

01/17/2021

HEY THERE! SPRING 2021 LECTURE NOTES FOR  
ECON62 WILL BE ON THIS DOCUMENT, PLEASE  
GET USED TO MY TERRIBLE HANDWRITING

01/21/2021

→ EXAMS AFTER THE SURVEY

LESSON SCAN:

- SELF-INTRODUCTION
- ECONOMETRICS:
  - ECON → ECONOMIC THEORY
  - METRICS → TEST IT WITH DATA.
- COURSE LOGISTICS +
- EXPECTATIONS +
- R + R STUDIO + GITHUB

SHOW CASE  
HOMEWORK ASSIGNMENT

DATA TALK

WHY DO WE NEED THESE?

BUSINESSES  
ARE CHANGING

→ MORE & MORE  
DATA DRIVEN

↳ DATA COLLECTION (MOBILE DEVICES)

STORING (CLOUD)

MANAGE (SQL, ETC)

→ NOT  
TABLES.

↳ MACHINE ALGORITHMS BECAME SMARTER

# EX. NATURAL LANGUAGE PROCESSING.

- PREDICTION → WHAT PREDICTS SALES?
- CAUSAL CHAIN → DOES SOCIAL MEDIA PRESENCE INCREASE SALES?

WHAT SHOULD A BUSINESS STUDENT LEARN?

LEARN BOTH AT AN OK LEVEL.  
(UNDERSTAND + BASELINE)  
KNOWLEDGE → MECHANICS OF  
REGRESSION ANALYSIS &  
DATA

LEARN ~~TEST WAYS TO~~  
~~INTER CAUSALITY~~  
RANDOMIZED  
EXPERIMENTS

01/25/2020 LECTURE 1

## LESSON PLAN:

### PREVIEW OF PROBABILITY - CH 2

#### RANDOM VARIABLES

##### PROBABILITY DISTRIBUTIONS

$$\cdot \mathbb{P}[X], \text{VAR}[X] = 0, \text{SD}_X^2, \text{SD}_X$$

##### STANDARDIZED RANDOM VAR.

$$\frac{X - \mu}{\text{SD}_X} \quad (\text{SD}_X)$$

##### APPLICATION IN R

#### RANDOM VARIABLE: NUMERICAL SUMMARY OF

-> FLIPPING A COIN IS AN EVENT WITH A RANDOM COMPONENT

-> YOUR WAGE AT AGE 40

- DISCRETE
- CONTINUOUS

INTERVALS

$\sum p_i$  → SUM OF PROBABILITIES FOR ALL THE

Possible Outcomes  
= 1

EXPECTED VALUE  $\rightarrow$  LONG RUN AVERAGE

RANDOM VARIABLE  $Y \rightarrow$   $\begin{cases} Y_i & \text{EXPECTED} \\ \mathbb{E}[Y] & \text{VALUE} \end{cases}$

$$\mathbb{E}[Y] = Y_1 \cdot P_1 + Y_2 \cdot P_2 + \dots + Y_k \cdot P_k$$

K OUTCOMES  $\downarrow$  DISCRETE PROBABILITIES

IF THE OUTCOME IS BINARY

DUMM - BERNOULLI - RANDOM VAR

IN GENERAL

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

$$\mathbb{P}(Y=0) = (1-p)$$

$$\mathbb{P}(Y=1) = p$$

$$\begin{aligned} \mathbb{E}[Y] &= 0 \times (1-p) + 1 \times p && \text{SURVIVE} \\ &= p \end{aligned}$$

Covid: MORTALITY  $Y \rightarrow \{0, 1\}$

$$\mathbb{P}(Y=1) = 0.005$$

$\frac{1}{200}$

$$\mu_Y = 0 \times 0.995 + 1 \times 0.005 \\ = 0.005$$

**VARIANCE:**  $\text{Var}(Y) = \sigma_Y^2$  - A MEASURE OF THE SPREAD OF POSSIBLE OUTCOMES AROUND THE EXPECTED VALUE

$$\sigma_Y^2 = (Y_1 - \mu_Y)^2 P_1 + (Y_2 - \mu_Y)^2 P_2 + \dots + (Y_k - \mu_Y)^2 P_k$$

$$\sigma_Y^2 = \sum_i (Y_i - \mu_Y)^2 P_i$$

$$\sigma_Y^2 \geq 0 \quad \sigma_Y = \sqrt{\sigma_Y^2}$$

IF  $Y$  IS BINARY

$$\Pr(Y=1) = p \quad \mu_Y = p$$

$$\Pr(Y=0) = (1-p)$$

$$\begin{aligned} \text{Var}(Y) &= (0-p)^2 \times (1-p) + (1-p)^2 \cdot p \\ &= p(1-p) [p - (1-p)] \\ &= p(1-p) \end{aligned}$$

$$\text{COVID MORTALITY : } \text{JAR}(Y) = 0.005 \times 0.995$$

$$\bar{J}_Y^2 = 0.004975$$

$$\bar{J}_Y \approx 0.0705$$

$\bar{J}$   $\rightarrow$  CASE FATALITY RATIO

TWO RANDOM VARIABLES

$\text{PR}(Y=y)$   $\rightarrow$  PROBABILITY THAT A RANDOM VARIABLE  $Y$  IS EQUAL TO A SPECIFIC VALUE  $y$

$\text{PR}(Y=1)$   $\rightarrow$  MORTALITY RISK (IN THE PREVIOUS EXAMPEL)

JOINT PROBABILITY  
DISTRIBUTION OF  
 $X$  &  $Y$

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

W PROB. THAT  $X=x$  AND

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

$\rightarrow$  SENIOR VS.  
NON SENIOR

$$Y = Y$$

$$\text{PR}(X=0, Y=1)$$

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

$\rightarrow$  MORTALITY

$$\text{PR}(X=1, Y=1)$$

$$\text{PR}(Y=y) = \sum_{i=1}^n \text{Pr}(X=x_i; Y=y)$$

(MARGINAL PR.)

(JOINT PROBABILITY)

$\bar{X}$ :

	Rain ( $X = 0$ )	No Rain ( $X = 1$ )	Total
Long commute ( $Y = 0$ )	0.15	0.07	0.22
Short commute ( $Y = 1$ )	0.15	0.63	0.78
Total	0.30	0.70	1.00

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

$$PR(X=0, Y=1) = 0.15$$

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

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


---

MUTUALLY  
EXCLUSIVE

$$PR(Y=1) = \sum_{i=1}^2 PR(X=x_i, Y=1)$$

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

$$= 0.78$$

$$PR(Y=0) = 0.22$$

$$PR(X=0) = 0.30$$

$$PR(X=1) = 0.70$$

BAYES  
THEOREM

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

PR(SHORT COMMUTE | RAIN)

$$PR(Y=1 | X=0) = \frac{PR(Y=1, X=0)}{PR(X=0)}$$

$$\text{or } \frac{0.15}{0.30}$$

$$= 50\%$$

01/28/2021

LESSON PLAN:

- REVIEW joint prob.
- CONDITIONAL EXPECTED VALUE
- LAW OF ITERATED EXPECTATIONS
- INDEPENDENCE, COVARIANCE,

CHAPTER 2

CORRELATION

SOLVE THE HANDOUT

QUESTIONS

MAY BE: INTRODUCE RANDOM SAMPLING

TWO RANDOM VARIABLES  $\rightarrow$  joint probability  
distribution

$\Pr(X=x, Y=y)$   $\rightarrow$  notation for joint probability

$$\Pr(Y=y) = \sum_{i=1}^l \Pr(X=x_i, Y=y)$$

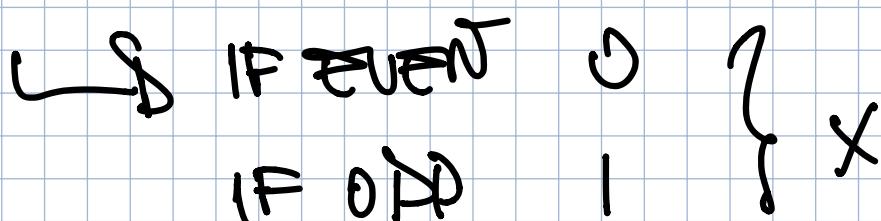
BAYES THEOREM:

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

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

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

roll A DIE  $\rightarrow$  outcome  $Y$



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

$$\begin{aligned} &= 1 \times 1/6 + 2 \times 1/6 + \dots + 6 \times 1/6 \\ &= 3.5 \end{aligned}$$

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

$$= 1 \times 1/3 + 2 \times 0 + 3 \times 1/3 + 0$$

$$+ 5 \times 1/3 + 0$$

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

LAW OF ITERATED EXPECTATIONS

$$E[Y] = \sum_{i=1}^l E[Y|X=x_i] \Pr(X=x_i)$$

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

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

$$E[Y|X=1] = 1 \times 0 + 2 \times 1/3 + 3 \times 0 + 5 \times 1/3 + 0 + 6 \times 1/3$$

$$= 4$$

$$E[Y] = 4 \times 1/2 + 3 \times 1/2$$

$$= 3.5$$

SHORT NOTATION FOR  
LAW OF ITERATED  
EXPECTATIONS

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

$$\text{INDEPENDENCE: } \Pr(Y=y|X=x) = \Pr(Y=y)$$

US KNOWING X PROVIDES  
NO INFORMATION ON THE

Possible outcome of Y

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

$$\Pr(Y=y) \Pr(X=x) = \Pr(X=x, Y=y)$$

ONLY HOLDS IF

$$X \text{ AND } Y \text{ ARE INDEPENDENT} \quad \text{COR}(x,y) = \overline{\sigma}_{xy} = \mathbb{E}[(x - \mu_x)(y - \mu_y)]$$

$$= \sum_i \sum_j (x_j - \mu_x)(y_i - \mu_y) \Pr(X=x_i, Y=y_i)$$

$$1, 0, -1$$

ONLY INDICATES

DIRECTION

$$\text{COR}(x,y) = \frac{\overline{\sigma}_{xy}}{\sqrt{\sigma_x^2 \sigma_y^2}}$$

$$-1 \leq \text{COR}(x,y) \leq 1$$

$\downarrow$

PERFECT  
NEGATIVE COR.

$\uparrow$   
 $\uparrow$   
 $\uparrow$

US PERFECT  
POSITIVE  
COR.

NO CORRELATION  
UNCORRELATED

### QUESTION 3:

$$\Pr(Y=1) = \sum_x \Pr(X=x, Y=1)$$

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

$$= 50\%$$

$$\Pr(A=1) = 50\%.$$

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

$$(0.35 / 0.5) = 70\%$$

$$E[M | A=0] = \sum_m \Pr(M=m | A=0)$$

$$0 \times 0.70 + 1 \times 0.13 + \dots = 0.56$$

$$0 \times 0.90 + 1 \times 0.17 + \dots = 0.17$$

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

$$= E[M | A=0] \times \Pr(A=0) +$$

$$E[M | A=1] \times \Pr(A=1)$$

$$= 0.56 \times 0.5 + 0.17 \times 0.5 = 0.35 = E[M]$$

9/11/2021 - LESSON PLAN - CH 8 REVIEW OF STATS

- RANDOM SAMPLING AND DISTRIBUTION OF THE SAMPLE MEAN
- NORMAL DISTRIBUTION

POPULATION DATA  $\rightarrow$  DOES NOT EXIST (IN MOST CASES)

( $\mu_x, \sigma_x$ )  $\rightarrow$  ALL UNKNOWN

$E[x]$   $\sigma_x^2$  COV, COR  $\rightarrow$  NOT

$\sigma_x$   $\rho_{xy}$   $\sigma_{xy}$  OBSERVED

STD. DEPT

INDIVIDUALS

RANDOM SAMPLING  $\rightarrow$  n OBJECTS RANDOMLY DRAWN FROM THE POPULATION

$y \rightarrow 10$  EXAMPLE

$\downarrow$  RANDOM SAMPLE OF n

ALL Y ARE DRAWN FROM THE SAME POP.

$\{y_1, y_2, y_3, \dots, y_n\} \rightarrow$

IDENTICALLY & INDEPENDENTLY

$$P(Y_1 \leq y) = P(Y_2 \leq y)$$

DISTRIBUTED

i.i.d  $\rightarrow$  KNOWING  $y_1$

$y_1$  &  $y_2$  ARE IDENTICALLY DISTRIBUTED.

PROVIDES NO INFO

KNOWING FIRST COIN /  $y_1 \rightarrow$  HEADS  
PROVIDES NO INFORMATION  $y_2 \rightarrow ?$

ON  $y_2$ .

$Y_1, Y_2, \dots, Y_N$  IID  $\rightarrow$  THROUGH RANDOM SAMPLING

$$\bar{Y} = \sum_i Y_i / n = (Y_1 + Y_2 + \dots + Y_N) / n$$

$\downarrow$  HAVE TO BE DONE  
THROUGH MACHINE  
ALGORITHM

USUALLY PROVIDES THE  
BEST INFO (INSTEAD OF

MEDIAN,  $Y_1, Y_{10}$   $n = 10,000$

$\hookrightarrow$  ASSUME

$\bar{Y} \rightarrow$  RANDOM VARIABLE WHY?

$\hookrightarrow$  EACH 10,000 RANDOM SAMPLE WILL  
GIVE YOU A DIFFERENT  $\bar{Y}$

$\bar{Y} \rightarrow$  HAS A RANDOM COMPONENT WITH  
AN EXPECTED VALUE AND PROBABILITY  
DISTRIBUTION.

$$Y \rightarrow E[Y] = \mu_Y \quad \text{JAR}(Y) = \sigma_Y^2$$

$$\bar{Y} \rightarrow E[\bar{Y}] = \mu_Y \quad \text{JAR}(\bar{Y}) = \sigma_{\bar{Y}}^2 / n$$

SEE THE TEXTBOOK  
FOR PROOF

$$\hookrightarrow \text{JAR}(\bar{Y}) = \sigma_{\bar{Y}}^2 = \sigma_Y^2 / n$$

$$Y \rightarrow E[Y] = \mu_Y$$

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

$$Y \rightarrow \text{JAR}(Y) \rightarrow \sigma_Y^2$$

$$\bar{Y} \rightarrow \text{JAR}(\bar{Y}) \rightarrow \sigma_{\bar{Y}}^2 \rightarrow \frac{\sigma_Y^2}{n}$$

$$Y \rightarrow \text{STD.DEV}(Y) \rightarrow \sigma_Y$$

$$\bar{Y} \rightarrow \text{STD.DEV}(\bar{Y}) \rightarrow \sigma_{\bar{Y}} \rightarrow \frac{\sigma_Y}{\sqrt{n}}$$

PROBLEM IS STILL THERE  $\rightarrow$  WE DON'T KNOW

$$\mu_Y, \sigma_Y^2$$

## LARGE SAMPLE APPROXIMATIONS

UNDER GENERAL CONDITIONS  $\bar{Y}$  IS AS GOOD AS

$$\bar{Y} \xrightarrow{P} \mu_Y \rightarrow \text{LAW OF LARGE NUMBERS}$$

1.  $y_1, y_2, \dots, y_n$  ARE iid

WITH  $E[y_i] = \mu_Y$  ALL DRAWS NO FREQUENT  
FROM SAME DIST.

2°  $\text{JAR}(y_i) = \sigma_Y^2 \perp \infty \rightarrow$  FINITE JARIANCE

**NORMAL DISTRIBUTION:**  $Y \sim N(\mu_Y, \sigma_Y^2)$

$Y \sim N(\mu_Y, \sigma_Y^2) \rightarrow$  NORMAL

$Z \sim N(0, 1) \rightarrow$  STANDARD NORMAL

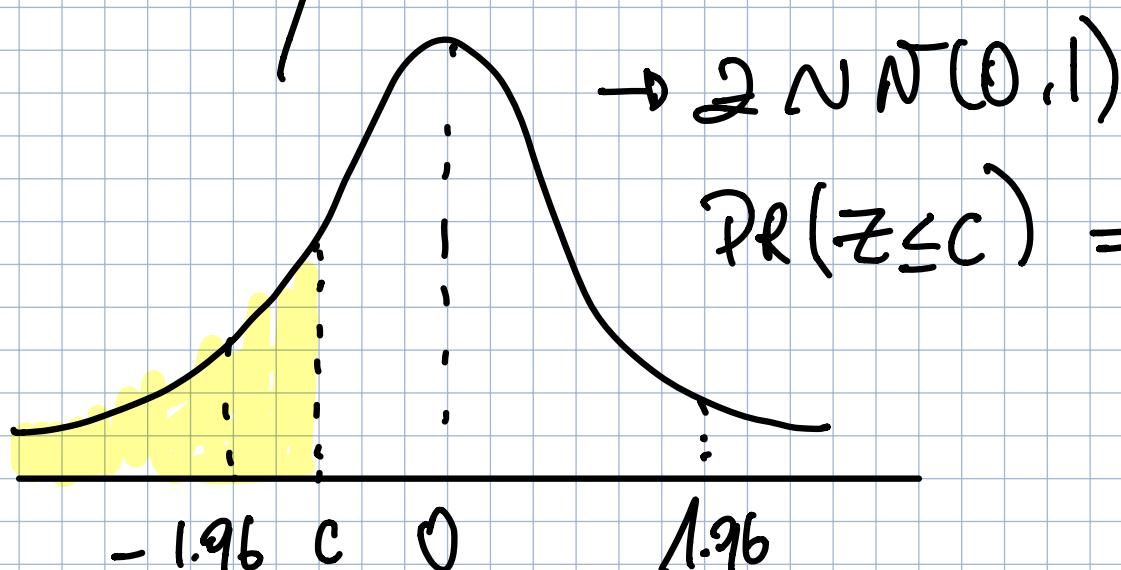
USUALLY BOTH ARE UNKNOWN  
(WE WILL PROXY THESE WITH SAMPLE METRICS)

- > ASYMETRIC
- > TAILS ASYMMETRIC
- >  $(-\infty, \infty)$

$$PR(\mu_Y - 1.96\sigma_Y < Y < \mu_Y + 1.96\sigma_Y) = 95\%$$

$$PR(Y < \mu_Y - 1.96\sigma_Y) = 2.5\%$$

$$PR(Y > \mu_Y + 1.96\sigma_Y) = 2.5\%$$



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

Φ PHIL,

IF  $Y \sim N(\mu_Y, \sigma_Y^2)$

THEN  $\frac{Y - \mu_Y}{\sigma_Y} = Z \sim N(0, 1)$

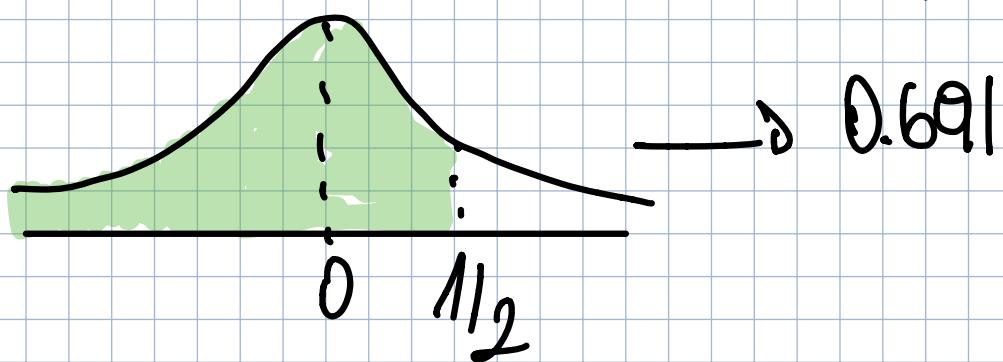
HANDBOOK QUESTION: IF  $Y \sim N(1, 4)$

$$\frac{Y - 1}{2} \sim N(0, 1)$$

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

$$P_R(Y \leq 2) \rightarrow P_R\left(\frac{Y-1}{2} < \frac{2-1}{2}\right)$$

$$\rightarrow P_R(Z \leq 1/2)$$



$$P_R(1 \leq Y \leq 2) \rightarrow P_R\left(\frac{1-1}{2} < \frac{Y-1}{2} < \frac{2-1}{2}\right)$$
$$P_R(0 \leq Z \leq 1/2)$$

02/04/2021

• **PECAP: LARGE SAMPLE APPROXIMATIONS TO SAMPLING APPROXIMATIONS**

- LAW OF LARGE NUMBERS
- CENTRAL LIMIT THEOREM

LOSSE  
 $Y \sim N(\lambda, \mu)$   
 $P(Y \leq Y \leq 2)$

POPULATION  $\rightarrow$  UNKNOWN BUT TRYING TO GATHER INFORMATION

↓ RANDOM SAMPLE  
n OBS

[ i.i.d ]

US IDENTICALLY & INDEPENDENTLY DISTRIBUTED

$\{Y_1, Y_2, \dots, Y_n\} \rightarrow Y_i$  - EACH OBSERVATION

$$\bar{Y} = \sum_i Y_i / n = \frac{1}{n} (Y_1 + Y_2 + \dots + Y_n)$$

LAW OF LARGE NUMBERS: AS n INCREASES  
iid SAMPLE  $\rightarrow$  FINITE VARIANCE

$$\bar{Y} \xrightarrow{P} E[Y]$$

$$n \rightarrow \infty$$

(N.B.)  $\rightarrow$  STARTS TO WORK

APPROXIMATION IN PROBABILITY

Will be very close to the population MEAN

CENTRAL LIMIT THEOREM: WITH  $n$  ENOUGH

$$\bar{Y} \sim N(\bar{\mu}_Y, \sigma_Y^2/n)$$

↳ THE INITIAL DISTRIBUTION OF  $Y$   
DOES NOT MATTER

$Y \rightarrow$  BINARY (COIN)  
UNIFORM (DIE)  
NORMAL (HEIGHT)

$$\left. \begin{array}{l} \\ \\ \end{array} \right\} \bar{Y} \sim N(\bar{\mu}_Y, \sigma_Y^2/n)$$

SIMULATION:  $Y \sim N(5, 625)$

$$E[Y] = 5 \quad \text{and} \quad \sigma_Y^2 = 625$$
$$\bar{Y} \sim N(5, 625/100)$$

$$\sigma_{\bar{Y}}^2 = \sigma^2/n = 625/100 = 6.25$$

$$\sigma_{\bar{Y}} = 2.5$$

$$n = 500 \quad \sigma_{\bar{Y}}^2 = 625/500 = 1.25$$

$$\sigma_{\bar{Y}} \approx 1.12$$

HANDOUT QUESTION:  $Y \sim N(1,4)$

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

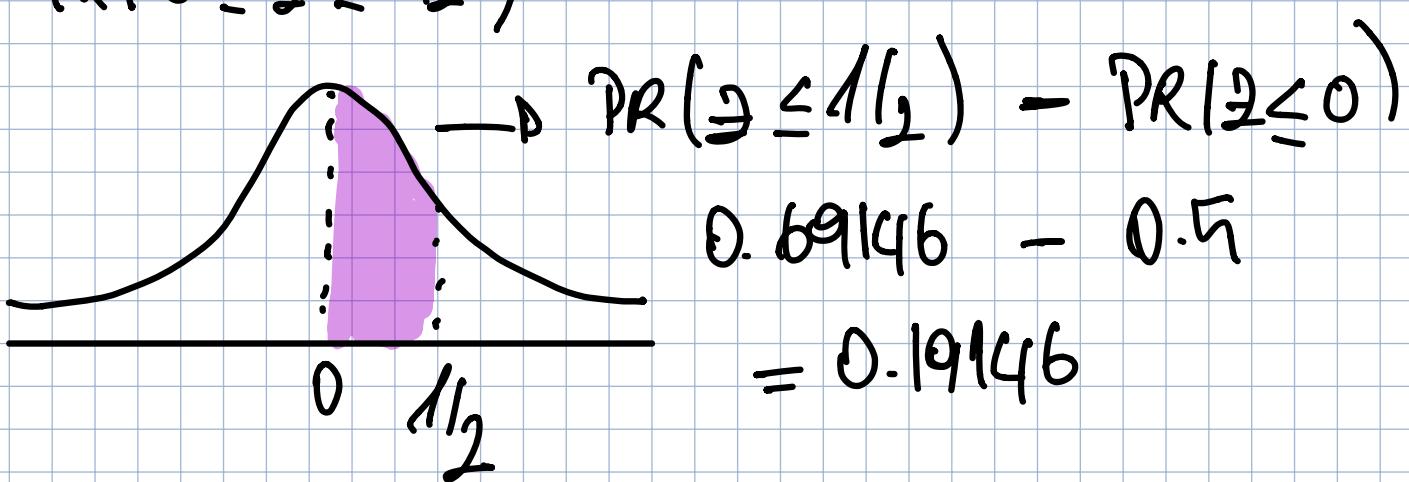
IN RANDOM VARIABLE ITSELF.

DRAW ONE OBSERVATION FROM THE  
DISTRIBUTION.

$$PR(1 \leq Y \leq 2) \leftrightarrow PR\left[\frac{1 - \mu_Y}{\sigma_Y} \leq \frac{Y - \mu_Y}{\sigma_Y} \leq \frac{2 - \mu_Y}{\sigma_Y}\right]$$

↑ IDENTICAL

$$PR(0 \leq Z \leq 1/2)$$



$$Y \sim ?(\mu_Y, \sigma_Y^2)$$

100, 43

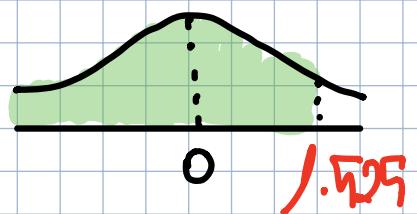
a.

CLT

$$\text{Is } \bar{Y} \sim N(100, 43/100)$$

$$PR(\bar{Y} < 101) \rightarrow PR\left(\frac{\bar{Y} - \mu_Y}{\sigma_{\bar{Y}}} < \frac{101 - 100}{\sqrt{43}}\right)$$

$$PR(\bar{Y} < 1.525) = 0.9364$$

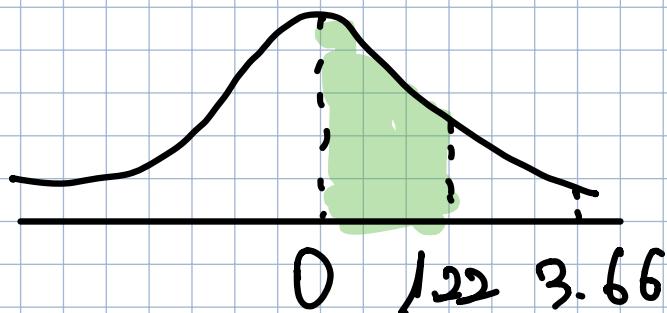


b.  $n=64 \quad PR(101 < \bar{Y} < 103)$

$$PR \left[ \frac{101 - 100}{\sqrt{0.6719}} < \frac{\bar{Y} - 100}{\sqrt{0.6719}} < \frac{103 - 100}{\sqrt{0.6719}} \right]$$

$\mu_Y = 100 \quad \Leftrightarrow PR(1.22 < Z < 3.66)$

$$\sigma_{\bar{Y}}^2/n = 0.6719$$



c. WHEN  $n=165$

$$\sigma_{\bar{Y}}^2 = \sigma_Y^2/n = 13/165 = 0.2606$$

$$PR(\bar{Y} > 98) = 1 - PR(\bar{Y} < 98)$$

$$1 - PR \left[ \frac{\bar{Y} - 100}{\sqrt{0.2606}} < \frac{98 - 100}{\sqrt{0.2606}} \right]$$

Almost 100%

# 02/08/2021 CH 3 STATISTICS REVIEW

↳ PROPERTIES OF A GOOD ESTIMATOR

↳ HYPOTHESES TESTING CONCERNING THE  
POPULATION MEAN

TRYING TO APPROXIMATE  $\mu_Y$ ,  $E[Y]$   
FOR RANDOM VARIABLE  $Y$

$\hat{\mu}_Y \rightarrow$  SAMPLE METRIC

CALCULATED USING  
AN i.i.d (RANDOM)  
SAMPLE

WHAT IS A GOOD ESTIMATOR?

1. UNBIASED  $E[\hat{\mu}_Y] = \mu_Y$

2. CONSISTENT  $\hat{\mu}_Y \xrightarrow{P} \mu_Y$  AS  $n \rightarrow \infty$

3. EFFICIENT (MINIMUM VARIANCE)

VAR( $\hat{\mu}_Y$ ) SHOULD BE AS SMALL AS  
POSSIBLE.

$\bar{Y} \rightarrow$  IS BLUE ESTIMATOR

UNBIASED

LINEAR

BEST



(MINIMUM VARIANCE - MOST PRECISE)

WHY?  $\rightarrow$  CHOOSE AN  $m$  THAT

$\sum_i (Y_i - m)^2$  MINIMIZE THIS

$$\begin{aligned} 1. \quad \sum_i Y_i / n &= \bar{Y} \quad E[\bar{Y}] = E\left[\sum_i Y_i / n\right] \\ &= 1/n \cdot E[\sum_i Y_i] \\ &= 1/n \cdot E[Y_1 + Y_2 + \dots + Y_n] \\ &= 1/n \cdot n \mu_Y \quad \text{So } E[\bar{Y}] = \mu_Y \\ &= \mu_Y \end{aligned}$$

$$\begin{aligned} 2. \quad \frac{d \sum_i (Y_i - m)^2}{dm} &= -2 \sum_i (Y_i - m) = 0 \\ &\quad -2n(\bar{Y} - m) = 0 \\ &\quad \bar{Y} = m \end{aligned}$$

TYPE I ERROR

TYPE II ERROR

SIGNIFICANCE LEVEL

CRITICAL VALUE

REJECTION ZONE

ACCEPTANCE ZONE

POWER OF THE TEST

P-VALUE

# TWO SIDED VS ONE-SIDED ALTERNATIVES

## CONFIDENCE INTERVALS

-HARD QUESTIONS

a.  $Y = \{0, 1\} \rightarrow$  FAIR COIN

$$\text{PR}(Y=0) = 0.5$$

$$\text{PR}(Y=1) = 0.5$$

b.  $E[Y] = \sum_i Y_i \text{ PR}(Y=Y_i)$

$$= 0 \times \text{PR}(Y=0) + 1 \times \text{PR}(Y=1)$$

$$E[Y] = \text{PR}(Y=1) = 0.5$$

$$\text{VAR}[Y] = \sum_i (Y_i - \mu_Y)^2 \text{ PR}(Y=Y_i)$$

$$= (0-0.5)^2 \times \text{PR}(Y=0) + (1-0.5)^2 \times \text{PR}(Y=1)$$

$$= 0.25 \times 0.5 + 0.25 \times 0.5$$

$$= 0.125 + 0.125 = 0.250$$

$$\sigma_Y^2 = 0.25 \quad E[Y] = \mu_Y$$

$$\sigma_Y = 0.5 \rightarrow \text{STD. DEF}$$

$$Y \sim B(0.5, 0.25)$$

c. CLT

$$\bar{Y} \sim N(\mu, \frac{\sigma^2}{n})$$

$$n = 100 \quad \bar{Y} \sim N(0.5, \frac{0.25}{100})$$

$$\sigma_{\bar{Y}}^2 = 0.0025 \quad E[\bar{Y}] = 0.5$$

$$\sigma_{\bar{Y}} = 0.05 \rightarrow \text{STD.DEV. OF SAMPLE MEAN}$$

d.  $n = 100$  61 HEADS  
39 TAILS

TEST IF  
THE COIN  
IS FAIR

NULL HYP.  
 $H_0$

$H_0: E[Y] = \mu_{y,0}$   $\xrightarrow{\text{VALUE IS DETERMINED BY THE RESEARCHER}}$   
DEFAULT ASSUMPTION THAT YOU ARE TRYING TO REJECT

$H_A: E[Y] \neq \mu_{y,0}$   $\xrightarrow{\text{TW-SIDED}}$   
HYPOTHESIS

ALTERNATIVE

$H_0: E[Y] > 0.5$   $H_A: E[Y] \neq 0.5$   
COIN IS FAIR COIN IS NOT FAIR

TWO POSSIBLE OUTCOMES: REJECT OR FAIL TO REJECT  
THE  $H_0$

DO NEVER 100% EXACTLY  
BECAUSE OF SAMPLE VARIATION

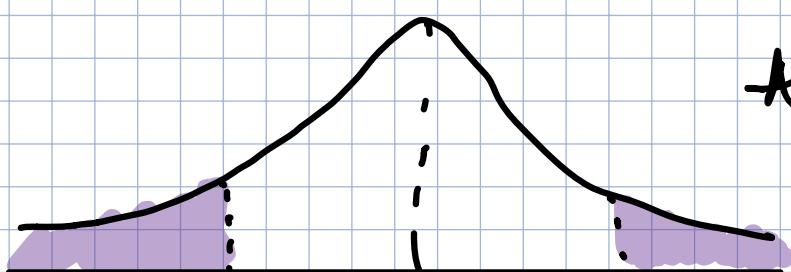
APRER  $n=100$  (PROVIDE A PROBABLE GUESS BASED ON SAMPLE)

HOW?

Q ASSUME Null IS CORRECT AND USE CLT

IF  $H_0: \mu_y = 0.5$  THEN  $\bar{Y} \sim N(0.5, 0.0025)$   
AND  $n=100$  → CALCULATE  $\bar{Y}$

AND P-VALUE



$$\mu_y - (\bar{Y} - \mu_y) / \sigma_{\bar{Y}} \quad \mu_y + (\bar{Y} - \mu_y) / \sigma_{\bar{Y}}$$

$$2 \times \Phi \left( - \left| \frac{\bar{Y} - \mu_y}{\sigma_{\bar{Y}}} \right| \right) = P\text{-VALUE}$$

$H_0: \mu_y = 0.5$

$$\bar{Y} = 0.61$$

$H_A: \mu_y \neq 0.5$

PROBABILITY TO DRAW A SAMPLE WITH A MEAN THAT IS

$|\mu_y - \bar{Y}|$  ABOVE OR BELOW  $\mu_y$  IF  $H_0$  IS TRUE

P-VAL : THE LIKELIHOOD OF DRAWING A SAMPLE  $n=100$  THAT HAS A SAMPLE MEAN THAT IS 0.11 ABOVE / BELOW

IF  $H_0$  IS CORRECT

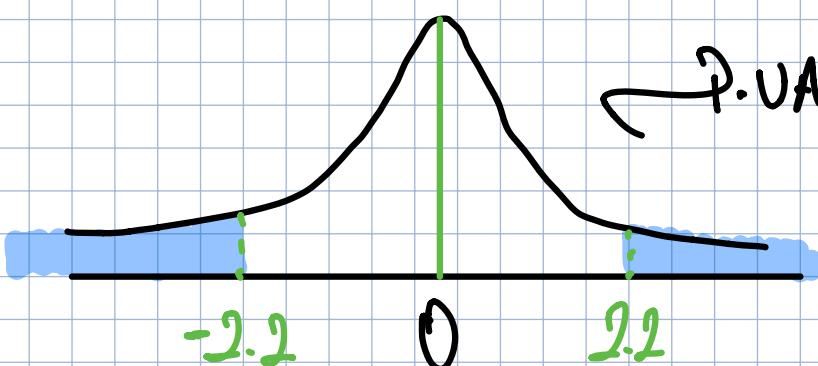
### P-VAL CALCULATION

- 1. •  $\sigma_y^2$  IS KNOWN
- 2. •  $\sigma_y^2$  IS UNKNOWN

$$1. \rightarrow Z\text{-SCORE} = \frac{\bar{Y} - \mu_{H_0}}{\sigma_{\bar{Y}}} \quad \begin{matrix} \text{TO HYP.} \\ \text{?} \end{matrix}$$

$$= \frac{0.61 - 0.50}{0.05} = \frac{0.11}{0.05} = 2.2$$

AS WE KNOW THIS  
B/C  $YNB(0.5, 0.25)$



$$P\text{-VAL} = 2 \times \Phi[-Z\text{-SCORE}]$$

IN STD-NORMAL.

DIST.

$$\begin{aligned} & 2 \times \Phi(-2.2) \\ & = 2 \times 0.0139 \\ & = 0.0278 \end{aligned}$$

IF THE COIN IS FAIR, THEN

THERE IS A 1.78% CHANCE

THAT WITH  $n=100$   $\bar{Y} > 0.61$  OR  $\bar{Y} < 0.39$

f. P-VAL: TYPE-I ERROR  $\rightarrow \alpha$  (ALPHA LEVEL)

↳

$H_0$  IS CORRECT

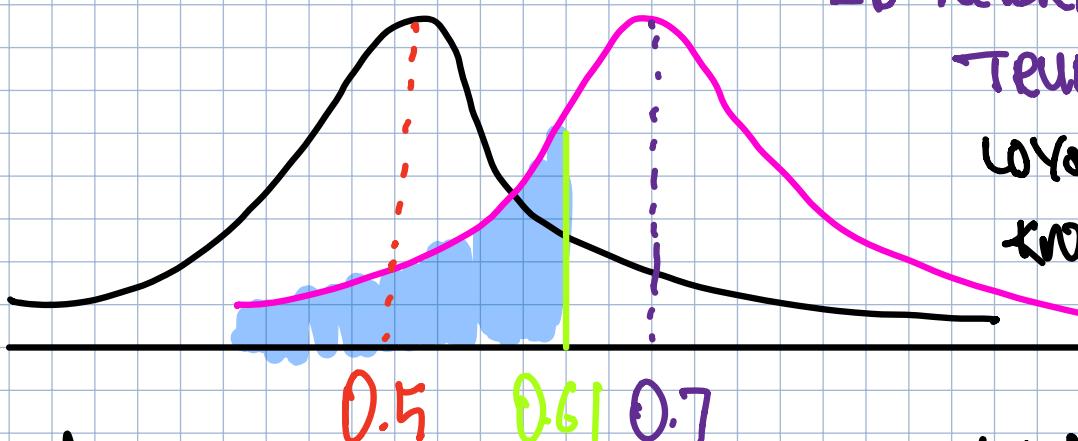
{ BUT YOU REJECT IT.

g. TYPE-II ERROR:  $H_0$  IS FALSE BUT YOU FAIL  
FALSE NEGATIVE TO REJECT IT

-> ALTERNATIVE  $H_A$

TRUE  $\mu_Y = 0.7$

(YOU NEED TO  
KNOW TO CALCULATE  
TYPE II.)



$$E[\bar{Y}_Y] = 0.70 \quad \text{VAR}(Y) = \sum_i (Y_i - \mu_Y)^2 \cdot \Pr(Y=Y_i)$$

$$\hat{\sigma}_Y^2 = (0-0.7)^2 \times 0.30 +$$

$$\bar{Y} \sim N(0.70, 0.0021) \quad (1-0.7)^2 \times 0.70$$

$$\Phi \left[ - \left| \frac{0.61 - 0.70}{\sqrt{0.0021}} \right| \right]$$

$$\hat{\sigma}_{\bar{Y}}^2 = 0.21$$

$$\hat{\sigma}_{\bar{Y}} = 0.21 / \sqrt{100}$$

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

$$= 0.0021$$

h. IF  $\sigma_Y^2$  IS UNKNOWN, USE  $s_Y^2$  AS AN ESTIMATOR FOR  $\sigma_Y^2$

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

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

$$\bar{Y}^2 = \frac{\bar{Y}_y^2}{n} \leftarrow \bar{Y}_{\bar{Y}} = \frac{\bar{Y}_y}{\sqrt{n}}$$

$$SE[\bar{Y}] = \bar{Y} = S_y / \sqrt{n}$$

STANDARD  
ERROR OF  $\bar{Y}$

$$t\text{-stat} = \frac{\bar{Y} - \mu_{Y,0}}{SE[\bar{Y}]}$$

$$S_y = 0.5$$

$$\bar{Y} = 0.049$$

$$t\text{-stat} = \frac{0.61 - 0.5}{0.049} \approx 2.044$$

$$P\text{-TWE} = 0.02706281$$

$\text{95\% CI} \rightarrow \bar{Y} \pm 1.96 SE[\bar{Y}]$

$$0.61 \pm 1.96 \times 0.049 [0.51 \quad 0.70]$$

95% PROBABILITY THAT THE TRUE PARAMETER  $\mu_Y$

WITHIN  
THIS  
INTERVAL

**POWER OF A TEST =** 1. PR (TYPE II ~~error~~)

↳ PROBABILITY THAT TEST CORRECTLY  
REJECTS THE  $H_0$

COMPARING MEANS FROM  
TWO DIFFERENT POPULATIONS:

$$\begin{array}{c}
 x \text{ } \& \text{ } y \rightarrow \text{two different distributions} \\
 \downarrow \qquad \downarrow \\
 \mu_y \qquad \mu_x \\
 \downarrow \qquad \downarrow \\
 \bar{x} \qquad \bar{y} \\
 \sigma^2_y \qquad \sigma^2_x \\
 \text{SE}[\bar{x}] \qquad \text{SE}[\bar{y}]
 \end{array}
 \quad
 \begin{aligned}
 \bar{x} &\sim N(\mu_x, \sigma^2_x / n_x) \\
 \bar{y} &\sim N(\mu_y, \sigma^2_y / n_y) \\
 \text{IF } x \text{ } \& \text{ } y \text{ ARE BOTH id} \\
 (\bar{x} - \bar{y}) &\sim N\left(\mu_x - \mu_y, \frac{\sigma^2_x}{n_x} + \frac{\sigma^2_y}{n_y}\right)
 \end{aligned}$$

UNKNOWN

$$\text{SE}[\bar{x} - \bar{y}] = \sqrt{\frac{\sigma^2_x}{n_x} + \frac{\sigma^2_y}{n_y}}$$

$$H_0: \mu_x - \mu_y = d_0$$

$$H_A: \mu_x - \mu_y \neq d_0$$

$$t = \frac{(\bar{x} - \bar{y}) - d_0}{SE[\bar{x} - \bar{y}]}$$

$$\text{P-VAL} = \sum_{x \in \Omega} (-H|)$$

CONFIDENCE INTERVAL:

$$(\bar{x} - \bar{y}) \pm 1.96 \times SE[\bar{x} - \bar{y}]$$

## QUESTION 2:

$$\mu_{Y_{1,0}} = 15.4 \rightarrow \text{HYPOTHESIS}$$

$$\bar{Y} = 14.6$$

$$\sigma_Y = 2.5$$

$$n = 35$$

$$\sigma_{\bar{Y}} = \frac{2.5}{\sqrt{35}}$$

$$Z = \frac{(\bar{Y} - \mu_{Y_{1,0}})}{\sigma_{\bar{Y}}}$$

02/18/2021 LESSON PLAN:

- COMPARING THE MEANS FROM 2 POP.
- SOLVE HANDOUT QUESTION
- ANSWER COMPETE'S QUIZ PPT
- INTRODUCTION TO REGRESSION ANALYSIS (CHAPTER 4)

### SAMPLE COMPARING MEANS FROM 2 DIFFERENT POPULATIONS

$$H_0: \mu_x - \mu_y = d_0 \quad X \text{ & } Y \quad \text{RANDOM VAR}$$
$$H_A: \mu_x - \mu_y \neq d_0 \quad \text{DIFF. POP.}$$

$$\bar{x}, \bar{y}$$

$$s_x^2, s_y^2$$

CLT:

$$SE[\bar{x}] \quad SE[\bar{y}]$$

$$\bar{x} \sim N(\mu_x, \sigma_x^2 / n_x) \quad \xrightarrow{\text{SAMPLE SIZE FOR } X}$$

$$\bar{y} \sim N(\mu_y, \sigma_y^2 / n_y) \quad \xrightarrow{\text{SAMPLE SIZE FOR } Y}$$

$$(\bar{x} - \bar{y}) \sim N(\mu_x - \mu_y, \frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y})$$

TYPICAL PROBLEM :  $\sigma_x^2, \sigma_y^2$   
NOT OBSERVED

ASSUMPTION: X & Y ARE BOTH iid

$$SE[\bar{x} - \bar{y}] = \sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}$$

$$(\bar{x} - \bar{y}) - d_0$$

$$t = \frac{(\bar{x} - \bar{y}) - d_0}{SE[\bar{x} - \bar{y}]}$$

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

$$\text{CONFIDENCE INTERVALS } (\bar{x} - \bar{y}) \pm 1.96 \times SE[\bar{x} - \bar{y}]$$

$$s_x^2 \rightarrow \text{SAMPLE VARIANCE} = \frac{\sum_i (x_i - \bar{x})^2}{n-1}$$

$$s_y^2 \rightarrow \text{SAMPLE VARIANCE} = \frac{\sum_i (y_i - \bar{y})^2}{n-1}$$

$s_{xy}$  → SAMPLE COV

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

$\downarrow$   
 $s_{xy}$  → POPULATION COVARIANCE

→ SAMPLE COR.

$$r_{xy} = -1 \leq r_{xy} \leq 1$$

# CHAPTER 4 : REGRESSION ANALYSIS

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

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

W SAMPLE SIZE

RESPONSE VARIABLE

LEFT-HANDED VARIABLE

$Y$  = INDEPENDENT VARIABLE

OUTCOME VARIABLE

$X$  → INDEPENDENT VARIABLE

RIGHT-HANDED VARIABLE

PREDICTOR

$\beta_0, \beta_1$  → POPULATION COEFFICIENTS

$\beta_0$  : INTERCEPT       $\beta_1$  : PARAMETERS

$\beta_1$  : SLOPE

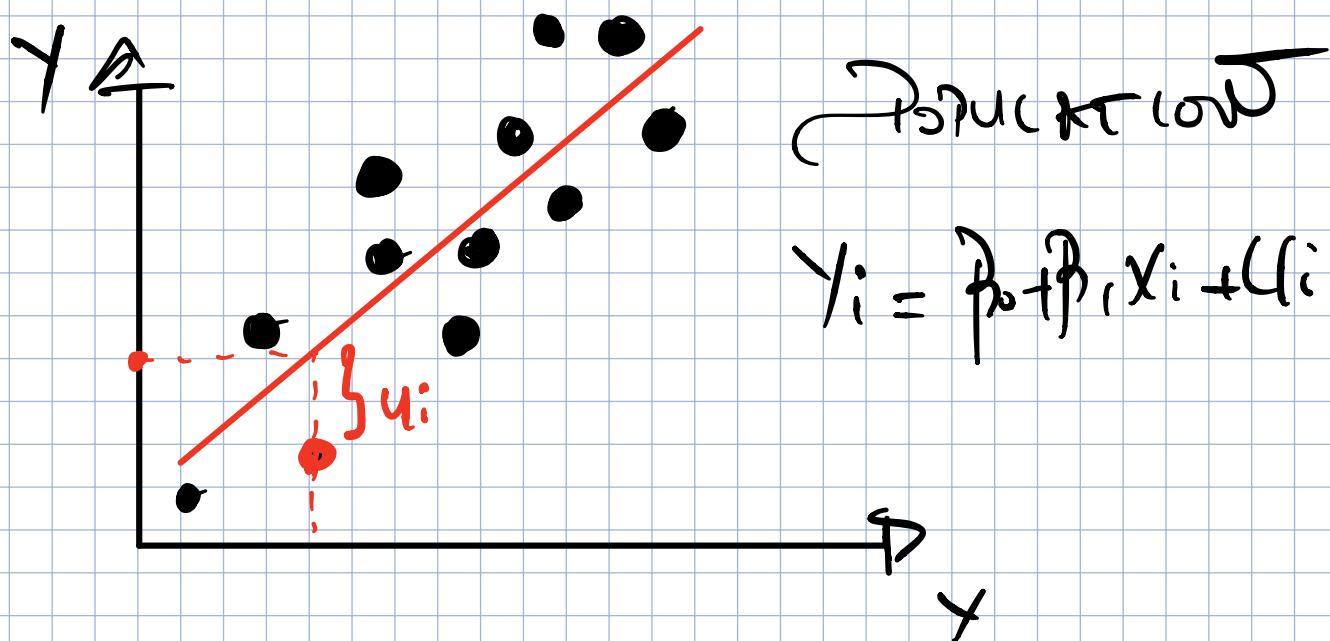
POPULATION PREDICTION

$\beta_0 + \beta_1 X_i$  → LINE

POPULATION PREDICTION  
FUNCTION

$U_i \rightarrow$  ERROR TERM

$\rightarrow \beta_0 + \beta_1 X_i$



But population parameters  $\beta_0, \beta_1$  can't be known because we don't observe the population data.

Two main problems:

1. Population parameters  $\rightarrow$  sample  
2. How to fit the best line?

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

↓

↳ Birth weight

EARNINGS

CHILD'S

SCHOOL TERM.

↳ SALES,  $\tau$

↳ PARENTS INCOME

↳ SPENDING ON SOCIAL MEDIA AD

using a sample of  $n$  observations,

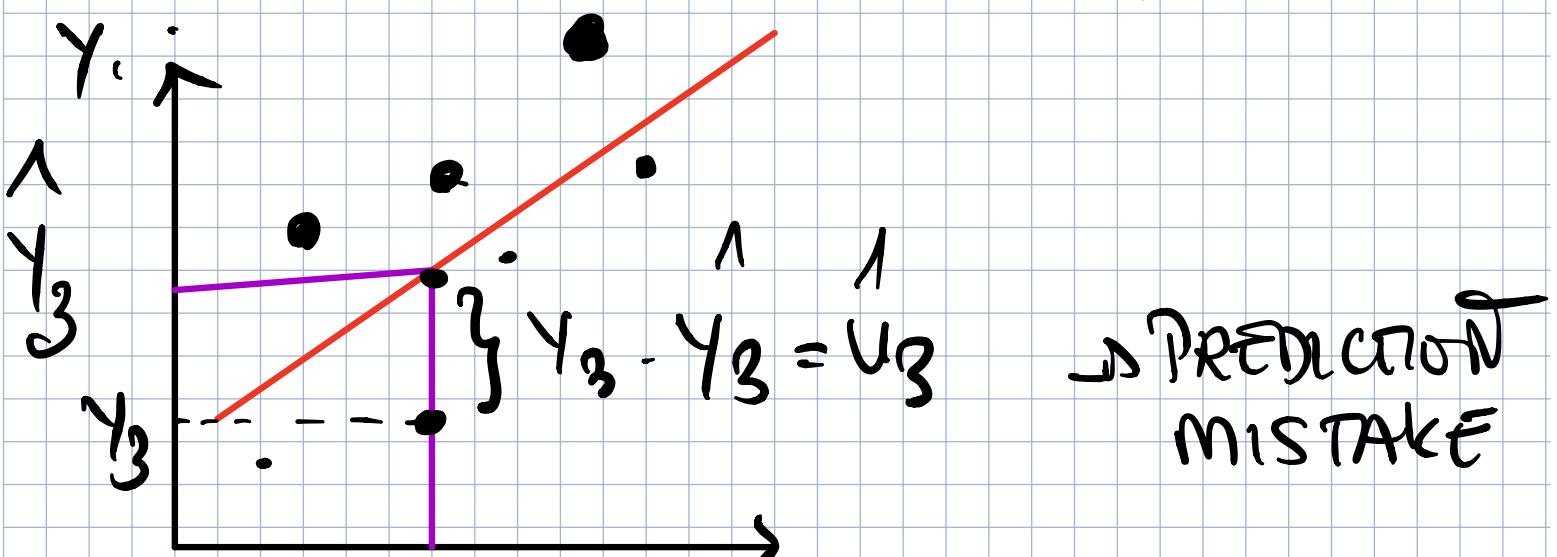
choose  $b_0 \rightarrow \beta_0$  [2 estimators]

$b_1 \rightarrow \beta_1$

that best describe the relationship  
between  $X_i$  &  $Y_i$ ?

$i$	$x_i$	$y_i$	$\hat{y}_i$	$u_i$
1	$x_1$	$y_1$	$\hat{y}_1$	$u_1$
2	$x_2$	$y_2$	$\hat{y}_2$	$u_2$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$x_n$	$y_n$	$\hat{y}_n$	$u_n$

THE LINE THAT  
MAKES THE LEAST  
POSSIBLE AMOUNT  
OF PREDICTION  
MISTAKE.



$$\sum_{i=1}^n u_i^2 = \sum_i (y_i - \beta_0 - \beta_1 x_i)^2$$

CHOOSE  $\beta_0, \beta_1$  TO MINIMIZE THE SUM OF SQUARED RESIDUALS

## ORDINARY LEAST SQUARES

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

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

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

$$U_i = Y_i - \hat{Y}_i \quad i = 1, \dots, n$$

## STANDARD ERROR OF THE REGRESSION

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

$$\text{SEE} = \underline{s_u} = \sqrt{\frac{\sum U_i^2}{(n-2)}} = \frac{\text{SSR}}{(n-2)}$$

MAGNITUDE OF A  
TYPICAL PREDICTION  
ERROR

*as  $\hat{\beta}_0, \hat{\beta}_1$  should  
BE KNOWN BEFORE  
SO LOSE 2 D.F*

MEASURES OF FIT: HOW GOOD THE PREDICTION LINE FITS THE DATA?

R<sup>2</sup> → R-SQUARED MOST COMMON

TSS → TOTAL SUM OF SQUARES

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

$\uparrow$   
in Y

SSR (RSS) = SUM OF SQUARED RESIDUALS  
RESIDUAL SUM OF SQUARES

$$\sum \hat{u}_i^2$$

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

UNEXPLAINED SUM OF SQUARES

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

$$R^2 = \frac{ESS}{TSS} \rightarrow \text{THE VARIATION IN } Y$$

EXPLAINED BY THE VARIATION IN X

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

SOME TRUE FACTS about  $\hat{y}_i$

$$1. \frac{1}{n} \sum_i \hat{y}_i = 0$$

$$2^{\circ} \frac{1}{n} \sum_i \hat{y}_i = \bar{Y} \quad \begin{array}{l} \text{REGRESSION LINE} \\ \text{CROSSES THROUGH} \end{array}$$

$$3. \sum_i \hat{y}_i x_i = 0 \quad \bar{x}, \bar{y}$$

$$4^{\circ} \sum_i \hat{y}_i^2 \text{ MINIMUM}$$

$$5^{\circ} TSS = ESS + SSR$$

STANDARD ERROR OF THE REGRESSION

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

$$SEE = \underline{\sum_i} = \sqrt{\frac{\sum_i \hat{y}_i^2}{(n-2)}} = \sqrt{\frac{SSR}{(n-2)}}$$

MAGNITUDE OF A  
PARTIAL PREDICTION  
ERROR

( $\hat{y}_i$  &  $\hat{y}_j$  should  
BE KNOWN BEFORE  
SO LOSE 2 D.F)

# LECTURE PLAN

03/08/2022

- i. OLS ASSUMPTIONS (SPECIFIC FOCUS ON THE CAUSAL INTERPRETATIONS OF OLS COEFFICIENTS  $\hat{\beta}_0, \hat{\beta}_1$ )
- ii. SAMPLING DISTRIBUTION OF THE OLS ESTIMATOR.
- iii. Go over a few points on the homework assignments

→ READING ON NEXT CLASS

i.5: CONFIDENCE INTERVALS & HYPOTHESES TEST USING A SINGLE REGRESSION

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

) RANDOM SAMPLE  
of n OBSERVATIONS

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

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{U}_i \rightarrow \text{SAMPLING VARIATION } (Y_i, \hat{Y}_i, \hat{U}_i)$$

THE LEAST SQUARES ASSUMPTIONS

1. OLS COEFFICIENTS  $\hat{\beta}_0$  &  $\hat{\beta}_1$  HAVE NO CAUSAL INTERPRETATION UNLESS  $E[U_i | X_i] = 0$

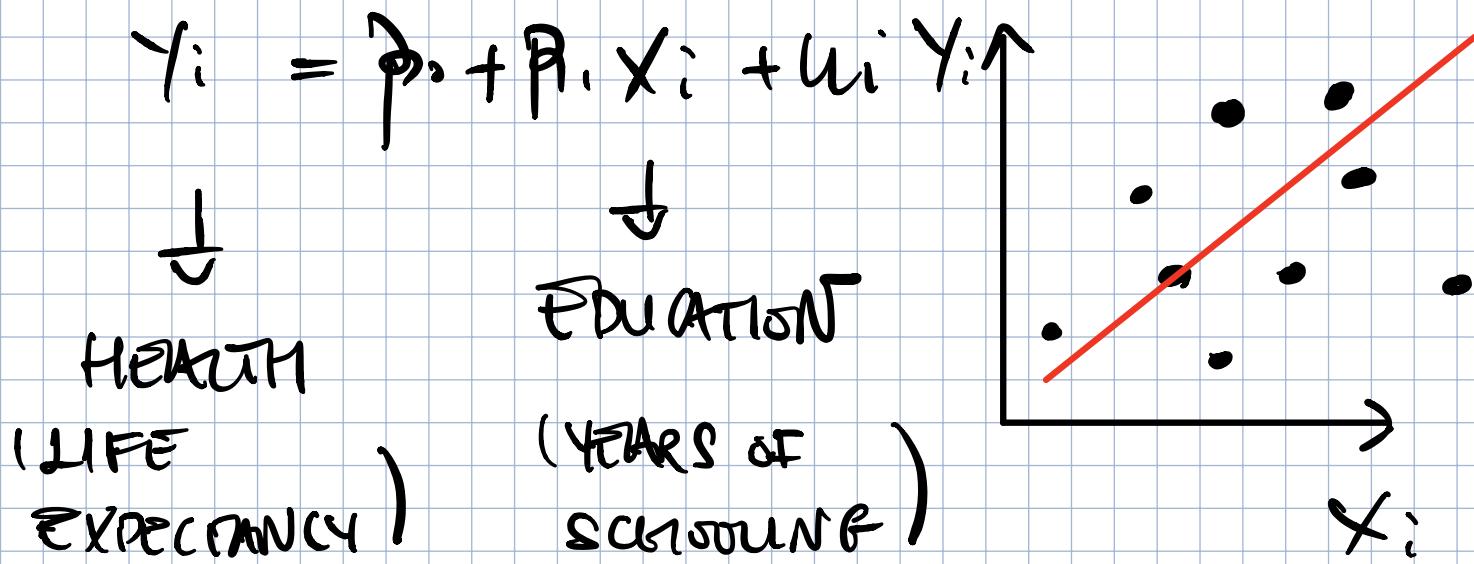
$$\mathbb{E}[u_i, x_i] = \text{COR}(u_i, x_i) = 0$$

$\hookrightarrow$  ASSUMPTION IS ABOUT THE

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

$\hookrightarrow$  POPULATION IS NOT PERFECT WAY TO TEST

$X_i$  &  $Y_i$  ARE ASSOCIATED BUT IT DOES  
 NOT MEAN THAT  $X_i \rightarrow Y_i$  )  
 (causes)



2 ISSUES:

$\beta_1 = 0.15$  AN ADDITIONAL OF EDUCATION IS ASSOCIATED WITH A 0.15 YEARS OF LIFE EXPECTANCY.

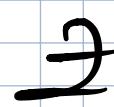
(WHEN CONSIDER CAUSALITY)

1. PERVERSE CAUSALITY



2. THIRD OMITTED

VARIABLE THAT IMPACT BOTH



HEALTH  $\rightarrow$  EDUCATION

$\rightarrow$  EDUCATION

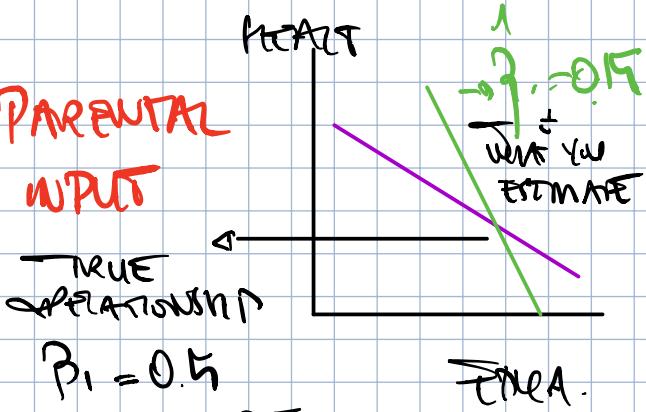
HEALTH

PARENTAL INFLUENCE

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

↓                  ↓  
HEALTH      EDUCATION

↑ PARENTAL INPUT



IF  $\text{COR}(x_i, U_i) = 0 \rightarrow \beta_1 \text{ CAN ONLY}$

N. OMITTED VARIABLE

BE DUE TO

BIAIS  $\rightarrow$  IF  $\text{COR}(x_i, U_i) \neq 0 \quad x_i \in E[U_i | x_i] = 0$   
WHEN NO LINEARITY

2° CAUSAL INTERPRETATION

OLS ASSUMPTION 2°  $x_i, y_i \quad i = 1, \dots, n$

ARE  $i, i, \dots$  DISTRIBUTED

↓ IDENTICALLY

INDEPENDENTLY

PANDEMONIUM  
SOMETHING MISWEEPS

THIS ASSUMPTION

3° LARGE OUTLIERS ARE UNLIKELY

(IF NOT OLS MIGHT BE MISLEADING)