEVEL

↳ REJECT  $H_0$

7. IF  $\sigma^2_y/n$  IS UNKNOWN?

USE

$\rightarrow$  SAMPLE VARIANCE

$$S_y^2 \rightarrow \hat{\sigma}_y^2$$

$$\bar{Y} = 0.61$$

$$S_y^2 = \frac{\sum (Y_i - \bar{Y})^2}{(n-1)}$$

$$S_y / \sqrt{n} \leftrightarrow SE[\bar{Y}] \rightarrow \hat{\sigma}_{\bar{Y}}$$

$$\hookrightarrow S_y \approx 0.49 \quad SE[\bar{Y}] = 0.0019$$

$$5. t\text{-stat} = \frac{\bar{Y} - \mu_{Y,0}}{SE(\bar{Y})} = \frac{0.61 - 0.5}{0.049} \approx 2.244$$

$$\hookrightarrow P\text{-TAUZ} = 2 \times \left[ -|t\text{-stat}| \right] \sim 2 \times \Phi(-|t\text{-stat}|)$$

IF  $n$  IS JAROB TYPE II ERROR : H<sub>0</sub> IS FALSE

$$P\text{-JAR} = 0.0270$$

YOU FAIL TO REJECT

YOU NEED AN ALTERNATIVE  
THAT IS MORE PRECISE

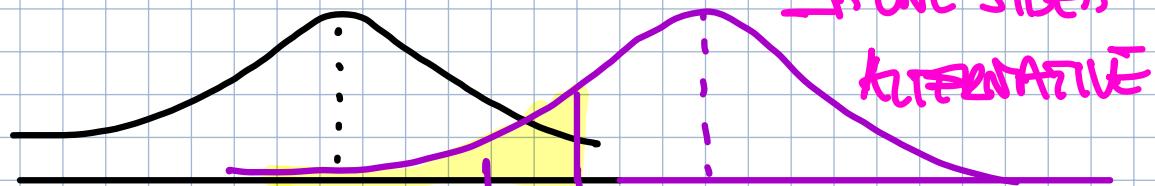
**HA:**  $\mu_Y = 0.7$   $JAR(Y) = \sigma_Y^2$

$$JAR(Y) = \sum_i (Y_i - \mu_Y)^2 \Pr(Y=Y_i)$$

$$= (0-0.7)^2 \times 0.3 + (1-0.7)^2 \times 0.70 \\ = 0.49 + 0.06 = 0.21$$

$$\sigma_{\bar{Y}}^2 = \sigma_Y^2/n = 0.21/100 = 0.0021$$

TYPE II  
ERROR :



$$\Phi \left[ - \left| \frac{0.61 - 0.7}{\sqrt{0.0021}} \right| \right]$$

TYPE-II  
ERROR

$$\Phi(-1.96) \approx 2.5\%$$

## CONFIDENCE INTERVALS:

$$90\%: \bar{Y} \pm 1.65 \times \text{SE}[\bar{Y}]$$

$$95\%: \bar{Y} \pm 1.96 \times \text{SE}[\bar{Y}]$$

$$99\%: \bar{Y} \pm 2.576 \times \text{SE}[\bar{Y}]$$

$$0.61 \pm 1.96 \times 0.049 \rightarrow [0.51, 0.7]$$

**INTERPRETATION:** 95% CONFIDENT THAT THE TRUE POPULATION MEAN IS WITHIN THE CI

## COMPARING MEANS FROM TWO POPULATIONS

$\mu_w \rightarrow$  POPULATION AVE FOR WOMEN

$\mu_m \rightarrow$  POPULATION AVE FOR MEN

$$H_0: \mu_m - \mu_w = d_0 \rightarrow \text{most often } d_0 = 0$$

$$H_A: \mu_m - \mu_w \neq d$$

But does not have to

RANDOM SAMPLE OF MEN  $\rightarrow \bar{Y}_m, n_m$

RANDOM SAMPLE OF WOMEN  $\rightarrow \bar{Y}_w, n_w$

$$\bar{Y}_m - \bar{Y}_w \rightarrow \text{Best Estimate for } \mu_m - \mu_w$$

LARGER THE INTERVAL GETS MORE CERTAIN THE ESTIMATE IS

$$\bar{Y}_{mN} \sim N(\mu_m, \frac{\sigma_m^2}{n_m}/\text{hm})$$

$$\bar{Y}_{wN} \sim N(\mu_w, \frac{\sigma_w^2}{n_w}/\text{hw})$$

$\hookrightarrow \bar{Y}_m - \bar{Y}_w \sim N(\mu_m - \mu_w, \frac{\sigma_m^2}{n_m} + \frac{\sigma_w^2}{n_w})$

$\sigma_m^2, \sigma_w^2 \rightarrow$  TYPICALLY UNKNOWN

You NEED TO ESTIMATE  $s_w^2, s_m^2 \rightarrow$  SAMPLE VARIANCE

$$SE[\bar{Y}_m - \bar{Y}_w] = \sqrt{\frac{s_m^2}{n_m} + \frac{s_w^2}{n_w}}$$

IF THE  $H_0$  IS TRUE &  $n$  IS LARGE

$$t\text{-stat} = \frac{(\bar{Y}_m - \bar{Y}_w) - d_0}{SE[\bar{Y}_m - \bar{Y}_w]} \sim N(0, 1)$$

$$P\text{-value} = 2 \Phi(-|t\text{-stat}|) / 2 + (-|t\text{-stat}|)$$

$\rightarrow$  SKIP CASE WHEN  $n < 30$

## TERMINOLOGY:

- TYPE I ERROR:  $H_0$  IS CORRECT BUT YOU FALSELY REJECT IT
- TYPE II ERROR:  $H_0$  IS INCORRECT BUT YOU FAIL TO REJECT IT
- SIGNIFICANCE LEVEL: PRESPECIFIED PROBABILITY OF TYPE-I ERROR
- THE CRITICAL VALUE: VALUE OF THE TEST STATISTIC FOR WHICH THE TEST REJECTS THE NULL
- REJECTION ZONE: AREA UNDER WHICH THE RESEARCHER REJECTS  $H_0$
- SIZE OF THE TEST: PROB. OF INCORRECTLY REJECTING  $H_0$
- POWER OF THE TEST: PROB. OF CORRECTLY REJECTING  $H_0$  WHEN  $H_1$  IS TRUE

		DECISION	
		FAIL TO REJECT	REJECT
$H_0$	TRUE	TRUE POSITIVE (1 - $\alpha$ )	FALSE POSITIVE ( $\alpha$ ) TYPE I ERROR
	FALSE	FALSE NEGATIVE TYPE II ERROR	TRUE NEGATIVE POWER OF A TEST (1 - $\beta$ )

10/07/2021

CONFIDENCE INTERVALS:  $(\bar{Y}_M - \bar{Y}_W) \pm 1.96 \text{SE}[\bar{Y}_M - \bar{Y}_W]$

a.  $Y = \xi_{0,1}^U$        $P(Y=0) = P$   
           $\leftarrow$  (INDEPENDENT)       $P(Y=1) = 1-P$

EMR CANCER

$$\bar{Y} = 215 / 400 = 0.5373$$

b.  $S_Y^2 = \sum_i (Y_i - \bar{Y})^2 / n-1$

$$S_Y^2 = \frac{1}{399} \left[ (1 - 0.5373)^2 \times 215 + (0 - 0.5373)^2 \times 185 \right]$$

$$= S_Y^2 \approx 0.249 \quad \text{SE}[\bar{Y}] = S_Y / \sqrt{n}$$

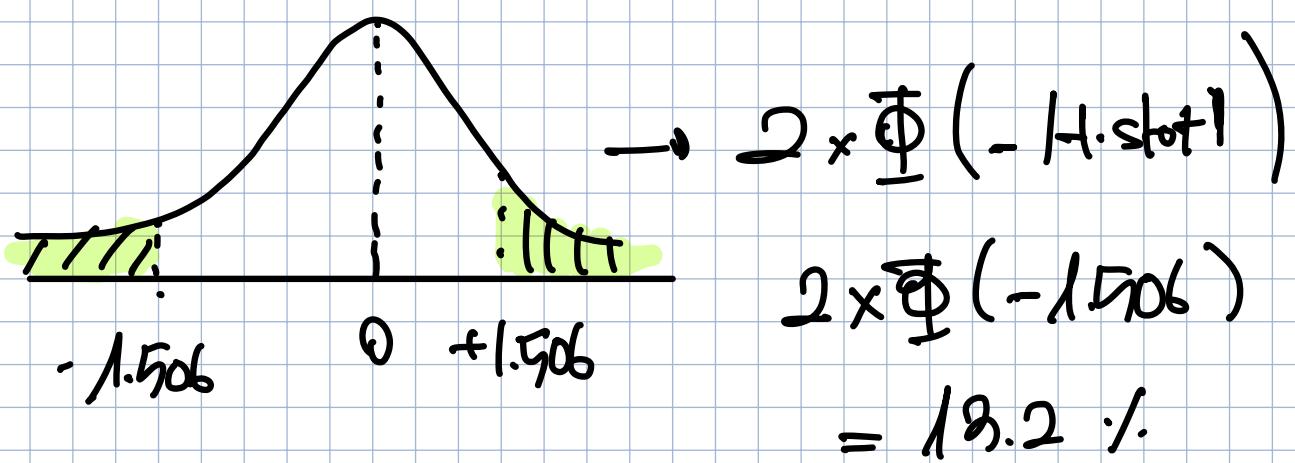
$$S_Y \approx 0.499 \quad \approx 0.0249$$

c.  $H_0: P = 0.5$

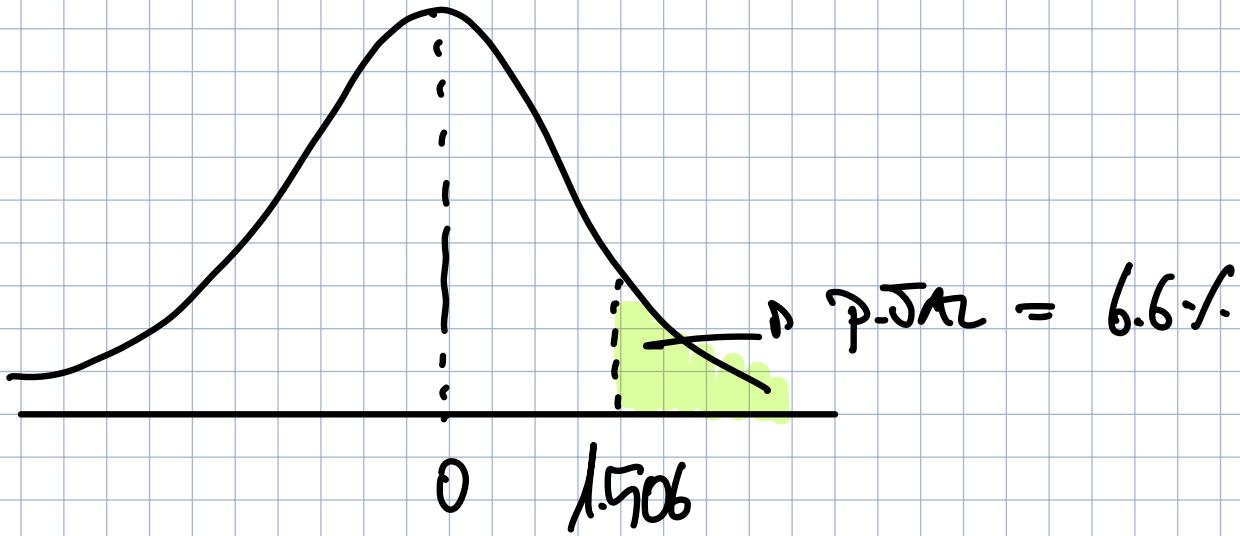
$$H_A: P \neq 0.5$$

$$-z_{\text{shf}} = \frac{\bar{Y} - \mu_{Y_0}}{\text{SE}[\bar{Y}]}$$

$$= \frac{0.5373 - 0.5}{0.0249}$$
$$= 1.506$$



d.



e.  $P\text{-JAR} > 0.05$  FAIL TO REJECT THE NULL

SURVEY DOES NOT CONTRADICT EVIDENCE

f. 95% CI  $\bar{Y} \pm 1.96 \times SE[\bar{Y}]$

$$[0.481 - 0.579]$$

95% OF CHANCE THAT TRUE  $\mu_y$  IS WITHIN

THE 95% CI  $\rightarrow$  MORE CONSERVATIVE

g.  $\bar{Y} \pm 2.576 \times SE[\bar{Y}]$   $[0.463 - 0.507]$

Quesion 2.  $\bar{N} = 420 \quad \bar{Y} = 646.2 \quad S_y = 19.5$

①.  $SE[\bar{Y}] = S_y / \sqrt{n} \approx 0.95$

$$\bar{Y} \pm 1.96 \times SE[\bar{Y}] \quad [644.34 \quad - 648.06]$$

②.  $\bar{Y}_1 - \bar{Y}_2 = 657.4 - 650$

$$= 7.4$$

$$SE[\bar{Y}_1 - \bar{Y}_2] = \sqrt{\frac{S_{Y_1}^2}{n_1} + \frac{S_{Y_2}^2}{n_2}}$$

$$= \sqrt{\frac{19.4}{238} + \frac{17.9}{182}}$$

$$= 1.8281$$

$$[7.4 \pm 1.96 \times 1.8281] = [3.82 \quad 10.98]$$

$H_0 : DIF = 0$

$$t\text{-stat} = \frac{7.4 - 0}{1.8281}$$

reject the null

$$= 4.0479$$

DISTRIBUTED WITH SMALLER CLASSES HAVE BETTER OUTCOMES

Q Question :  $\bar{Y}_1 - \bar{Y}_2 = -253.2284$

$$\bar{Y}_1 = 3178.832$$

$$S_{Y_1} = 580.0068$$

$$\bar{Y}_2 = 3432.06$$

$$S_{Y_2} = 584.622$$

$$SE[\bar{Y}_1 - \bar{Y}_2] = \sqrt{\frac{S_{Y_1}^2}{n_1} + \frac{S_{Y_2}^2}{n_2}}$$

$$= 28.82106$$

$$f.\text{stat} = \frac{-253.2284}{28.82106}$$

$$= -9.44$$

$$\rightarrow P\text{-VAL} = 0$$

## LINEAR REGRESSION

SCATTER PLOT : X: Y: n OBS PUTTING DATA

$$S_{XY} = \frac{1}{n-1} \sum_i (x_i - \bar{x})(y_i - \bar{y})$$

$\hookrightarrow$  SAMPLE COVARIANCE

$\hookrightarrow$  SAMPLE VARIANCE

$$S_x^2 = \sum_i \frac{(x_i - \bar{x})^2}{(n-1)}$$

$$r_{xy} = \frac{S_{xy}}{S_x S_y} \rightarrow \text{SAMPLE CORRELATION}$$

$-1 \leq r_{xy} \leq +1 \rightarrow$  TELLS YOU HOW MUCH  
X, Y ARE RELATED

P → MEASURE OF LINEAR ASSOCIATION  
 X Y → (DO NOT IMPLY CAUSAL INFERENCE)

↳ SHOW STRONG ASSOCIATION B/W X & Y

PREDICTION NOTATION:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \rightarrow \text{PREDICTION FUNCTION}$$

Y → OUTCOME VARIABLE  
 RESPONSE VARIABLE  
 LEFT HANDSIDE VARIABLE  
 INDEPENDENT VARIABLE

X → INDEPENDENT VAR.  
 RIGHT HANDSIDE  
 PREDICTION

$\beta_0, \beta_1 \rightarrow$  POPULATION COEFFICIENTS / PARAMETERS

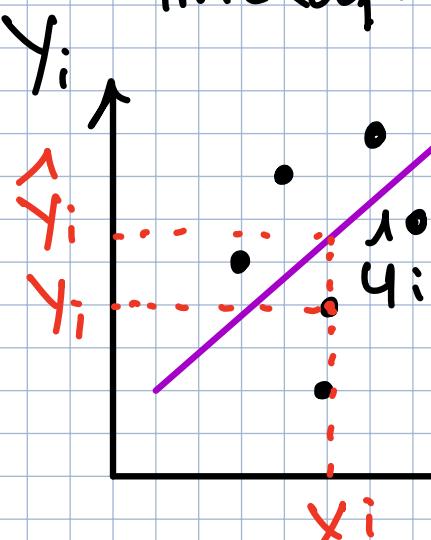
$$i = \{1, 2, \dots, n\}$$

$\beta_0 + \beta_1 X_i \rightarrow$  POPULATION PREDICTION LINE

$\beta_0 \rightarrow$  SLOPE

$\beta_1 \rightarrow$  INTERCEPT

$\epsilon_i \rightarrow$  ERROR TERM



$\rightarrow \beta_0 + \beta_1 X_i$   
 $Y_i \rightarrow$  OBSERVED  
 $\hat{Y}_i \rightarrow$  PREDICTED

$$\hat{Y}_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$X_i$	$\hat{Y}_i$	$M$
5	7	1
6	10	2
12	12	3
⋮	⋮	⋮

## Two Main Problems:

I. WE DON'T OBSERVE THE DISTRIBUTION

LIBRARY CAN USE RANDOM SAMPLES

2. Which line to fix? → First focus on this

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \rightarrow \text{LEFTOVER}$$

A hand-drawn diagram on lined paper. At the top right, the word "EDUCATION" is written in large, bold, black capital letters. A bracket is drawn under the letter "T". To the left of this bracket, the word "PARENTAL INPUT" is written in large, bold, black capital letters. Below these two concepts, the word "SOCIAL NETWORKS" is written in large, bold, black capital letters. A bracket is drawn under the letter "S". To the left of this bracket, the word "EARNINGS" is written in large, bold, black capital letters. The entire diagram is drawn in black ink on white paper.

$$y_i = \beta_0 + \beta_1 x_i + u_i$$

$$\begin{array}{r} \text{id} \\ \hline 1 \\ x_1 \\ \vdots \\ n \end{array} \quad \begin{array}{r} X \\ \hline x_1 \\ y_2 \\ \vdots \\ x_n \end{array} \quad \begin{array}{r} Y \\ \hline y_1 \\ y_2 \\ \vdots \\ y_n \end{array} \quad \begin{array}{r} Z \\ \hline z_1 \\ z_2 \\ \vdots \\ z_n \end{array}$$

$$Y_i = \beta_0 + \beta_1 X_i$$

$\hookrightarrow$  PREDICTED VALUE     $\hookrightarrow$  OBSERVED VALUE

$$Y_i = Y_i + C_i \rightarrow \text{PREDICTED}$$

$\downarrow$   
D PREDICATED VALUE

$U_i \rightarrow$  ERROR TERM

$y_i \rightarrow$  PREDICTED VALUE

PROBLEM: CHOOSE  $\hat{\beta}_0, \hat{\beta}_1 \rightarrow \hat{\beta}_0 + \hat{\beta}_1 X_i \rightarrow Y_i$

How?

$$\text{MINIMUM PREDICTION ERROR} = \min \sum_i \hat{u}_i^2$$

$$\sum_i \hat{u}_i^2 = \hat{e}_1^2 + \hat{e}_2^2 + \dots + \hat{e}_n^2$$

Sum of the squared residuals

$$\min \sum_i \hat{u}_i = \min \sum_i (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

$$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{u}_i$$

$$\hat{u}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i$$

$$\hat{u}_i^2 = (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

$$\hat{\beta}_1 = \frac{\sum_i (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_i (X_i - \bar{X})^2}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

$$\left. \begin{aligned} \sum_i \hat{u}_i^2 &= \sum_i (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2 \end{aligned} \right\}$$

$$= \frac{s_{xy}}{s_x^2}$$

WE CAN DEBATELY  
ESTIMATE FROM  
THE SAMPLE

# 10/30/2021 LECTURE PLAN

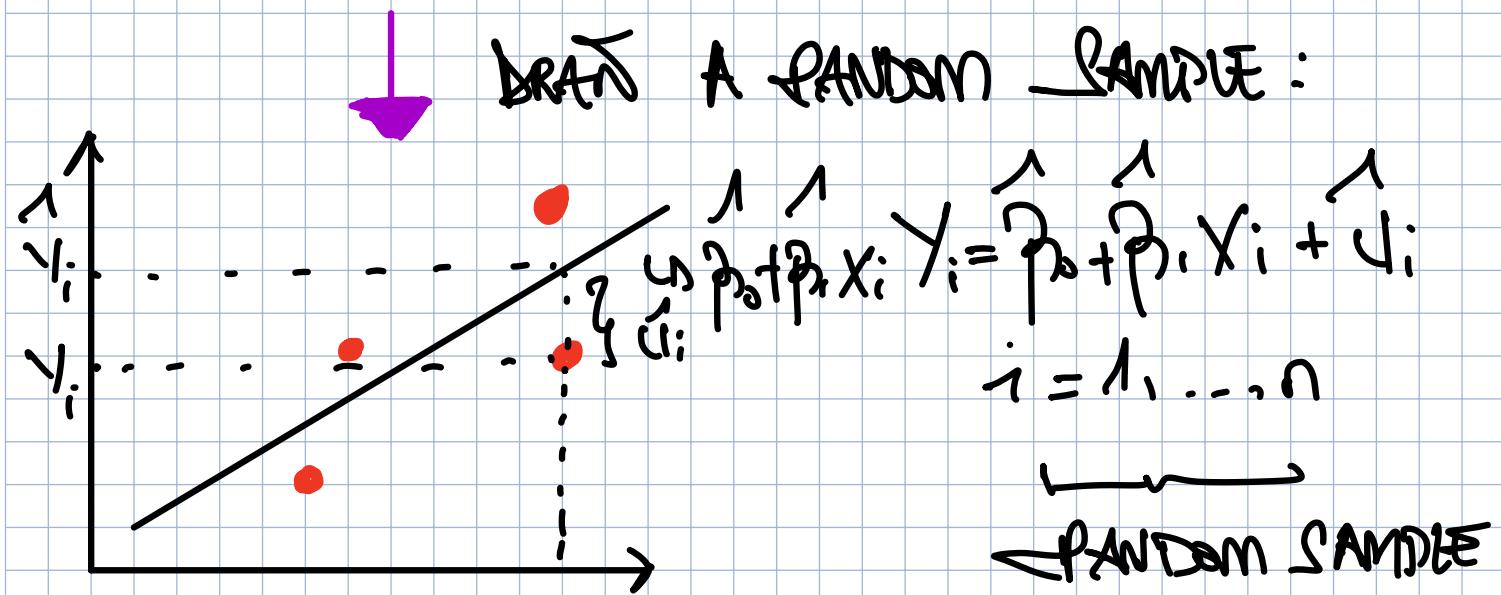
- REVIEW OF OLS → ORDINARY LEAST SQUARES
- MEASURES FOR GOODNESS OF FIT
- ONE PLS EXAMPLE

$\hat{Y}_i$  → PREDICTION ERROR

$$Y_i = \beta_0 + \beta_1 X_i + U_i \rightarrow \text{REGRESSION FUNCTION}$$

PREDICTS A LINEAR & IMPERFECT  $\beta_0 + \beta_1 X_i + U_i$

$\beta_0 \rightarrow$  INTERCEPT  
 $\beta_1 \rightarrow$  SLOPE } REGRESSION PARAMETERS  
 $Y$  IS A LINEAR FUNCTION OF  $X$  }  $\beta_0, \beta_1$ . ARE UNKNOWN



$$Y_i = \hat{Y}_i + U_i$$

UNOBSERVED

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i \rightarrow \text{OBSERVED}$$

L → PREDICTED VALUE

$\sum_i u_i^2 \rightarrow \text{SSR} \rightarrow \text{MIN}$

$$\hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} = \frac{\text{S}_{xy}}{\text{S}_{xx}}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \rightarrow \text{EVERYTHING IS OBSERVED}$$

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + u_i$$

WHAT DO THESE ESTIMATED PARAMETERS  
 $\hat{\beta}_0, \hat{\beta}_1$  TELL US?

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + u_i \quad 1^\circ$$

$$y_i + \Delta y = \hat{\beta}_0 + \hat{\beta}_1 (x_i + \Delta x) + u_i \quad 2^\circ$$

$$\Delta y = \hat{\beta}_1 \Delta x \rightarrow \hat{\beta}_1 = \Delta y / \Delta x$$

$$\hat{\beta}_1 = \Delta y \quad \text{WHEN } \Delta x = 1$$

CHANGE IN  $\bar{Y}$  ASSOCIATED WITH ONE  
UNIT INCREASE IN  $\bar{X}$

$$\hat{Y}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

$\hat{Y}_0 = \bar{Y}$  WHEN  $\bar{X} = 0$  WHY?

$$\bar{X} = \sum_i X_i / n = 0 \quad \sum_i X_i = 0$$

$\hat{Y}_0$  → SOMETIMES IT MAKES SENSE  
SOMETIMES IT DOES NOT

HANDBOUT ANSWERS

$\hat{Y} \rightarrow$  YOU CAN PREDICT  
 $\bar{Y}$  FOR ANY

JACQUE

$\hat{Y}_0$  HAS NO MEANING

$\hat{\beta}_1$  → AN INCREASE IN THE STR BY ONE

UNIT IS ASSOCIATED WITH A 2.83 POINTS

$$\text{IF } X_i = 24 \quad \hat{Y}_i = \hat{Y}_0 + \hat{\beta}_1 X_i$$

$$= 695.35 - 2.83 \cdot 24 \\ = 627.38$$

IS THE ASSOCIATION LARGE OR SMALL?

SUMMARISE THE OUTCOME:  $\bar{Y} = 628.8$

$$2.83 / 695.3 = 0.4\%$$

$$2.83 / 23.748 \approx 0.12 \text{ SD}$$

MEASURES OF FIT

$R^2$  - STD. ERROR

OF THE

PREDICTION

$$R^2: 0 \leq R^2 \leq 1$$

US SHARE OF VARIATION IN Y EXPLAINED

BY THE VARIATION IN X.

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + U_i$$

$$Y_i = \hat{Y}_i + U_i$$

UTSS: TOTAL SUM OF SQUARES

$$\text{UTSS} = \sum_i (Y_i - \bar{Y})^2 \rightarrow \text{TOTAL VARIATION IN Y}$$

RSS : RESIDUAL SUM OF SQUARES

$$RSS = \sum_i u_i^2$$

ESS : EXPLAINED SUM OF SQUARES

$$ESS = TSS - RSS$$

$$ESS = \sum_i (\hat{y}_i - \bar{y})^2 \rightarrow \begin{matrix} \text{TOTAL VARIATION} \\ \text{IN } y \end{matrix}$$

$$R^2 = 0 \leq R^2 \leq 1$$

$$R^2 = ESS / TSS = 1 - RSS / TSS$$

FUN FACTS ABOUT RSS

1. SAMPLE AVERAGE OF THE RSS RESIDUALS = 0

$$\hat{u}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\text{L} \quad \hat{u}_i = y_i - \bar{y} + \hat{\beta}_1 \bar{x} - \hat{\beta}_1 x_i$$

$$\hat{u}_i = (y_i - \bar{y}) - \hat{\beta}_1 (x_i - \bar{x})$$

$$\sum_i \hat{u}_i = \sum_i (y_i - \bar{y}) - \hat{\beta}_1 \sum_i (x_i - \bar{x}) \\ = 0 \quad \text{LDO} \quad \text{L}$$

$$\sum_i (Y_i - \bar{Y}) = \sum_i Y_i - \sum_i \bar{Y} \\ = n \bar{Y} - n \bar{Y} = 0$$

$$\sum_i Y_i = Y_1 + Y_2 + \dots + Y_n \\ \overbrace{\qquad\qquad\qquad}^n \times \bar{n}$$

$$\sum_i Y_i = \bar{Y} \cdot n \quad \text{IF } \sum_i U_i = 0$$

$$1/n \sum_i U_i = 0$$

$$2^o \quad 1/n \sum_i \hat{Y}_i = \bar{Y}$$

$$Y_i = \hat{Y}_i + U_i \rightarrow \sum_i Y_i = \sum_i \hat{Y}_i + \sum_i U_i$$

$$\sum_i Y_i = \sum_i \hat{Y}_i + 0 \\ \bar{n} \bar{Y} = \sum_i \hat{Y}_i$$

$$\sum_i \hat{Y}_i / \bar{n} = \bar{Y} \quad \xrightarrow{\text{OMITTED}}$$

$$3^o \quad \sum_i \hat{U}_i X_i = \sum_i \hat{U}_i (X_i - \bar{X}) = 0$$

$$\sum_i \hat{U}_i X_i = \sum_i [(Y_i - \bar{Y}) - \hat{U}_i (X_i - \bar{X})] (X_i - \bar{X})$$

$$\Leftrightarrow \sum_i (Y_i - \bar{Y})(X_i - \bar{X}) - \hat{U}_i \sum_i (X_i - \bar{X})^2 = 0$$

**REMEMBER**  $\hat{y}_i = \frac{\sum_i (y_i - \bar{y})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2}$

$$\begin{aligned} \text{SSTSS} &= \sum_i (Y_i - \bar{Y})^2 = \sum_i (Y_i - \hat{Y}_i + \hat{Y}_i - \bar{Y})^2 \\ &= \sum_i (Y_i - \hat{Y}_i)^2 + \sum_i (\hat{Y}_i - \bar{Y})^2 + 2 \sum_i (Y_i - \hat{Y}_i)(\hat{Y}_i - \bar{Y}) \\ &\Leftrightarrow \text{SSR} + \text{ESS} + 2 \underbrace{\sum_i u_i \hat{Y}_i}_{\text{SSE}} \end{aligned}$$

$$\begin{aligned}\sum_i \hat{u}_i y_i &= \sum_i \hat{u}_i (\hat{\beta}_0 + \hat{\beta}_1 x_i) \\ &= \hat{\beta}_0 \sum_i \hat{u}_i + \hat{\beta}_1 \sum_i \hat{u}_i x_i\end{aligned}$$

2

$\hat{\beta}_0, \hat{\beta}_1 \rightarrow$  FIND THESE GIVE THAT  
BEST FIT THE DATA!

WHAT DOES THIS MEAN?

LEAST SUM OF THE SQUARED MISTAKES  
WHEN YOU USE  $X$  TO PREDICT  $Y$ .

ANSWER TO THE QUESTION

$$1. \hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} = \frac{s_{xy}}{s_x^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = -150.78 / 63.83$$

$$= -2.83$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\begin{aligned} \hat{\beta}_0 &= 628.8 - (-2.83 \times 23.5) \\ &= 628.8 + 66.55 \\ &= 695.35 \end{aligned}$$

$$1 \quad 1 \quad 1 \\ y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i, i = 1, \dots, n$$

$$1 \quad 1 \\ u_i = y_i - \hat{y}_i$$

$$1 \quad 1 \quad 1 \\ u_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \quad \text{OLS} \quad \sum_i u_i^2$$

$$\hat{\beta}_0, \hat{\beta}_1 \rightarrow \hat{\beta}_0, \hat{\beta}_1 \quad u_i = 0$$

SAMPLE PARAMETERS RESIDUAL  $\leftarrow$

WEERKTAF

$$\widehat{\text{TEST SCORE}} = 695.3 - 2.83 \times \text{STR}$$

WHAT MEANS  
PREDICTED

IN STUDENT  
TO TEACHER

1. You can do prediction by plugging <sup>RATIO</sup> STR assuming that the relationship is linear.
2. Use coefficients to interpret the relationship.
  - $\beta_1$  has no meaning
  - $\beta_1$  — an increase in the student to teacher ratio by 1 unit is associated with a 2.83 points

$$\text{IF } X_i = 24 \quad \widehat{Y}_i = \beta_{0i} + \beta_{1i} X_i \\ = 695.3 - 2.83 \times 24 \\ = 627.38$$

IS THE ESTIMATE SMALL / LARGE?

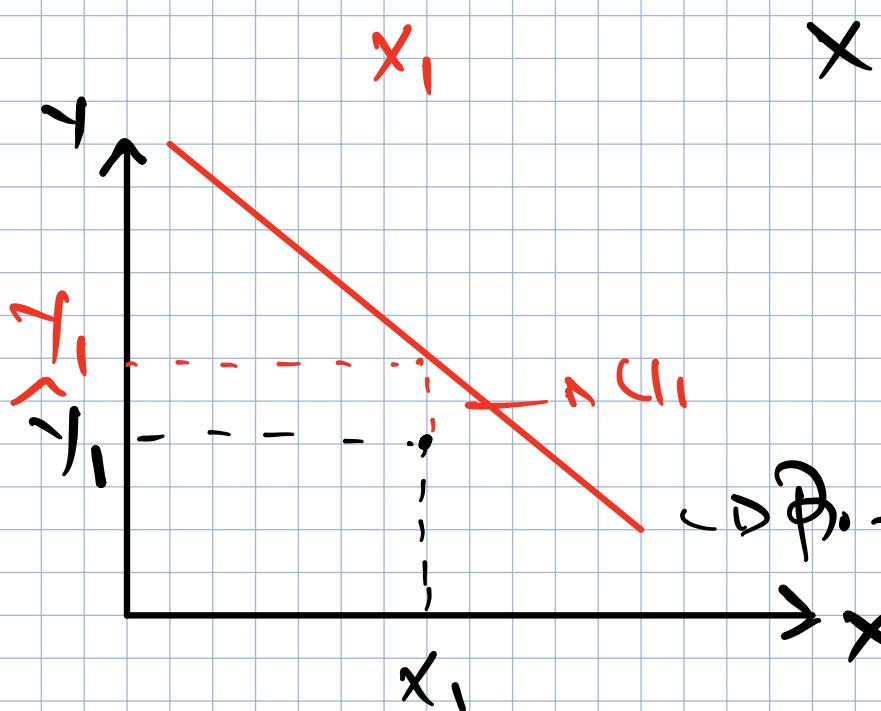
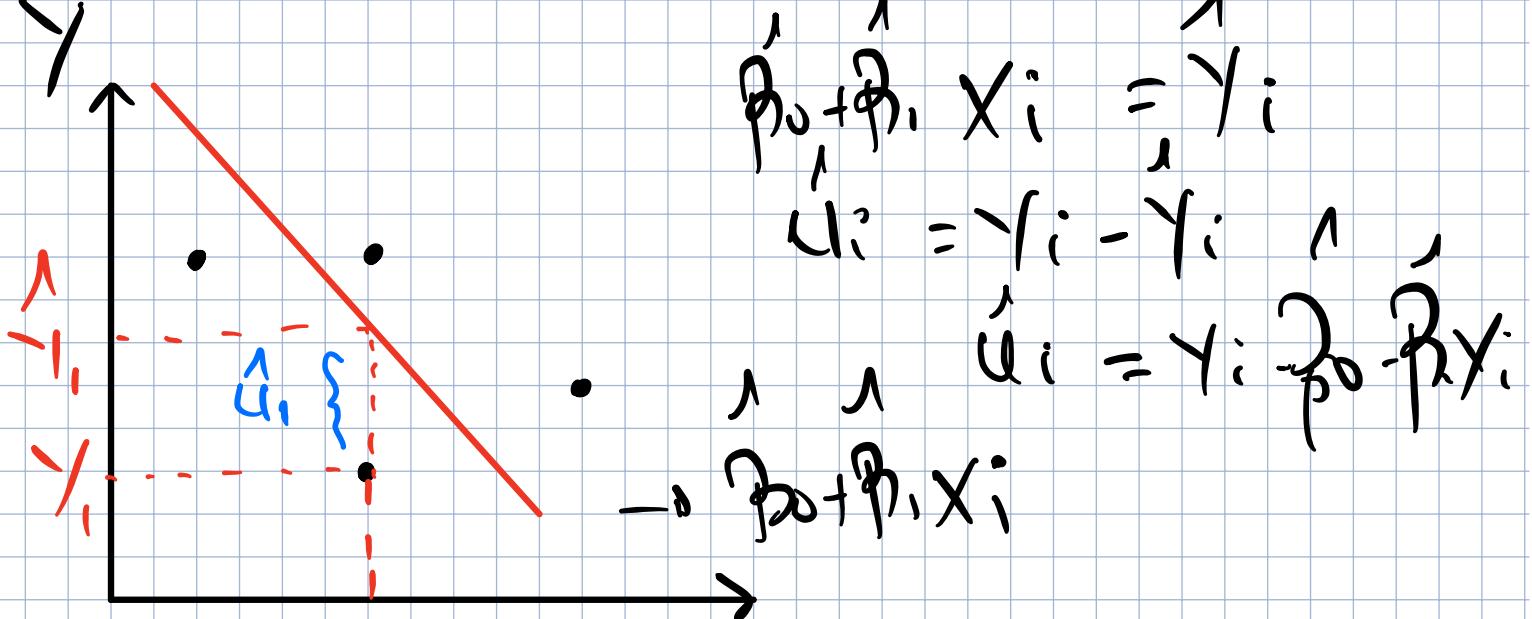
2 STUDENTS  $\rightarrow$  Almost 1SD  $\widehat{\sigma}_1$

$$2 \times 2.83 = 5.66$$

A BETTER WAY

$$\widehat{Y} = 694 \rightarrow S_{\widehat{Y}} = 19.05$$

STAND.  
TEST SCORES



MEASURES OF FIT :  $R^2$  - STANDARD ERROR OF THE REGRESSION

$R^2$  : SHARE OF THE VARIATION IN  $Y$  EXPLAINED BY THE VARIATION IN  $X$ .

$$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{u}_i$$

$$Y_i = \hat{Y}_i + \hat{u}_i$$

## EXPLAINED SUM OF SQUARES (ESS)

→ SUM OF THE SQUARED DEVIATIONS OF THE  $y_i$  FROM THE AVERAGE

$$ESS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \rightarrow \text{TOTAL VARIATION IN } \hat{y}$$

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2 \rightarrow \text{TOTAL VARIATION IN } y$$

$$R^2 = \frac{ESS}{TSS} \rightarrow \text{TOTAL SUM OF SQUARES}$$

SSR → SUM OF SQUARED RESIDUALS

$$SSR = \sum_{i=1}^n u_i^2 \quad R^2 = 1 - \frac{SSR}{TSS}$$

X CAN EXPLAIN  $\Leftrightarrow 0 \leq R^2 \leq 1$   
 NO VARIATION IN Y  $\Leftrightarrow$  X CAN FULLY EXPLAIN THE VARIATION IN Y

$$1. \frac{1}{n} \sum_i c_i u_i = 0$$

$$4. \sum_i u_i^2 \text{ minimum}$$

$$2. \frac{1}{n} \sum_i \hat{y}_i = \bar{y}$$

$$5. TSS = ESS + SSR$$

$$3. \sum_i \hat{u}_i x_i = 0$$

1. SAMPLE AVERAGE OF THE LS RESIDUALS IS ZERO.

$$\hat{u}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

$$\sum \hat{u}_i = (Y_i - \bar{Y}) - \hat{\beta}_1 (X_i - \bar{X}) \quad \text{so } \bar{Y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{X}$$

$$\sum \hat{u}_i = \underbrace{\sum (Y_i - \bar{Y})}_{0} - \hat{\beta}_1 \underbrace{\sum (X_i - \bar{X})}_{0}$$

$$\sum \hat{u}_i = 0$$

$$1/N \sum \hat{u}_i = 0$$

$$2. Y_i = \hat{Y}_i + \hat{u}_i \rightarrow \sum Y_i = \sum \hat{Y}_i + \sum \hat{u}_i$$

$$\sum Y_i = \sum \hat{Y}_i + 0$$

$$\bar{Y} = \sum \hat{Y}_i$$

$$\sum \hat{Y}_i / N = \bar{Y}$$

→ OMITTED

$$3. \sum \hat{u}_i X_i = \sum \hat{u}_i (X_i - \bar{X})$$

$$\sum \hat{u}_i X_i = \sum [(Y_i - \bar{Y}) - \hat{\beta}_1 (X_i - \bar{X})] (X_i - \bar{X})$$

$$\therefore \sum (Y_i - \bar{Y})(X_i - \bar{X}) - \hat{\beta}_1 \sum (X_i - \bar{X})^2 = 0$$

$$\text{REMEMBER } \hat{\beta}_1 = \frac{\sum (Y_i - \bar{Y})(X_i - \bar{X})}{(X_i - \bar{X})^2}$$

$$\begin{aligned}
 \text{LSS} &= \sum_i (Y_i - \bar{Y})^2 = \sum_i (Y_i - \hat{Y}_i + \hat{Y}_i - \bar{Y})^2 \\
 &= \sum_i (Y_i - \hat{Y}_i)^2 + \sum_i (\hat{Y}_i - \bar{Y})^2 + 2 \sum_i (Y_i - \hat{Y}_i) \hat{Y}_i - \bar{Y} \\
 \hookrightarrow \text{SSR} + ESS + 2 \sum_i \hat{u}_i \hat{Y}_i
 \end{aligned}$$

$$\begin{aligned}
 \sum_i \hat{u}_i \hat{Y}_i &= \sum_i \hat{u}_i (\hat{\beta}_0 + \hat{\beta}_1 X_i) \\
 &= \hat{\beta}_0 \sum_i \hat{u}_i + \hat{\beta}_1 \sum_i \hat{u}_i X_i \\
 &\quad \text{---} \quad \text{---} \\
 &= 0
 \end{aligned}$$

## STANDARD ERROR OF THE REGRESSION

A MEASURE OF THE SPREAD OF THE OBSERVATIONS AROUND THE REGRESSION LINE.

$$\text{SER} = S_{\hat{u}} = \sqrt{\frac{\sum \hat{u}_i^2}{n-2}}$$

SER = STANDARD DEVIATION OF THE RESIDUALS  
 Magnitude of a typical regression error expressed in the same unit as the outcome

$S_{\hat{u}}$  =  $\sqrt{\frac{\text{SSR}}{n-2}}$   
 IN  $\hat{\beta}_0, \hat{\beta}_1$ , SHOULD BE ESTIMATED  
 BECAUSE YOU CAN COMPUTE  $S_{\hat{u}}$

# THE LEAST SQUARES ASSUMPTIONS

$\hat{\beta}_0, \hat{\beta}_1$  ARE "GOOD" ESTIMATES FOR  $\beta_0, \beta_1$   
WHEN

THE CONDITIONAL DISTRIBUTION OF  $u_i$  GIVEN

$X_i$  HAS A MEAN OF ZERO  $\rightarrow$  CAUSAL INTERPREATION OF  $X$

$$E[u_i | X_i] = 0 \leftrightarrow \text{Corr}(u_i, X_i) = 0$$

1 STR

2.63 POINTS AUTOMATICALLY SATISFIED WHEN  $X_i$  IS

new assumption:  $STR = 16$

DIFFERENT

RANDM

$STR = 17$

DISTRICTS

BUT  $\beta_0, \beta_1$  ARE BIASED

WEAKER

$$Y_i = \beta_0 + \beta_1 X_i + u_i + \text{CORR}(u_i, X_i) \neq 0$$

INCLUDES ALL

THE OTHER

STUFF THAT

DETERMINES THE OUTCOME

$\text{Cor}(poverty, str) > 0$   
MORE POVERTY  $\rightarrow$  LESS

SMALL LARGE

• - •

$\rightarrow$  -1 POINT (IF THEY NEVER EQUAL)

SMALL - LARGE

$\rightarrow$  -2.63 POINTS (ALSO  $\propto$ )  
BECAUSE (DYNAMIC)

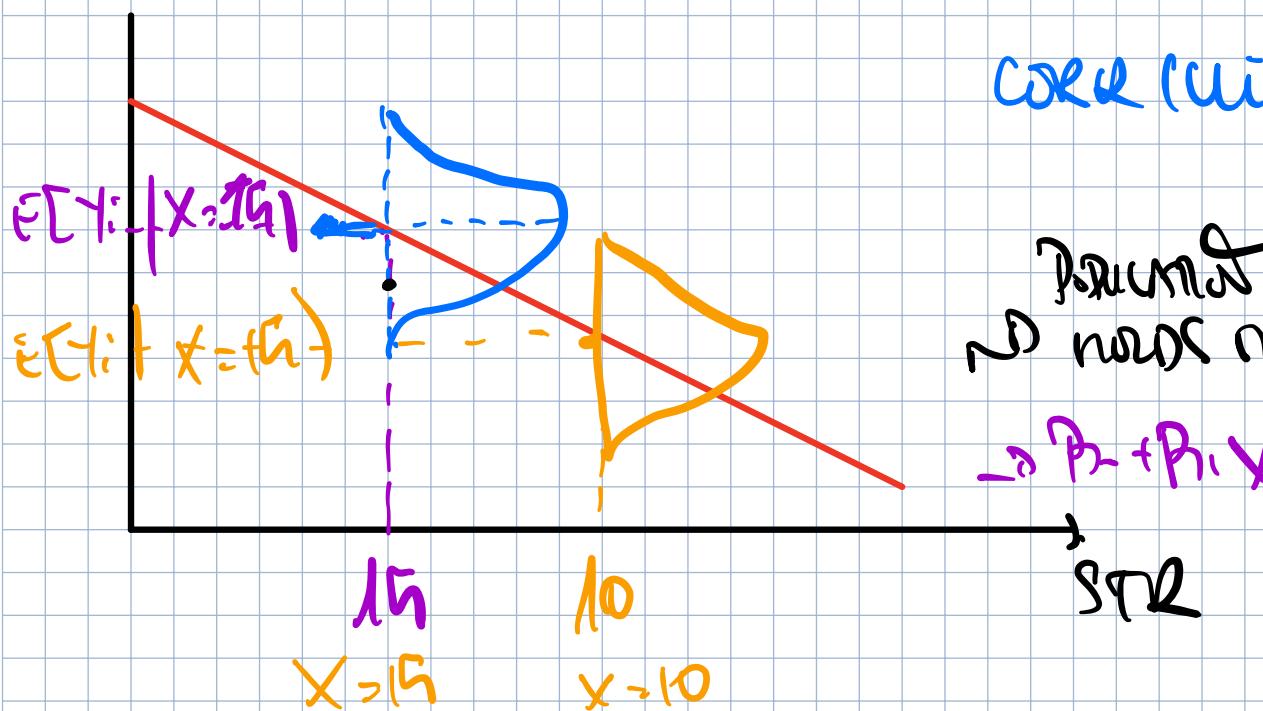
↳ CLASSROOM  $\Leftarrow$

$$Y_i = \beta_0 + \beta_1 X_i + U_i$$

$$E[Y_i] = \beta_0 + E[X_i] + E[U_i]$$

$$E[Y_i | X_i] = \beta_0 + \beta_1 X_i + E[U_i | X_i]$$

TERMS SCORED  
NO



CORR ( $U_i, X_i$ ) = 0  
No linear relationship  
No words in sentence

$\rightarrow \beta_0 + \beta_1 X$

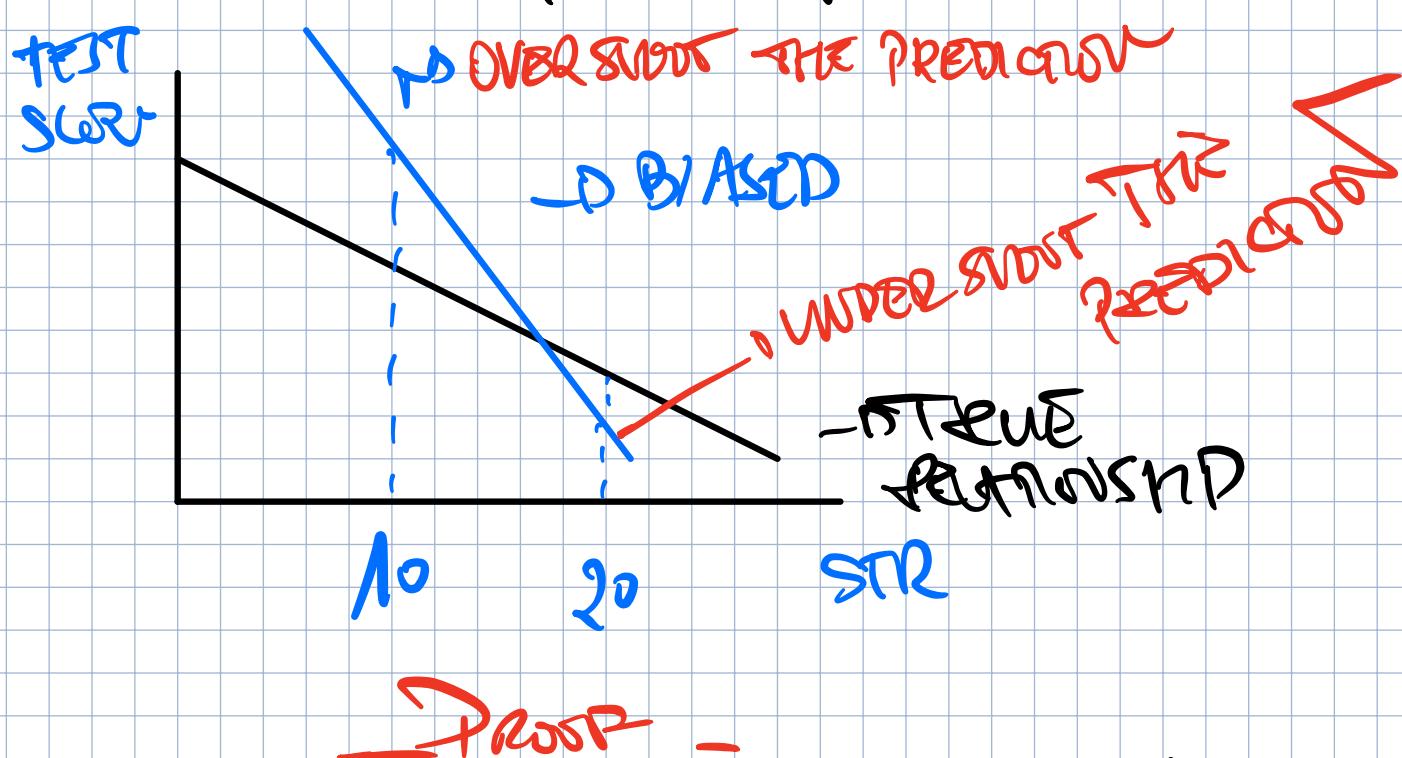
↳ AT A GIVEN VALUE OF  $X$ ,  $Y$  IS DISTRIBUTED ALONG THE PREDICTION LINE, AND THE ERROR  $U = Y - (\beta_0 + \beta_1 X)$

↳ HAS A CONDITIONAL MEAN OF ZERO

PROOF : THE DIFFERENCE B/W

$$U_i | X_i \quad Y_i | X_i \text{ & } E[Y_i | X_i]$$

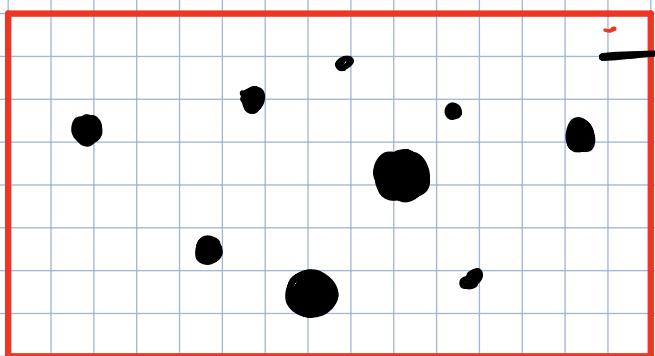
ASSUMPTION IS IF  $\text{COR}(u_i, x_i) = 0$   
 THEN YOU SOMETIMES UNDERPREDICT  
 & SOMETIMES OVER PREDICT BUT ON AVERAGE  
 $E[u_i | x_i]$  MEANS OF THE ERROR IS ZERO.



$$\begin{aligned}
 \text{Cov}(x, u) &= E[(x - E[x])(u - E[u])] \\
 &= E[xu] - E[x]E[u] - E[x]E[u] + E[x]E[u] \\
 &= E[xu] - E[x]E[u] \\
 &= E[xu] - E[x]E[\mathbb{E}(u|x)] \\
 &\stackrel{\text{def}}{=} E[\mathbb{E}(xu|x)] - E[x]E[\mathbb{E}(u|x)] \\
 &= E[\mathbb{E}(xu|x)] - E[x]E[0] \\
 &= E[\mathbb{E}(xu|x)] = E[x\mathbb{E}(u|x)] = 0
 \end{aligned}$$

$$\frac{\text{cov}(x_i u)}{\text{var}(y)} = \text{cov}(x_i u) = 0$$

**ASSUMPTION 2:**  $(x_i, y_i) \quad i=1, \dots, n$  ARE INDEPENDENTLY & IDENTICALLY DISTRIBUTED.  $\rightarrow$  ASSUMPTION ABOUT HOW SAMPLE IS DRAWN  
RANDOM SAMPLING INSURES THIS.



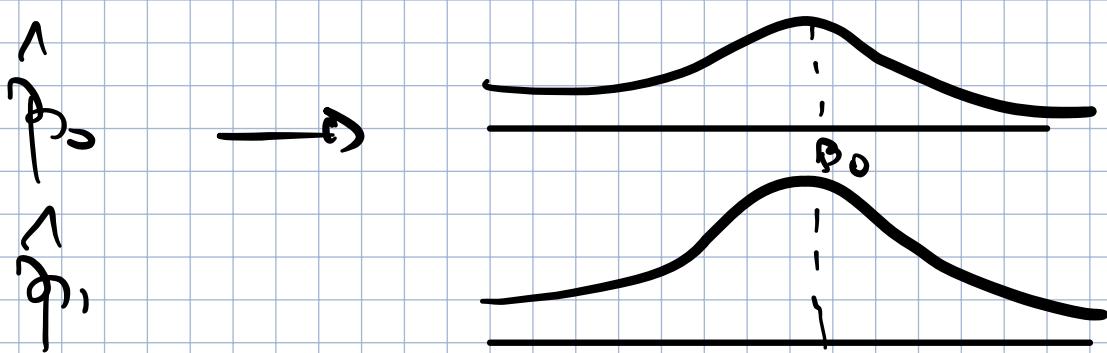
$\rightarrow$  IDENTICALLY DISTRIBUTED  
(COME FROM THE SAME)  
PROBABILITY DISTRIBUTION  
 $\rightarrow$  INDEPENDENT: DRAWS  
HAVE NO MEMORY.

**ASSUMPTION 3:** LARGE OUTLIERS ARE UNLIKELY  
OLS MIGHT BE MISLEADING IF THERE ARE LARGE OUTLIERS  
(NEED FOR CONSISTENCY  $\sigma^2$ )

$$S_2 Y \underset{\text{P}}{\rightarrow} Q_Y^2 \rightarrow \left[ \text{REGULARITY ASSUMPTION} \right]$$

SAMPLING DISTRIBUTION OF THE OLS ESTIMATOR





WHEN  $\gamma$  IS LARGE  
ALL OLS ASSUMPTIONS HOLD

$$E(\hat{\beta}_0) = \beta_0$$

$$E(\hat{\beta}_1) = \beta_1$$

BY  
CENTRAL LIMIT  
THEOREM

$\rightarrow$  ASSUMPTIONS HOLD

1.  $E(u_i | x_i) = 0$
2.  $\text{cov}(u_i, x_i) = 0$
3.  $(x_i, y_i)$  iid.

### Large-Sample Distributions of $\hat{\beta}_0$ and $\hat{\beta}_1$

#### KEY CONCEPT

#### 4.4

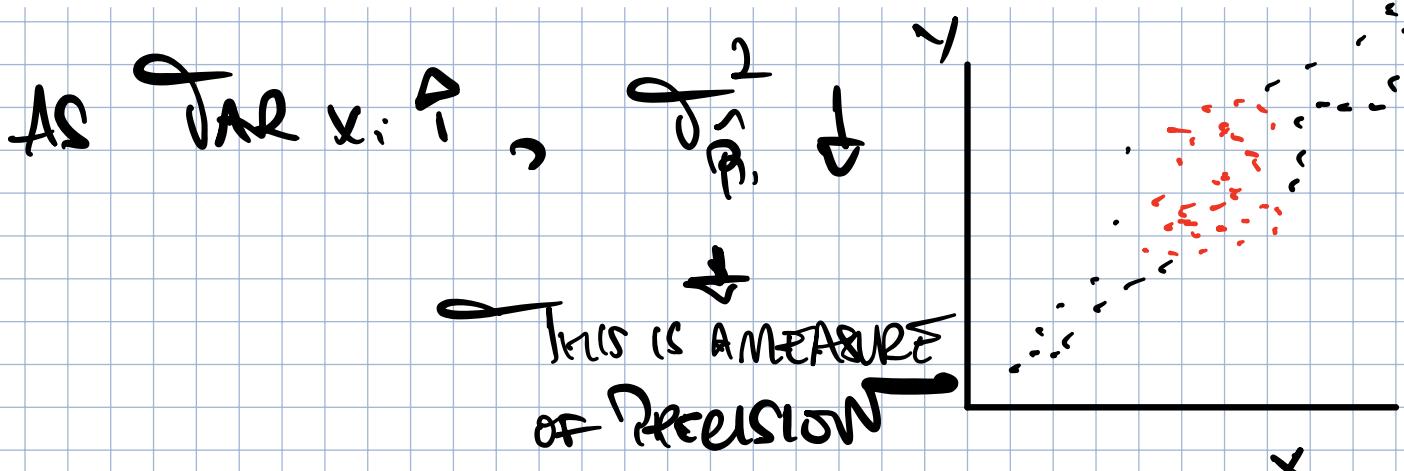
If the least squares assumptions in Key Concept 4.3 hold, then in large samples  $\hat{\beta}_0$  and  $\hat{\beta}_1$  have a jointly normal sampling distribution. The large-sample normal distribution of  $\hat{\beta}_1$  is  $N(\beta_1, \sigma_{\hat{\beta}_1}^2)$ , where the variance of this distribution,  $\sigma_{\hat{\beta}_1}^2$ , is

$$\sigma_{\hat{\beta}_1}^2 = \frac{1}{n} \frac{\text{var}[(X_i - \mu_X)u_i]}{\text{var}(X_i)^2}. \quad (4.21)$$

IMPORTANT

The large-sample normal distribution of  $\hat{\beta}_0$  is  $N(\beta_0, \sigma_{\hat{\beta}_0}^2)$ , where

$$\sigma_{\hat{\beta}_0}^2 = \frac{1}{n} \frac{\text{var}(H_i u_i)}{E(H_i^2)^2}, \text{ where } H_i = 1 - \left[ \frac{\mu_X}{E(X_i^2)} \right] X_i. \quad (4.22)$$



- ANSWER TO QUESTIONS -

a.  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

$$Y_i = \hat{Y}_i + \hat{U}_i$$

$$\hat{\beta}_0 = 509.384 \quad X_i = 22$$

$$\hat{\beta}_1 = -5.614$$

$$\hat{Y}_i = 509.384 - 5.614 \times 22$$

$$\hat{Y}_i \approx 385.9$$

$$\rightarrow \hat{\beta}_1 = \frac{\Delta Y_i}{\Delta X_i}$$

b.

$$\Delta X_i = 23.19$$

$$= 4$$

$$\Delta Y_i = 4 \times -5.614$$

$$= -22.5$$

c.  $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$

$$\bar{Y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{X}$$

$$= 509.384 - 5.614 \times 21 = 389.2$$

$$\bar{X} = ?$$

**Q.** If we use CS to predict test scores, then we will be off by 105.8 points

$$Q. \text{ S.E.R} = \sqrt{\frac{SSR}{n-2}}$$

$$105.8 = \sqrt{\frac{SSR}{98}}$$

$$11193.64 = SSR / 98$$

$$SSR = 1096977$$

$$R^2 = 1 - \frac{SSR}{TSS} \rightarrow 0.0849 = 1 - \frac{1096977}{TSS}$$

$$S^2_y = \frac{TSS}{(n-1)} \rightarrow 0.9151 = \frac{1096977}{TSS}$$

$$TSS = 1198791$$

$$S^2_x = 12109 \quad S_y = 110$$

# CHAPTER 5. HYPOTHESIS TESTS & CONFIDENCE INTERVALS

$$\left\{ \begin{array}{l} Y_i = \beta_0 + \beta_1 X_i + u_i \\ Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{u}_i \\ Y_i = \hat{Y}_i + \hat{u}_i \end{array} \right.$$

USING A SINGLE REGRESSION - SAMPLING DISTRIBUTION

ESTIMATOR - HYPOTHESES TEST VALUE

$$t = \frac{\text{ESTIMATOR} - \text{HYPOTHESES TEST VALUE}}{\text{STANDARD ERROR OF THE ESTIMATOR}}$$

IN MOST COMMON CASE:  $H_0: \beta_1 = 0$

GENERAL CASE:

$$H_0: \beta_1 = \beta_{1,0}$$

$$\beta_1 \neq \beta_{1,0}$$

$$H_A: \beta_1 \neq 0$$

THREE STEPS

1.  $\text{SET } \hat{\beta}_1 \rightarrow \text{ESTIMATOR FOR } \sigma_{\hat{\beta}_1}$

$$\text{SE}[\hat{\beta}_1] = \sqrt{\sigma_{\hat{\beta}_1}^2}$$

LINE

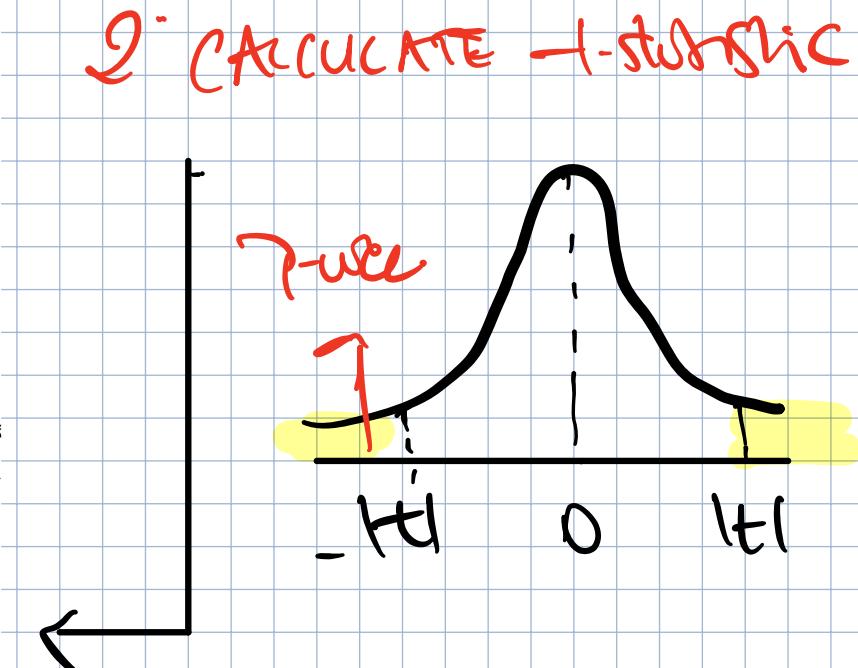
(STANDARD DEVIATION OF THE SAMPLING DISTRIBUTION)

$$SE(\hat{\beta}_1) = \sqrt{\hat{\sigma}_{\hat{\beta}_1}^2}$$

$$\hat{\sigma}_{\hat{\beta}_1}^2 = \frac{1}{n} \times \frac{\frac{1}{n-2} \sum_{i=1}^n (X_i - \bar{X})^2 \hat{u}_i^2}{\left[ \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right]^2}$$

the variance in Equation (5.4) is discussed  
formula for  $\hat{\sigma}_{\hat{\beta}_1}^2$  is complicated, in applications the  
ession software so that it is easy to use in pra  
ep is to compute the *t*-statistic,

$$t = \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)}$$



**3. COMPUTE P-value**: THE PROBABILITY OF  
OBSERVING A VALUE OF  $\hat{\beta}_1$  UNDER THE  
ASSUMPTION THAT  $\beta_1 = \beta_{1,0}$  *RESERVED*

ANY BETA  $\Rightarrow$  DEPENDS  
↑ ACTUAL ON THIS VALUE  
*t*-STAT

$$\text{P-value} = \Pr_{H_0} \left[ \left| \hat{\beta}_1 - \beta_{1,0} \right| > \left| \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)} \right| \right]$$

$$\hookrightarrow \Pr_{H_0} \left[ \left| \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)} \right| > \left| \frac{\hat{\beta}_1 - \beta_{1,0}}{SE(\hat{\beta}_1)} \right| \right]$$

$$\hookrightarrow \Pr_{H_0} (|t| > |t|)$$

Absolute =  $\frac{\text{abs value}}{2}$   $\Rightarrow$  2/C  
(TEST)

$\Pr(Z \geq 1.77)$   
 If we sample  
 $\Rightarrow 2\Phi(-1.77)$

Reject the Null if

$t\text{-stat} < \text{pre-specified threshold}$

$\rightarrow 1\%, 5\%, 10\%$

Center  
+ - values

$$\begin{array}{c} \downarrow \\ 2.58 \end{array} \quad \begin{array}{c} \downarrow \\ -1.96 \end{array}$$

$\rightarrow 1.645$

~~Answer to question~~

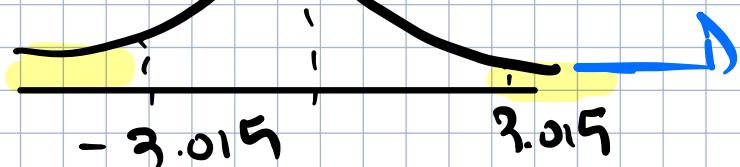
Q.  $H_0: \beta_1 = 0$

$H_A: \beta_1 \neq 0$

1.  $\text{SE}[\hat{\beta}_1] = 1.862$

$-5.614 - 0$   
 2.  $t\text{-stat} = -3.015$

3.  $P\text{-value} = 2\Phi(-3.015) \Rightarrow 1.86$



27

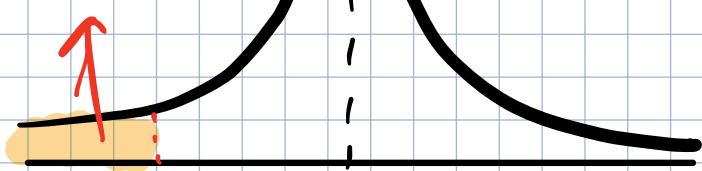
$P\text{-value} < 0.01 \rightarrow$  REJECT THE NULL HYPOTHESIS.

$$H_0: \bar{p}_1 = 0.162 \quad t\text{-stat} = \frac{-5.614}{1.862}$$

$$P\text{-value} = \Phi(-3.015) = -3.015$$

$$P\text{-value} = 0.001635$$

P-value.



-3.015

$$H_0: \bar{p}_1 = 0$$

$$H_A: \bar{p}_1 < 0$$

$t\text{-stat} = 9.5 \rightarrow ?$  ~~Fail to reject the null.~~  
~~DO NOT REJECT THE null.~~

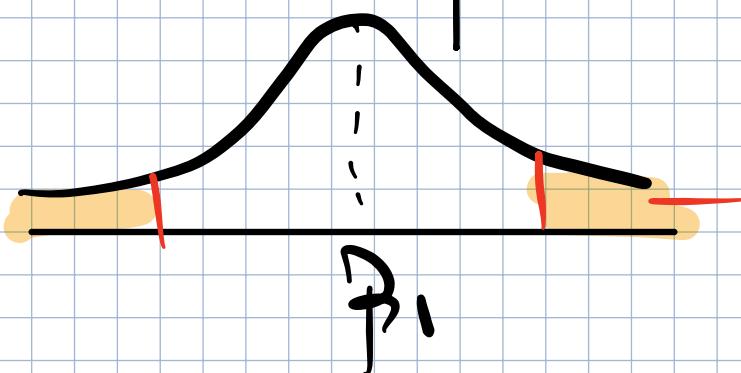
TESTING HYPOTHESES ABOUT THE POPULATION  
 $P_0 \rightarrow$  KUMAR'S USE AN INCORRECT  
EVEN IF YOU FAIL TO REJECT THE NULL

# CONFIDENCE INTERVAL FOR $\hat{P}_1$ (95%)

2 DEFINITIONS:

1. SET OF VALUES THAT CAN'T BE  
PREDICTED USING A TWO-SIDED CI.

2° 95% CONTAINING THE TRUE  
VALUE OF  $\hat{P}_1$



→ 95% OF ALL  
POSSIBLE  
SAMPLES OF  
 $\hat{P}$  WILL BE  
PREDICTED

$$\hat{P}_1 - 1.96 \text{SE}[\hat{P}_1]$$

$$\hat{P}_1 + 1.96 \text{SE}[\hat{P}_1]$$

$$99.7\% \text{ CI} \rightarrow \hat{P}_1 \pm 2.576 \times \text{SE}[\hat{P}_1]$$

$$-5.614 \pm 2.576 \times 1.162$$

$$[-10.50 \quad -0.7231]$$

~~CONFIDENCE IN PREDICTION FOR CHANGING X~~

$$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + u_i$$

we consider  $\Delta x \rightarrow \text{CHANGE IN } X$   
THE PREDICTED CHANGE IN Y ( $\Delta Y$ )

(1) THE REGRESSION SLOPE (slope)

$$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + u_i$$

$$\Delta Y_i = \hat{\beta}_1 \Delta x_i$$

(2) 99% CONFIDENCE INTERVAL  
FOR  $\Delta x$ :  $\hat{\beta}_1$

$$[\hat{\beta}_1 \Delta x \pm 1.96 \text{ SE}[\hat{\beta}_1] \times \Delta x]$$

$$\Delta x - [\hat{\beta}_1 \mp 1.96 \text{ SE}[\hat{\beta}_1]]$$

$$\Delta x [+10.50, +0.7231]$$

$$[+10.50, +0.7231]$$

(3)  $[21.01, 1.45]$  ACROSS IN  
TEST SCORE

# REGRESSION WHEN $X$ IS A BINARY VARIABLE

ONE OF THE MOST COMMON CASES:

- > GENDER (MALE FEMALE)  
(1 MALE vs NO-MALE)

$$\left\{ \begin{array}{l} X \\ 0 \\ 1 \end{array} \right\}$$

INDICATOR, DUMMY, BINARY, INDEX VARIABLE

$$D_i = \{ 0, 1 \}$$

$$X_i \rightarrow D_i = \begin{cases} \text{SMALL STR} < 20 \\ \text{MORE STR} \geq 20 \end{cases}$$

$$Y_i = \beta_0 + \beta_1 D_i + u_i \quad i=1, \dots, n$$

$\hat{\beta}_1$  -> SLOPE -> HOW TO INTERPRET THIS

DEPERCENT?

DO NOT ASSUME ANYMORE

$$Y_i = \beta_0 + u_i \quad (\bar{D}_i = 0)$$

$$Y_i = \beta_0 + \beta_1 + u_i \quad (\bar{D}_i = 1)$$

EASIEST WAY

$$\mathbb{E}[Y_i | D_i=0] = \mathbb{E}[\beta_0] + \sum u_i | D_i=0$$

$$\mathbb{E}[Y_i | D_i=1] = \mathbb{E}[\beta_0] + \mathbb{E}[\beta_1] + \sum u_i | D_i=1$$

DLS ASSUMPTION 1.

$$\mathbb{E}[u_i | X_i] = 0$$

POPULATION MEAN  
OF  $Y_i$  WHEN

$$\mathbb{E}[Y_i | D_i=0] = \beta_0$$

$$\uparrow \quad D_i=0$$

$$\mathbb{E}[Y_i | D_i=1] = \beta_0 + \beta_1$$

$$\beta_1 \rightarrow \mathbb{E}[Y_i | D_i=1] - \mathbb{E}[Y_i | D_i=0]$$

→ POPULATION MEAN  
OF  $Y_i$  WHEN  $D_i=1$

$\beta_1 \rightarrow$  DIFFERENCE IN MEAN TEST SCORES  
IN STUDENT-TO-TEACHER RATIO AND

HIGH STUDENT-TO-TEACHER RATIO.

$\hat{\beta}_1 \rightarrow$  DIFFERENCE IN SAMPLE  
AVERAGE OF  $\bar{Y}$  WHEN  $D_i=1$

VS  $D_i=0$

$\bar{Y} | D_i=0 \rightarrow 369.92$  → AVERAGE TEST SCORE FOR STUDENTS IN LARGE CLASSROOM.

$\bar{Y} | D_i=1 \rightarrow 369.42 + 44.45$   
 $\rightarrow 414.37$  → AVERAGE TEST SCORE FOR STUDENT IN SMALL CLASSROOM.

$\bar{Y} | D_i=1 - \bar{Y} | D_i=0 \rightarrow 44.45$  POINTS  
 OR DIFFERENCE BETWEEN SMALL & LARGE CLASSROOMS:

$$44.45 \pm 1.96 \times 22.19$$

$$[0.96 \quad 87.94] \rightarrow 95\% \text{ CI}$$

(D DOES NOT INCLUDE ZERO.)

$$1.5\text{SE} = \frac{\hat{\beta}_1 - \hat{\beta}_{1,0}}{\text{SE}(\hat{\beta}_1)} = \frac{44.45}{22.19} = 2.003$$

# CH 6. LINEAR REGRESSION WITH MULTIPLE REGRESSORS

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + u_i$$

$\checkmark$  - independent variables

IF WE ARE INTERESTED IN THE RELATIONSHIP BUT ONE X & Y; WHY DO WE NEED MORE THAN ONE INDEPENDENT VARIABLE?

THE IMPACT OF CLASSROOM SIZE  $\rightarrow$  TEST SCORES (CENSUS) IN CA

OMITTED VARIABLE BIAS:

CA HAS A LARGE IMMIGRANT POPULATION

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

TEST RATES

TEST SCORE

$\rightarrow \beta_1$  MIGHT NOT REFLECT THE TRUE PROBABILITY WELL  
NOT PERFECTLY  
TRUE IMPACT OF  
THE STR ON TEST SCORE

NEW FACTOR: % OF ENGLISH LEARNERS  
IN A DISTRICT.

↳ DISTRICTS WITH A HIGHER SHARE OF  
ENGLISH LEARNERS PERFORM WORSE IN  
ENGLISH TEST

↳ DISTRICTS WITH LARGE CLASSES MIGHT  
ALSO HAVE A HIGHER SHARE OF ENGLISH  
LEARNERS

SO THE RELATIONSHIP THAT IS DESCRIBED  
BY THE REGRESSION RESULTS MIGHT BE DRIVEN  
BY THE % OF ENGLISH LEARNERS (MIGRANTS)  
WHILE THE TRUE RELATIONSHIP IS MUCH  
WEAKER OR EVEN NON-EXISTENT.

**OMITTED VARIABLE BIAS:** AN IMPORTANT PROPERTY  
THAT i. INFLUENCES THE OUTCOME  
ii. CORRELATED WITH ONE OR MORE INDEPENDENT  
VARIABLES

OMITTED FROM THE REGRESSION.

TWO CONDITIONS: ↳ OMITTED VARIABLES SHOULD  
HAVE AN IMPACT ON THE OUTCOME  
↳ CORRELATED WITH A REGRESSOR.

$E[\text{Wilx}_i]$  to  $\text{Corr}(u_i, x_i)$  to

$\text{Corr}(x_i, u_i) = \rho_{xu}$  1st BIAS & ERROR TERM

$E[\hat{\rho}_i] = \rho_i + \rho_{xu} \frac{\sigma_u}{\sigma_x} \rightarrow$  BIASED DUE  
TO OMITTED  
VARIABLE

$E[\hat{\rho}_i] = \rho_i \rightarrow$  UNBIASED

1. ↑ SAMPLE SIZE WILL NOT HELP.

2. WHETHER THE OMISS IS SMALL OR LARGE

DEPENDS ON  $\text{Corr}(u_i, x_i)$  &  $\text{Corr}(x_i, y_i)$

CORR THE COEFF, INVERSE THE OJB.

ALSO THE DIRECTION DEPENDS ON  
THE THREE WAY RELATIONSHIP B/W

$x_{1i}, x_{2i}, y_i$  → OUTCOME

↓  
OMITTED  
VAR

COMPLEXITY OF THIS  
CASE

# ADDRESSING THE OMITTED VARIABLE BIAS:

MEASURE THE IMPACT OF SIR ON PERTAINING % OF ENGLISH LEARNERS CONSTANT.  
IDEA: COMPARE SMALL VS LARGE CLASS SIZES AMONG DISTRICTS WITH SIMILAR % OF ENGLISH LEARNERS

## THE MULTIPLE REGRESSION MODEL

31/11/2020

TO SIR;

IDEA: ESTIMATE THE IMPACT OF  $X_{1i}$  ON TEST SCORE  $Y_i$  WHILE HOLDING  $X_{2i}$

PREDICTION REGRESSION FUNCTION LINE ( $\hat{Y}_i$ )  
% OF ENGLISH LEARNERS

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$$

$$\bar{E}[Y_i | X_{1i} = x_1, X_{2i} = x_2] = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

PREDICTION → PLACE OF TEST SCORES FOR DISTRICTS WITH  $X_1$ , SIR &  $X_2$  % OF ENG. LEARNERS

$\beta_0 \rightarrow$  INTERCEPT

$\beta_1 \rightarrow$  SLOPE COEFFICIENT FOR  $X_1$ :

$\beta_2 \rightarrow$  SLOPE COEFFICIENT FOR  $X_2$ :

$X_1, X_2 \rightarrow$  CONTROL VARIABLES

$\beta_1:$  CHANGE IN  $Y$  INDUCED BY A UNIT  
CHANGE IN  $X_1$ , HOLDING  $X_2$  CONSTANT.

AFTER ACCOUNTING FOR  $X_2$ :  
AFTER CONTROLLING FOR  $X_2$ :

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

$$\underline{Y + \Delta Y = \beta_0 + \beta_1 (X_1 + \Delta X_1) + \beta_2 X_2}$$

$$\Delta Y = \beta_1 \Delta X_1 \quad \beta_1 = \frac{\Delta Y}{\Delta X_1}$$

( $\downarrow$ ) FINAL EFFECT  
OF  $X_1$  ON  $Y$

$$E[Y | X_1=0, X_2=0] = \beta_0$$

POPULATION AVERAGE  
WHEN  $X_1=0, X_2=0$

# BRUCAN ON MULTIPLE REGRESSION MODEL:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$$

↗ different  
 ↗ for each person

SPECIFIC PARTS ↓  
 OF  $Y_i$  THAT CAN  
 BE EXPLAINED BY  
 $X_{1i} + X_{2i}$

$$i = 1, 2, \dots, n$$

$$Y_i = \beta_0 X_{0i} + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + u_i$$

Constant  
Depressur  
(CONSTANT TERM)

$\underbrace{\beta_0 + \beta_1 + \dots + \beta_k}_{n \times 1}$

MORE GENERAL

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + u_i$$

homoskedasticity Assumption:

$$\text{Var}(u_i | X_{1i}, X_{2i}, \dots, X_{ki}) = \sigma^2$$

↗ constant

$$E(Y_i | X_{1i} = x_1, X_{2i} = x_2, \dots, X_{ki} = x_k)$$

$$\approx \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

$\beta_1 \rightarrow \frac{\Delta Y}{\Delta X_1}$  THE CHANGE IN  $Y$  INDUCED BY  
A UNIT CHANGE IN  $X$  HOLDING  
ALL ELSE CONSTANT.

### THE LS ESTIMATOR

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + u_i$$

$$\text{FIND } \hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k \rightarrow \min \sum_{i=1}^n u_i^2$$

$\underbrace{\qquad}_{K+1 \text{ PARAMETERS}}$   
 $\rightarrow \text{ESTIMATE}$

$$\text{WHERE } Y_i = \hat{Y}_i + \hat{u}_i$$

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_k X_{ki}$$

$$\hat{u}_i = Y_i - \hat{Y}_i \quad \min \sum_{i=1}^n \hat{u}_i^2 \quad i(Y_i - \hat{Y}_i)^2$$

DISTRICTS WITH A HIGH % OF ENGLISH LEARNERS  
TEND TO HAVE NOT ONLY LOW TEST SCORES BUT ALSO  
A HIGH STUDENT/TEACHER RATIO WHEN MEASURED  
ENGLISH LEARNERS (ARMED ESTIMATED COEFFICIENT)

TRUE MODEL:

MODEL 1:  $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + u_i$

↳ WHAT YOU ARE

MODEL 2:

INTERESTED IN,  $\wedge \wedge$

$$\hookrightarrow \beta_0, \beta_1, \beta_2$$

ESTIMATED MODEL:

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1i} + u_i$$

$$\sim \sim \\ \hat{\beta}_0, \hat{\beta}_1$$

$x_2$  IS OMITTED  $\rightarrow$  NOT AVAILABLE  
IN THE DATA

$$\hat{\beta}_1 = -2.6210$$

AUXILIARY MODEL:

$$\hat{\beta}_1 = -1.28970$$

$$x_{2i} = \alpha_0 + \alpha_1 x_{1i} + \epsilon_i$$

$$\hat{\alpha}_2 = -0.73403$$

$\hat{\epsilon}_{2i}$   
 $\hat{\epsilon}_{2i}$

$$\hat{\alpha}_1 = 1.8137$$

$$\hat{\beta}_1 = \hat{\beta}_1 + \hat{\beta}_2 \times \hat{\alpha}_1$$

$$-2.6210 = -1.28970 - 0.73403 \times 1.8137$$

$$= -2.6210$$

$$\text{PIAS} = \hat{\beta}_2 \times \hat{\alpha}_1 + 0$$

$\underbrace{\phantom{000}}$

-1.93

↑

# OLS ESTIMATOR

OBJECTIVE: Estimate

$$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_K X_{Ki} + U_i$$

FIND  $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_K$  THAT MINIMISE

SIMULTANEOUSLY

( $K+1$ ) PARAMETERS

$K$ -INDEPENDENT VARIABLES

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_K X_{Ki}$$

$$\hat{U}_i = Y_i - \hat{Y}_i$$

$$\text{MIN } \sum \hat{U}_i^2 \leftarrow \text{MIN } \sum_i (Y_i - \hat{Y}_i)^2$$

SSR or  
LSS

ADJUSTED  
 $\uparrow R^2$

MEASURES OF FIT: SEK + RMSE + Adj. R<sup>2</sup>

STANDARD  
ERROR OF  
THE REGRESSION

$\downarrow$   
Root Mean  
Squared Error

SER & RMSE  $\rightarrow$  VERY SIMILAR

BOTH MEASURE THE SPREAD OF  $y_i$  AROUND  $\hat{y}_i$

$$SER = \sqrt{\frac{1}{(n-k-1)} \sum_i \epsilon_i^2}$$

↑  
PREDICTED  
OUTCOME  
VS  
OBSERVED  
OUTCOME

$$RMSE = \sqrt{\frac{1}{n} \sum_i \epsilon_i^2}$$

}   
 n - DEGREES OF  
 FREEDOM  
 }   
 n - NO. OF  
 INDEPENDENT  
 VARIABLES

BOTH ARE MEASURED IN THE UNITS OF OUTPUT.

$$\hat{y}_i = \bar{y}_i + \epsilon_i$$

RESIDUAL

} AVERAGE DISTANCE  
B/W OBSERVED  
VALUES JS

REGRESSION LINE

AVERAGE PREDICTION ERROR

WHEN  $n$  IS LARGE RMSE & SER ARE CLOSE TO EACH OTHER.

ADJUSTED R-SQUARED:  $R^2$  ADJUSTED FOR THE  
NB OF PREDICTORS IN THE EQUATION

$$R^2 = \frac{ESS}{TSS} \quad ESS = \sum_i (\hat{Y}_i - \bar{Y})^2$$

$$TSS = \sum_i (Y_i - \bar{Y})^2$$

$$SSR = \sum_i \hat{U}_i^2$$

You ADD A NEW

VARIABLE  $R^2$  WILL  
GROW INDEPENDENT

$$TSS = ESS + SSR$$

OF THAT VARIABLE IS USEFUL OR NOT.

SOLUTION: ADJUSTED R-SQUARED

$$\bar{R}^2 = 1 - \left[ \frac{(n-1)}{(n-k-1)} \right] \frac{SSR}{TSS}$$

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

↑ PENALIZE FOR  
EACH PREDICTOR

$$\frac{(n-1)}{(n-k-1)} > 1 \rightarrow \bar{R}^2 \leq R^2$$

IF  $n \rightarrow \infty$   $\widehat{R}^2 \rightarrow R^2$

$\widehat{R}^2 \leq 0 \rightarrow$  THEORETICALLY

LEAST SQUARES ASSUMPTIONS FOR MULTIVARIATE

REGRESSION:

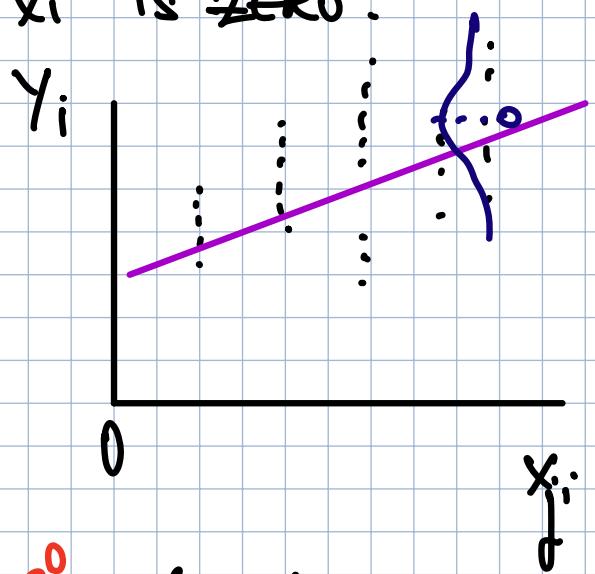
$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + u_i$$

$$i = 1, \dots, n$$

1<sup>o</sup>  $E[u_i | X_{1i} = x_1, X_{2i} = x_2, \dots, X_{ki} = x_k] = 0$

↳ THE CONDITIONAL DISTRIBUTION OF  $u_i$  GIVEN

$X_i$  IS ZERO.



} NO OMITTED  
VARIABLE      TWO CONDITIONS  
BAS              ↳ FOR THE  
OMV

1<sup>o</sup>  $\text{COR}(u_i, X_j) \neq 0$

2<sup>o</sup>  $u_i \rightarrow y_i$

? RANDOM  
SHOCKING  
ENSURES

3<sup>o</sup>  $X_{1i}, X_{2i}, \dots, X_{ki}, u_i \text{ i.i.d.}$

IDENTICALLY & INDEPENDENTLY DISTRIBUTED INSURES

↳ AUTOMATICALLY SATISFIED WITH A RANDOM THIS

3<sup>o</sup> LARGE OUTLIERS ARE UNLIKELY.

↳ CHECK MIN, MAX, MEAN, MEDIAN

↳ No perfect multicollinearity

↗ GAP IS LARGER WHEN THERE ARE OUTLIERS

↳ ONE REGRESSOR IS AN EXACT LINEAR FUNCTION OF ANOTHER REGRESSOR

$$\text{Ex: } Y_i = \beta_0 + \beta_1 X_{1i} + u_i$$

↑  
WAGE

EDUCATION IN YEARS

$$Y_i = \beta_0 + \beta_1 Y_{1i} + \beta_2 X_{2i} + u_i$$

$$Y_{2i} = 2 \times Y_{1i}$$

↖ PERFECT LINEAR FUNCTION

```
> mo5 <- lm(data=CPS1988, formula=wage~education)
> summary(mo5)

call:
lm(formula = wage ~ education, data = CPS1988)

Residuals:
    Min      1Q  Median      3Q     Max 
-786.0   -264.9   -54.8    174.5 18035.1 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -12.8281    11.8970  -1.078   0.281    
education    47.1810     0.8888   53.085  <2e-16 ***  
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 432.4 on 28153 degrees of freedom
Multiple R-squared:  0.09099, Adjusted R-squared:  0.09096 
F-statistic: 2818 on 1 and 28153 DF, p-value: < 2.2e-16

>
> CPS1988$educ2 <- CPS1988$education^2
> mo6 <- lm(data=CPS1988, formula=wage~education + educ2)
> summary(mo6)

call:
lm(formula = wage ~ education + educ2, data = CPS1988)

Residuals:
    Min      1Q  Median      3Q     Max 
-786.0   -264.9   -54.8    174.5 18035.1 

Coefficients: (1 not defined because of singularities)
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -12.8281    11.8970  -1.078   0.281    
education    47.1810     0.8888   53.085  <2e-16 ***  
educ2          NA        NA       NA      NA      
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 432.4 on 28153 degrees of freedom
Multiple R-squared:  0.09099, Adjusted R-squared:  0.09096 
F-statistic: 2818 on 1 and 28153 DF, p-value: < 2.2e-16
```

JUMPY VARIABLE TRAP: SUPPOSE THAT YOU WANT TO INCLUDE AN INDICATOR VARIABLE  $D_{ji}$  IN YOUR REGRESSION:

$$\text{FEMALE} = \begin{cases} 0, & \text{MALE} \\ 1, & \text{FEMALE} \end{cases}$$

RESPONSE OUTCOME VARIABLE  $\rightarrow$  WAGE

$$D_{ji} = \begin{cases} 0, & \text{MALE} \\ 1, & \text{FEMALE} \end{cases}$$

$$Y_i = \alpha_0 + \alpha_1 D_{ji} + u_i$$

$\frac{1}{\alpha_1}$  = AVERAGE WAGE DIFFERENCE B/W MALES AND FEMALES

$$Y_i = \beta_0 + \beta_1 D_{ji} + u_i$$

$\beta_1$  = AVERAGE WAGE

DIFFERENCE

B/W FEMALES

AND MALES

$$\frac{1}{\alpha_1} = -\beta_1$$

$$Y_i = \beta_0 + \beta_1 D_{ji} + \beta_2 D_{ji} + u_i$$

DO WHAT IS THE ISSUE?

$$D_{1i} + D_{2i} = 1 \quad D_{1i} = D_{2i} - 1$$

PERFECT LINEAR  
FUNCTION

$$Y_i = \beta_0 + \beta_1 D_{1i} + \beta_2 (1 - D_{1i}) + u_i$$

$$Y_i = \beta_0 + \beta_1 D_{1i} + \beta_2 - \beta_2 D_{1i} + u_i$$

$$Y_i = (\beta_0 + \beta_2) + (\beta_1 - \beta_2) D_{1i} + u_i$$

↳ SINGLE PREDICTOR  
TWO PARAMETERS

**IMPOSSIBLE ESTIMATE**

```

> CPS1988$white <- 0
> CPS1988$white[CPS1988$ethnicity=='cauc'] <- 1
> m05 <- lm(data=CPS1988, formula=wage~education+white)
> summary(m05)

Call: lm(formula = wage ~ education + white, data = CPS1988)
      ↳ Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 D_{1i} + u_i -
```

Residuals:

Min	1Q	Median	3Q	Max
-792.0	-256.3	-57.1	173.6	18027.3

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-123.2581	14.2508	-8.649	<2e-16 ***
education	46.2504	0.8882	52.070	<2e-16 ***
white	133.1458	9.5332	13.967	<2e-16 ***

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 430.9 on 28152 degrees of freedom  
Multiple R-squared: 0.09724, Adjusted R-squared: 0.09718  
F-statistic: 1516 on 2 and 28152 DF, p-value: < 2.2e-16

```

> CPS1988$non.white <- 0
> CPS1988$non.white[CPS1988$white==0] <- 1
>
> mo6 <- lm(data=CPS1988, formula=wage~ education + non.white + white)
> summary(mo6)

call:
lm(formula = wage ~ education + non.white + white, data = CPS1988)

Residuals:
    Min      1Q  Median      3Q     Max 
-792.0  -256.3   -57.1   173.6 18027.3 

Coefficients: (1 not defined because of singularities)
              Estimate Std. Error t value Pr(>|t|)    
(Intercept)  9.8877   11.9672   0.826   0.409    
education   46.2504   0.8882  52.070   <2e-16 ***  
non.white -133.1458   9.5332 -13.967   <2e-16 ***  
white        NA        NA        NA        NA      
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 430.9 on 28152 degrees of freedom
Multiple R-squared:  0.09724, Adjusted R-squared:  0.09718 
F-statistic: 1516 on 2 and 28152 DF,  p-value: < 2.2e-16

```

*(white)* ~~REMOVED~~

## IMPERFECT multicollinearity

MAKES A BIG DIFFERENCE.

$X_1$  &  $X_2$  ARE HIGHLY CORRELATED  
BUT NOT A PERFECT FUNCTION OF EACH OTHER

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$$

THE OLS WILL WORK BUT POORLY

$\hat{\beta}_1, \hat{\beta}_2$  WILL BE VERY IMPRECISE  
(LARGE STANDARD ERRORS)

$\hat{\beta}_1$  → RELATIONSHIP WITH  $X_1; Y Y$

HOLDING  $X_2$ ; CONSTANT.

IF  $X_1$  &  $X_2$  ARE HIGHLY CORRELATED,  
THEN IT IS DIFFICULT TO ESTIMATE THE  
NET EFFECT  $\hat{\beta}_1$ . → ~~VERY LITTLE VARIATION~~  
~~ESTIMATE EXPLOD~~

RESULT →  $SE[\hat{\beta}_1] = \text{large}$

95% CI  $\hat{\beta}_1$  IS LARGE

HOW TO DETECT?

1 COR | 710.80 → NOT A GOOD SIGN

2<sup>2</sup> HIGH BUT INDIVIDUAL COEFF. ARE NOT  
STATISTICALLY SIGNIFICANT DUE TO LARGE  
STANDARD ERRORS.

HIGH VARIANCE INFLATION FACTOR :

→ DISCARDED

$$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_k X_{ki} + U_i$$

$$Y_{t,i} = \alpha_1 + \alpha_2 X_{j,i} + \alpha_3 X_{\theta,i} + \dots + \alpha_k X_{r,i} + \epsilon_i$$

$$JIF(\beta_1) = 1 / (1 \cdot R_1^2)$$

DO CALCULATE  
R<sup>2</sup>

$JIF(\beta_1) \rightarrow$  INDICATES SEVERE MULTICOLLINEARITY

How To Fix?

- DO NOTHING
- DROP REDUNDANT VARIABLE
- TRANSFORM MULTICOLLINEATE VARIABLES

GDP / PSP CORRECTED  $\rightarrow$  GDP PER CAPITA

LARGER SAMPLE ALSO HELPS.

```

>
> m07 <- lm(data=CASchools, formula=read~str+english+lunch+expenditure)
>
> summary(m07) → -0.203 + 1.96 × 0.286
call:
lm(formula = read ~ str + english + lunch + expenditure, data = CASchools)

Residuals:
    Min      1Q  Median      3Q     Max 
-31.181  -5.353  -0.028   5.247  31.735 
[ -0.76 + 0.36 ]  
95.1.

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 662.0903950  9.0716578 72.984 < 0.0000000000000002 *** 
str          -0.2031297  0.2860344  -0.710   0.478    
english      -0.2105934  0.0304482  -6.916   0.0000000000176 *** 
lunch         -0.5502143  0.0203238  -27.072 < 0.0000000000000002 *** 
expenditure   0.0046660  0.0008406   5.551   0.000000508161 *** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 8.544 on 415 degrees of freedom
Multiple R-squared:  0.8212,   Adjusted R-squared:  0.8195 
F-statistic: 476.5 on 4 and 415 DF,  p-value: < 0.0000000000000022

> cor(CASchools)
      read      str      english      lunch expenditure
read 1.0000000 -0.2465930 -0.69028587 -0.87880769  0.21792682
str  -0.2465930 1.0000000  0.18764237  0.13520340 -0.61998216
english -0.6902859  0.1876424  1.00000000  0.65306072 -0.07139604
lunch  -0.8788077  0.1352034  0.65306072  1.00000000 -0.06103871
expenditure 0.2179268 -0.6199822 -0.07139604 -0.06103871  1.00000000
> vif(m07)
      str      english      lunch expenditure
str 1.680861  1.779490  1.744355  1.630108
>

```

variable of interest: STR

control variables: % OF ENGLISH LEARNERS

% OF FREE LUNCH

EXPENDITURE PER STUDENT

control variables do not need to have a causal interpretation but still can be used.

# HYPOTHESIS TESTING & CONFIDENCE INTERVALS IN MULTIPLE REGRESSIONS

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_K X_{Ki} + u_i$$

→ CAN BE ANY NUMBER

$$H_0: \beta_j = \beta_{j,0} \quad j = 1, \dots, K$$

$$H_A: \beta_j \neq \beta_{j,0}$$

$$\hat{\beta}_j - E[\hat{\beta}_j]$$

CLT

→

$$\frac{\hat{\beta}_j - E[\hat{\beta}_j]}{\sqrt{JAR(\hat{\beta}_j)}} \sim N(0, 1)$$

$$t\text{-stat} = \frac{\hat{\beta}_j - \beta_{j,0}}{SE[\hat{\beta}_j]}$$

$$P\text{-JAR} = 2 \times \Phi(-|t\text{-stat}|)$$

If  $P\text{-val} < \text{CRITICAL}$  → Reject

$$\hat{\beta}_j \pm 1.96 \times SE[\hat{\beta}_j] \rightarrow 95\% \text{ CONFIDENCE INTERVAL}$$

## TEST OF JOINT HYPOTHESES:

NULL REQUIRES TO TEST THE HYPOTHESES JOINTLY

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \beta_4 X_{4i} + u_i$$

$\downarrow$        $\downarrow$        $\downarrow$        $\downarrow$

TEST SCORE      % SF IN EXpenditure      % EDUCABLE PER MUNICIPALITY      % EDUCABLE PERSON

NULL HYPOTHESIS: SCHOOL RESOURCES DON'T MATTER

ALTERNATIVE: THEY DO

$$H_0: \beta_1 = 0 \quad \text{and} \quad \beta_4 = 0$$

$$H_A: \text{either } \beta_1 \neq 0 \text{ or } \beta_4 \neq 0$$

$$G = 2$$

LOTS OF RESTRICTIONS  
IN A JOINT HYPOTHESIS

INDIVIDUAL T-TEST

WON'T WORK

WHY?

ASSUME  $\beta_1 = 0 \rightarrow$  AT 5% SIGNIFICANCE  
LEVEL 95% CHANCE OF FAILING TO REJECT.

ASSUME  $\rho_{44} = 0 \rightarrow$  AT 5% SIGNIFICANCE LEVEL QM.S. FAIL TO REJECT

$$\text{PROP. FAIL TO REJECT} = 0.05 \times 0.05 \\ = 0.0025$$

$$\text{PROP. REJECT} \quad \hat{\rho}_{11} = 0 \times \hat{\rho}_{44} = 0 \\ 1 - 0.0025 = 0.9975.$$

POTENTIALLY A BIGGER

PROBLEM:

$$\hat{\rho}_{11} = 0 \times \hat{\rho}_{44} = 0$$

REJECT TOO OFTEN  
IF YOU USE INDIVIDUAL  
t-tests

ARE DEPENDENT.

$$\text{SCHWARTZ TEST} = \frac{1}{2} \left[ \frac{-\hat{t}_{11}^2 + \hat{t}_{44}^2 - 2\hat{\rho}_{11,44}^1}{1 - \hat{\rho}_{11,44}^{1,2}} \right]$$

$$\text{ESTIMATED CORRELATION B/W} \quad \hat{\rho}_{11,44}^1 \rightarrow \text{ACOUNTS} \quad \frac{\text{COR}(x,y)}{\text{VAR}(x)}$$

$\hat{t}_1, \hat{t}_2$   $\hat{\rho}_{11,44}^1$  → ACCOUNTS FOR CORRELATION B/W  $X_1 \text{ & } X_2$

IF  $F$  IS LARGE, REJECT THE NULL

IF  $t, \chi^2$  ARE INDEPENDENT

$$\bar{F} = \frac{1}{2} [t_1^2 + t_4^2] \rightarrow \text{ONLY IF } t, \chi^2 \text{ ARE INDEPENDENT}$$

+ LARGE SAMPLE  $\rightarrow \bar{F} \sim \chi_q^2$

If  $q = 1 \rightarrow 5\% \text{ CRITICAL-}F \approx 3.84$

$q = 2 \rightarrow 5\% \text{ CRITICAL-}F \approx 3.50$

$q = 3 \rightarrow 5\% \text{ CRITICAL-}F \approx 2.60$

```
> LinearHypothesis(mo7, c("str=0", "expenditure=0"))
Linear hypothesis test
```

Hypothesis:

str = 0  
expenditure = 0

Model 1: restricted model

Model 2: read ~ str + english + lunch + expenditure

	Res.Df	RSS	Df	sum of sq	F	Pr(>F)
1	417	34575				
2	415	30292	2	4283.2	29.34	0.0000000000001207 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

>

SIMPLE FORMULA: HOLDS IN SPECIFIC CONDITIONS  
(HOMOSCEDASTICITY) BUT USEFUL TO UNDERSTAND  
WHAT F-TEST IS DOING.

PUT TWO PROGRESSIONS:

- RESTRICTED  $\rightarrow \hat{\beta}_1 = 0, \hat{\beta}_4 = 0$

- UNRESTRICTED  $\rightarrow$  FULL MODEL

CALCULATE  $R^2$  FROM BOTH

$$F = \frac{(R^2_{\text{UNR}} - R^2_{\text{RES}}) / q}{(1 - R^2_{\text{UNR}}) / (n - k - 1)}$$

IF DIFFERENCES IN  
 $R^2$  IS BIG,  
} F IS LARGE  
MORE LIKELY  
TO REJECT  
H<sub>0</sub>

IF DIFF IS NOT BIG, THEN MAY BE  
THE COEFS JOINTLY DO NOT ADD MUCH  
PREDICTION POWER TO THE MODEL.

$$R^2_{\text{UNR}} = 0.8212 \quad q = 2 \quad (n - k - 1) = 420 - 4 - 1 = 415$$

$$R^2_{\text{RES}} = 0.7959 \quad F = \frac{(0.8212 - 0.7959)}{(1 - 0.8212) / 415}$$

$$= 29.36 \quad \text{reject the null} > 2.30$$

REJECT THE NULL. CLASSROOM RESOURCES ARE USEFUL

```
> summary(mo8)

Call:
lm(formula = read ~ str + english + lunch + expenditure, data = CASchools)

Residuals:
    Min      1Q  Median      3Q     Max 
-31.181  -5.353  -0.028   5.247  31.735 

Coefficients:
            Estimate Std. Error t value     Pr(>|t|)    
(Intercept) 662.0903950  9.0716578 72.984 < 0.000000000000002 ***  
str          -0.2031297  0.2860344 -0.710     0.478    
english      -0.2105934  0.0304482 -6.916     0.0000000000176 ***  
lunch         -0.5502143  0.0203238 -27.072 < 0.000000000000002 ***  
expenditure   0.0046660  0.0008406  5.551     0.000000508161 ***  
---
Signif. codes:  0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 8.544 on 415 degrees of freedom
Multiple R-squared:  0.8212, Adjusted R-squared:  0.8195 
F-statistic: 476.5 on 4 and 415 DF,  p-value: < 0.000000000000022

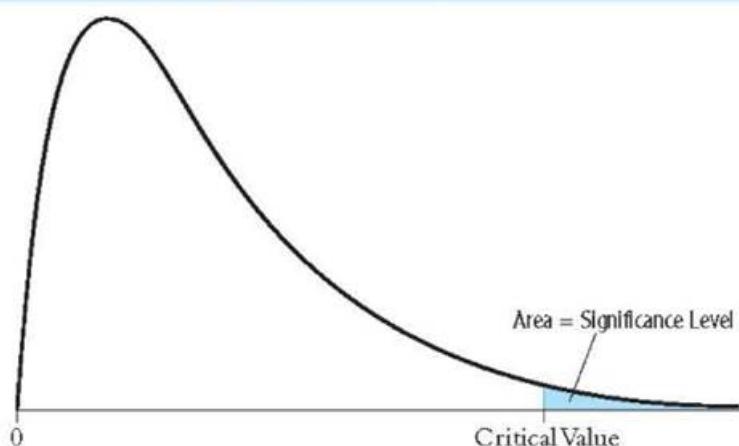
> summary(mo9)

Call:
lm(formula = read ~ english + lunch, data = CASchools)

Residuals:
    Min      1Q  Median      3Q     Max 
-27.130  -5.718  0.185   5.572  33.894 

Coefficients:
            Estimate Std. Error t value     Pr(>|t|)    
(Intercept) 683.22276  0.86649 788.497 < 0.000000000000002 ***  
english      -0.22313  0.03212 -6.946     0.0000000000145 ***  
lunch        -0.55327  0.02166 -25.547 < 0.000000000000002 ***  
---
Signif. codes:  0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 9.106 on 417 degrees of freedom
Multiple R-squared:  0.7959, Adjusted R-squared:  0.7949 
F-statistic: 813.1 on 2 and 417 DF,  p-value: < 0.000000000000022
```

**TABLE 4** Critical Values for the  $F_{m,\infty}$  Distribution

Degrees of Freedom	Significance Level		
	10%	5%	1%
1	2.71	3.84	6.63
2	2.30	3.00	4.61
3	2.08	2.60	3.78
4	1.94	2.37	3.32
5	1.85	2.21	3.02
6	1.77	2.10	2.80
7	1.72	2.01	2.64
8	1.67	1.94	2.51
9	1.63	1.88	2.41
10	1.60	1.83	2.32
11	1.57	1.79	2.25
12	1.55	1.75	2.18
13	1.52	1.72	2.13
14	1.50	1.69	2.08
15	1.49	1.67	2.04
16	1.47	1.64	2.00
17	1.46	1.62	1.97
18	1.44	1.60	1.93
19	1.43	1.59	1.90
20	1.42	1.57	1.88
21	1.41	1.56	1.85
22	1.40	1.54	1.83
23	1.39	1.53	1.81
24	1.38	1.52	1.79
25	1.38	1.51	1.77
26	1.37	1.50	1.76
27	1.36	1.49	1.74
28	1.35	1.48	1.72
29	1.35	1.47	1.71
30	1.34	1.46	1.70

This table contains the 90<sup>th</sup>, 95<sup>th</sup>, and 99<sup>th</sup> percentiles of the  $F_{m,\infty}$  distribution. These serve as critical values for tests with significance levels of 10%, 5%, and 1%.

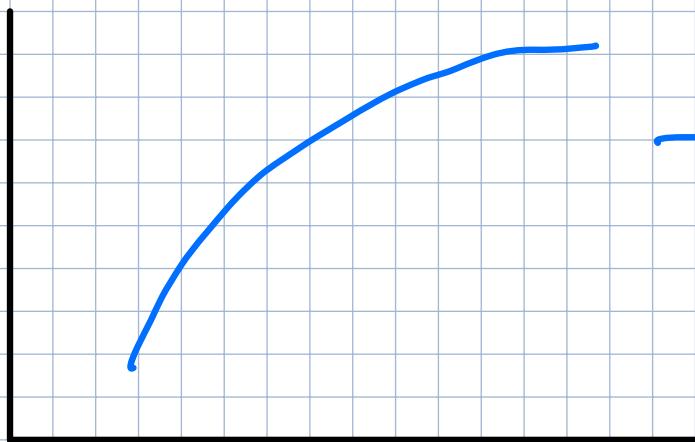
## NON-LINEAR MODELS

FORM OF THE RELATIONSHIP B/W  $Y$  &  $X$  ARE

NON-LINEAR

↳ THE IMPACT OF  $X$  ON  $Y$  DEPENDS ON  
THE LEVEL OF  $X \rightarrow ?$  IS NOT CONSTANT.  
VALUE

AND IS A FUNCTION OF  $X$ . SOMETIMES THEORETICALLY  
JUSTIFIED.



→ WAGE ↑ FASTER  
EARLY IN THE  
CAREER

GENERAL NON-LINEAR FORMULATION  
REGRESSION FUNCTION

$$Y_i = f(X_{1i}, X_{2i}, \dots, X_{ri}) \quad i = 1, 2, \dots, n$$

↳ NON-LINEAR

SOME CASES: CAN USE OLS AFTER TRANSFORMATION

# POLYNOMIALS

↳ POPULATION - PREDICTION FUNCTION CAN BE

PROXIED BY A QUADRATIC, CUBIC OR HIGHER ORDER POLYNOMIAL

**LOGARITHMIC TRANSFORMATION**: X, Y OR BOTH TRANSFORMED TO LOG. WHICH MAKES THEM EASIER TO INTERPRET

**POLYNOMIALS:**

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \dots + \beta_k X_i^k + u_i$$

USING ONE X

ALL REGRESSIONS

ARE POWER OF X

$$Y_i = f(X_i)$$

THE MODEL IS LINEAR IN PARAMETERS SO YOU CAN USE OLS. COEFFICIENTS ARE HARD TO INTERPRET

1.  $Y_i = \beta_0 + \beta_1 X_i + u_i \rightarrow \text{LINEAR}$

2.  $Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 \rightarrow \text{QUADRATIC}$

3.  $Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \beta_3 X_i^3 + u_i \rightarrow \text{CUBIC}$

# HANDOUT ANSWERS:

1.  $\hat{\beta}_1 = 1.815 \quad P\text{-TAC} < 0.01$

$SE[\hat{\beta}_1] = 0.091$

\$1000 ↑ IN INCOME WITH A 1.82 POINTS  
INCREASE IN TEST SCORES.

2. TEST SCORES ↑ WITH INCOME BUT  
THE IMPACT OF INCOME IS SMALLER AS THE  
LEVEL OF INCOME ↑ BUT THE EFFECT OF  
INCOME IS SMALLER AS THE LEVEL OF  
INCOME ↑

$\hat{\beta}_1 > 0 \quad SE[\hat{\beta}_1] = 0.210 \quad \hat{\beta}_1 = 0.173$

$\hat{\beta}_2 < 0 \quad SE[\hat{\beta}_2] = 0.006 \quad \hat{\beta}_2 = -0.036$

$$Y_i = \beta_0 + \hat{\beta}_1 X_i + \hat{\beta}_2 X_i^2 + U_i$$

$$\frac{dy_i}{dx_i} = \hat{\beta}_1 + 2\hat{\beta}_2 X_i$$

+      -      x  
                ↑ LEVEL OF

$$\bar{x} \rightarrow \frac{dy_i}{dx_i} = 3.473 + [2 \times (-0.036) \times 15.32] \\ = 3.473 - 1.103 \\ = 2.37$$

A \$100 IN INCOME WHEN INCOME  $\approx \$15,317$   
IS ASSOCIATED WITH A 2.37 POINTS INCREASE  
IN TEST SCORES.

$$\left. \begin{array}{l} \frac{dy_i}{dx_i} \text{ When } x = \$15,000 \\ x = \$15,000 \\ x = \$30,000 \end{array} \right\} \begin{array}{l} 2.11 \\ 2.39 \\ 1.31 \end{array}$$

**CAUTION:** NEVER EXTRAPOLATE (MAKE PREDICTIONS)  
OUTSIDE THE DATA RANGE  
OF  $x$ .

T-TEST: UNRESTRICTED MODEL:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + u_i$$

RESTRICTED

$$y_i = \beta_0 + \beta_1 x_i + u_i$$

$$H_0: \beta_2 = 0 \quad \text{vs} \quad \beta_2 \neq 0$$

$H_A:$  either  $\beta_2 \neq 0$  or  $\beta_3 \neq 0$

$$F\text{-stat} = 15.762 \quad p\text{-value} < 0.01$$

reject  $H_0 \rightarrow$  model should be non-linear.

$$H_0: \beta_3 = 0 \quad Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + u_i$$

$$H_A: \beta_3 \neq 0 \quad Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \beta_3 X_i^3 + u_i$$

$$F\text{-stat} = 0.3768 \quad p\text{-value} = 0.539 > 0.01$$

FAIL TO REJECT THE NULL.  $\rightarrow$  OVER-FITTING

$X^3$  SHOULD NOT BE IN THE MODEL.

↳ EQUIVALENT TO AN INDIVIDUAL T-TEST ON  $\beta_3$

↳ SUMMARY: POLYNOMIAL MODELS CAN BE ESTIMATED USING LS

↳ INDIVIDUAL COEFFICIENTS ARE HARD TO INTERPRET.

↳ BEST OPTION: TAKE PREDICTION + EVALUATE THE IMPACT AT A SPECIFIC  $X$ .

↳ DECIDE ON THE APPROPRIATE FORM USING F-TEST OR T-TEST.

# LOGARITHMIC FUNCTIONS OF Y & X

$\log(x) = \ln(x)$   $\rightarrow$  THE NATURAL

LOG OF X

IS VERY USEFUL FOR MODELLING RELATIVE (%) CHANGES.

$$x + \Delta x - x = \Delta x$$

$$\ln(x + \Delta x) - \ln(x) = \ln\left(\frac{x + \Delta x}{x}\right)$$

$$\approx \frac{\Delta x}{x} \rightarrow \text{RELATIVE CHANGE}$$

1<sup>o</sup> LINEAR-COF  $\Delta x/x \cdot 100 \rightarrow \%$   
CHANGE

$$Y_i = \beta_0 + \beta_1 U(X_i) + u_i$$

2<sup>o</sup> LOG-LINEAR

$$U(Y_i) = \beta_0 + \beta_1 X_i + u_i$$

$$3^o \text{ LOG-COF } U(Y_i) = \beta_0 + \beta_1 U(X_i) + u_i$$

$\beta_1$  IS INTERPRETED DIFFERENTLY  
IN EACH CASE. USE "BEFORE VS  
AFTER" RULE

### 1. LINEAR LOG

BEFORE :  $y_i = \beta_0 + \beta_1 \ln(x_i) + u_i$

AFTER :  $y_i + \Delta y = \beta_0 + \beta_1 \ln(x_i + \Delta x) + u_i$

---

$$\Delta y = \beta_1 [\ln(x_i + \Delta x) - \ln(x_i)]$$

$$\Delta y = \beta_1 \frac{\Delta x}{x_i} \quad \Delta \rightarrow \text{SMALL}$$

$$\beta_1 = \frac{\Delta y}{\Delta x / x}$$

$\beta_1 \rightarrow$  A 1% INCREASE IN  $X$  IS  
ASSOCIATED WITH A  $\Delta y / 100$   
CHANGE IN  $Y$ .

## 2. LOG-LINEAR FORM

BEFORE :  $l_1(y) = \beta_0 + \beta_1 x_i + u_i$

AFTER :  $l_1(y + \Delta y) = \beta_0 + \beta_1 (x_i + \Delta x) + u_i$

$$l_1(y + \Delta y) - l_1(y) = \beta_1 \Delta x$$

$$\frac{\Delta y}{y} = \beta_1 \Delta x$$

$$\beta_1 = \frac{\Delta y / y}{\Delta x}$$

$\beta_1 \rightarrow$  % CHANGE IN Y ASSOCIATED

WITH A 1 UNIT INCREASE IN X

## 3. LOG-WG FORM

$$l_1(y_i) = \beta_0 + \beta_1 l_1(x_i) + u_i$$

$$l_1(y_{i+1}) = \beta_0 + \beta_1 l_1(x_{i+1}) + u_i$$

$$l_1(y_{i+1}) - l_1(y_i) = \beta_1 [l_1(x_{i+1}) - l_1(x_i)]$$

$$\beta_1 = \frac{\Delta Y / Y}{\Delta X / X}$$

$\beta_1 \rightarrow$  elasticity

∴  $\Delta$  IN  $Y$  ASSOCIATED WITH A

$\Delta$  IN  $X$

→ ESTIMATED USING OLS

→ HYPOTHESIS TEST + CONFIDENCE

→ INTERVALS INTERPRETED AS USUAL

→ WHEN TO USE WFF FORM?

CONINUED READERS → SKewed —  
DISTRIBUTION OUTCOMES

- PLOT RELATIONSHIPS

- DIAGNOSTIC TESTS

$$E(Y_i) = \beta_0 + \beta_1 D_i + u_i$$

$$D_i = 0 \rightarrow E(Y_i | D_i = 0) = \beta_0 + u_i$$

$$D_i = 1 \rightarrow E(Y_i | D_i = 1) = \beta_0 + \beta_1 + u_i$$

$$P_i = \mu \left( \frac{Y_i | D_i=1}{Y_i | D_i=0} \right)$$

$\therefore \Delta \ln Y \text{ WHEN } D_i=0 \rightarrow D_i=1$

## HANDOUT ANSWERS :

Question 1.

Interpret the coefficients on the following regression results.

```
> data(CASchools)
>
> CASchools$ln.income<- log(CASchools$income)
> CASchools$ln.math <- log(CASchools$math)
> CASchools$str <- CASchools$students/CASchools$teachers
> CASchools$small <- 0
> CASchools$small[CASchools$str<20] <- 1
>
>
> mo1 <- lm(data=CASchools, formula=math ~ ln.income)
> mo2 <- lm(data=CASchools, formula=ln.math ~ income)
> mo3 <- lm(data=CASchools, formula=ln.math ~ ln.income)
> mo4 <- lm(data=CASchools, formula=math ~ small)
> mo5 <- lm(data=CASchools, formula=ln.math ~ small)
>
> stargazer::stargazer(mo1,mo2,mo3,mo4,mo5, type='text')

=====
Dependent variable:

```

	math (1)	ln.math (2)	math (4)	ln.math (5)
ln.income	34.664*** (1.610)		0.053*** (0.002)	
income		0.003*** (0.0001)		
small			5.599*** (1.830)	0.008*** (0.003)
Constant	561.661*** (4.304)	6.440*** (0.002)	6.342*** (0.007)	650.156*** (1.380)

```

Observations          420      420      420      420      420
R2                   0.526     0.481     0.522     0.022     0.022
Adjusted R2           0.525     0.480     0.521     0.020     0.019
Residual Std. Error (df = 418) 12.928    0.021    0.020    18.570    0.028
F Statistic (df = 1; 418)   463.806*** 387.338*** 456.883*** 9.364*** 9.256***
```

```

Note: *p<0.1; **p<0.05; ***p<0.01
>
```

- (1) : 1 PERCENT IN INCOME IS ASSOCIATED WITH A 0.347 POINTS INCREASE IN MATH SCORES.
- (2) : \$1000 IN INCOME IS ASSOCIATED WITH A 0.8% IN TEST SCORES
- (3) : 1%. IN INCOME IS ASSOCIATED WITH A 0.053% INCREASE IN TEST SCORES.
- (5) : ON AVERAGE SMALL CLASSROOMS PERFORM 0.8% HIGHER.

**WARNING!**

APPROXIMATION ONLY WORKS WHEN :

$\Delta Y / Y$  IS SMALL IF YOU WANT THE PERCENT CHANGE :

$$M Y_i = \beta_0 + \beta_1 X_i + u_i$$

$$M(Y_i + \Delta Y_i) = \beta_0 + \beta_1 (X_i + \Delta X_i) + u_i$$

$$= \beta_1 \Delta X_i$$

$$\ln\left(\frac{\Delta Y_i + Y_i}{Y_i}\right) = \beta_0 + \beta_1 \Delta X_i$$

} regression  
output

$$\ln\left(\frac{\Delta Y_i + 1}{Y_i}\right) = \beta_0 + \beta_1 \Delta X_i$$

$\underbrace{\quad}_{\Delta Y_i}$

$$\approx \frac{\Delta Y_i}{Y_i} = \therefore \Delta Y_i$$

ONLY MATTERS  
WHEN  $\beta_1$  IS  
LARGE.

~~EXACT~~  $\rightarrow \therefore \Delta Y_i =$

$$\ln Y_{i0} = \beta_0 + \beta_1 X_{i0} + u_i; \quad X = \exp(\ln(x))$$

$$\ln Y_{ii} = \beta_0 + \beta_1 X_{ii} + u_i$$

$$\frac{\Delta Y_i}{Y_i} = \frac{Y_{ii} - Y_{i0}}{Y_{i0}} = \frac{Y_{ii}}{Y_{i0}} - 1$$

$$\hookrightarrow \exp\left(\ln\left(\frac{Y_{ii}}{Y_{i0}}\right)\right) - 1$$

$$\hookrightarrow \exp(\beta_1 X_{ii} - \beta_1 X_{i0}) - 1 \rightarrow \exp(\underbrace{\Delta X \beta_1}_{-1})$$

$$\text{IF } \hat{\beta}_1 = 0.3 \rightarrow \therefore \Delta Y = 0.35$$

REGRESSION WITH A BINARY DEPENDENT VAR:

$$Y = \{0, 1\} \quad Y \rightarrow \text{CONTINUOUS}$$

LATEST SCORES, TRAFFIC FATALITY RATE

GIVING INTO COLLEGE:  $\hat{\beta}_1 = 20.13$  SMOKING:  $\hat{\beta}_1 = 20.13$

OBESITY =  $\hat{\beta}_1 = 20.13$  MORTGAGE APPLICATION  $\hat{\beta}_1 = 20.13$

$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + U_i$   
 HOW TO INTERPRET  $\hat{\beta}_1$  WHEN  $X_i$  IS  
 CONTINUOUS &  $Y_i$  IS BINARY.

LOGISTIC PROBABILITY MODEL:

$$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + U_i$$

$$E[Y_i | X_i] = 1 \times \text{PR}(Y_i = 1 | X_i) + 0 \times \text{PR}(Y_i = 0 | X_i)$$

$$E[Y_i | X_i] = \hat{P}_R(Y_i | X_i)$$

↳ **HANDOUT** : PROBABILITY OF DENIAL CONDITIONAL ON PI PATHO

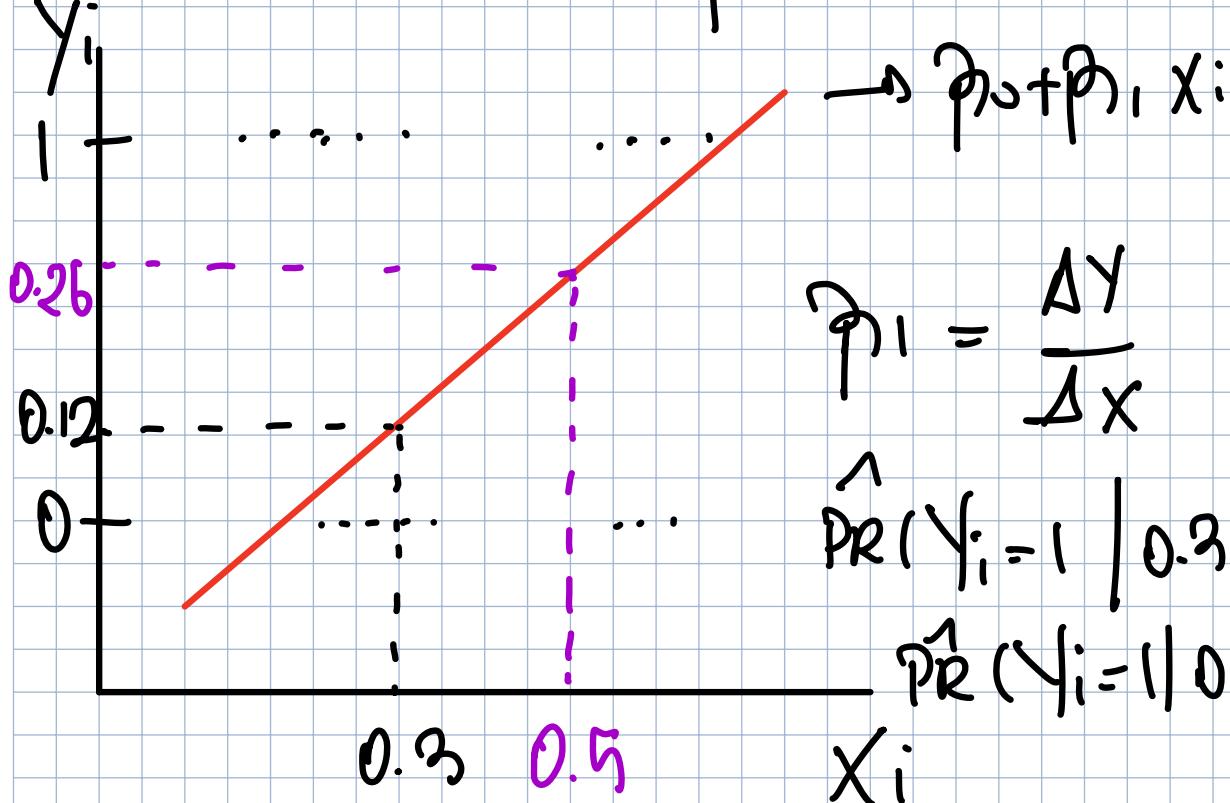
$$Y_i = \beta_0 + \beta_1 X_i + U_i$$

$$E[Y_i | X_i] = \hat{P}_R(Y_i = 1 | X_i) = E[\beta_0 + \beta_1 X_i + U_i | X_i]$$

$$E[\beta_0] + E[\beta_1 X_i | X_i] + E[U_i | X_i] = 0$$

$$\hat{P}_R(Y_i = 1 | X_i) = \beta_0 + \beta_1 X_i$$

$$\hat{P}_R(Y_i = 1 | X_i) = \hat{\beta}_0 + \hat{\beta}_1 X_i \rightarrow \text{SAMPLE}$$



$$\hat{\beta}_1 = \frac{\hat{P}_R(Y_i=1|0.5) - \hat{P}_E(Y_i=1|0.3)}{0.2}$$

PROBABILITY OF DENIAL GOES UP BY 0.14

OR 14 PERCENTAGE POINTS AS THE  
PI RATIO INCREASED BY 0.2.

$$1^0 \quad \begin{aligned} & -0.07991 + 0.1 \times 0.60353 = 16.2\% \\ & -0.07991 + 0.3 \times 0.60353 = 10.1\% \end{aligned}$$

$$\underline{0.060353}$$

IF PI ↑ BY 0.1

THE PROB. DENIAL GOES UP BY 6 PERCENTAGE  
POINTS.

$$2^0 \quad Y_i = \beta_0 + \beta_1 D_i + u_i$$

$\downarrow$

$$\{0,1\}$$

$\downarrow$

$$\{0,1\}$$

$$E[Y_i | D_i] = P_R(Y_i=1 | D_i) = \beta_0 + \beta_1 D_i$$

$$PR(Y_i=1 | D_i=0) = \beta_0 \rightarrow \text{PROBABILITY OF DENIAL IN BEING NON-BLACK}$$

$$PR(Y_i=1 | D_i=1) = \beta_0 + \beta_1 \rightarrow \text{PROBABILITY OF DENIAL OF BEING BLACK}$$

$$\therefore PR(Y_i=1 | D_i=0) = 9.1\%.$$

$$\overline{PR(Y_i=1 | D_i=1)} = 28.4\%.$$

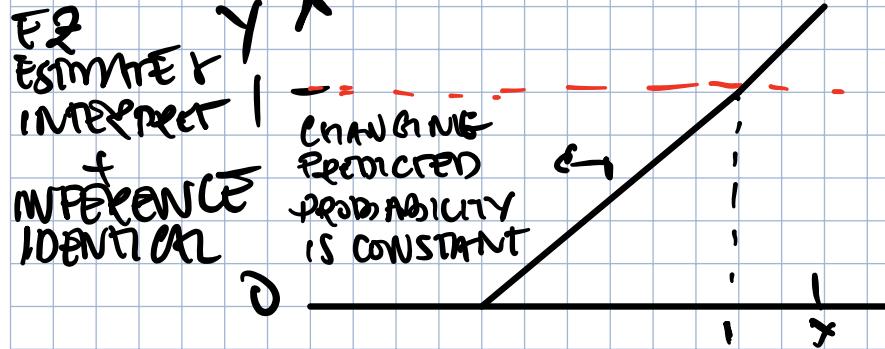
---


$$\hat{\beta}_1 = 0.191$$

$\hat{\beta}_1 = 0.191 \rightarrow$  BEING BLACK IS ASSOCIATED WITH A 19.1% DIFFERENCE DECLINE IN MORTGAGE APPROVAL.

3rd MODEL: BEING BLACK IS ASSOCIATED WITH A 17.1% DECLINE IN APPROVAL HAVING THE RATE CONSTANT.

**ISSUES:**  $PR(Y=1|x)$  INCREASING WITH  $x$   $\beta_1 > 0$



$\rightarrow$  NOT PRACTICABLE  
EXAMPLE: 20m

$PR(Y=1|x) = \text{LINEAR FUNCTION OF } x$   
 $(0, 1)$   
CONSEQUENCE!