

January 21, 2023

DSML: CC Fundamentals.

# Experiment Design

Recap:

- (a) Hypothesis testing -  $\alpha$ , p-value, Null, Alternate
- (b) Z-test { 1 sample, 2 sample, independent, paired,
- (c) t-test { left-tailed, right-tailed, 2-tailed.
- (d) R-S test
- (e) ANOVA
- (f) Kendall
- (g) Pearson r
- (h)  $\chi^2$  test



Class begins @ 9:05 p.m.

## Agenda:

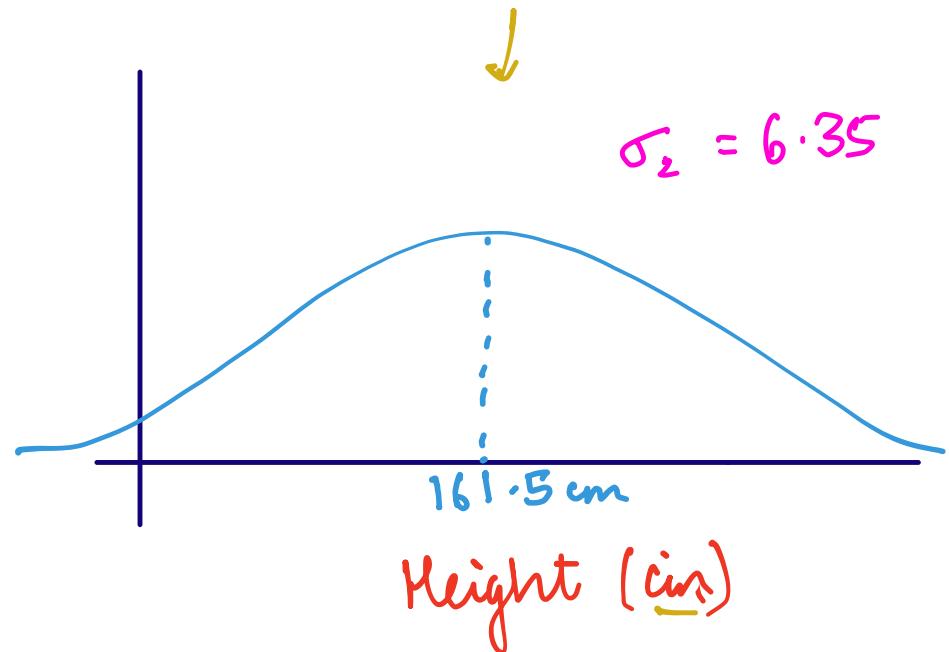
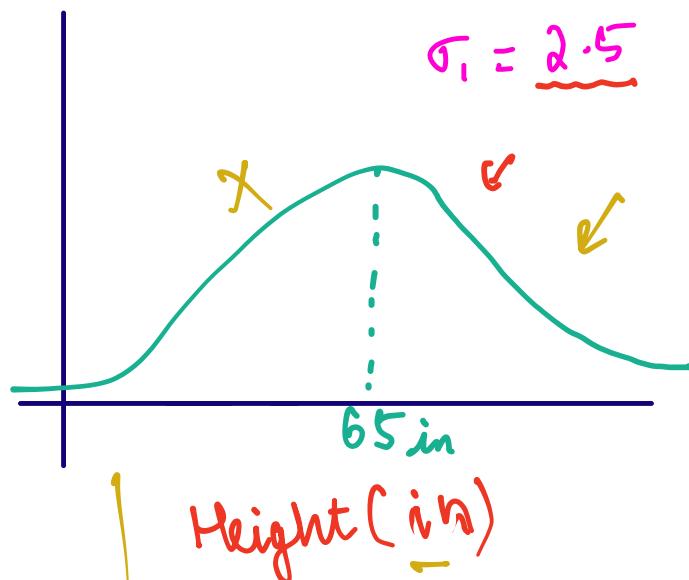
### Part 1: Normalization and Standardization.

- \* sklearn.preprocessing - StandardScaler, MinMaxScaler.

### Part 2: Experiment Design.

- \* Defining  $\alpha$  - significance level.
- \* Defining  $\beta$  - power of the test.
- \* Obtaining sample size.

## Standardization : (Column Standardization)



$$1 \text{ inch} = 2.54 \text{ cm.}$$

$$\rightarrow X = \left\{ \bar{x}_i \in \mathbb{R} \right\}_{i=1}^n \rightarrow Y = \left\{ \bar{y}_i \in \mathbb{R} \right\}_{i=1}^n$$

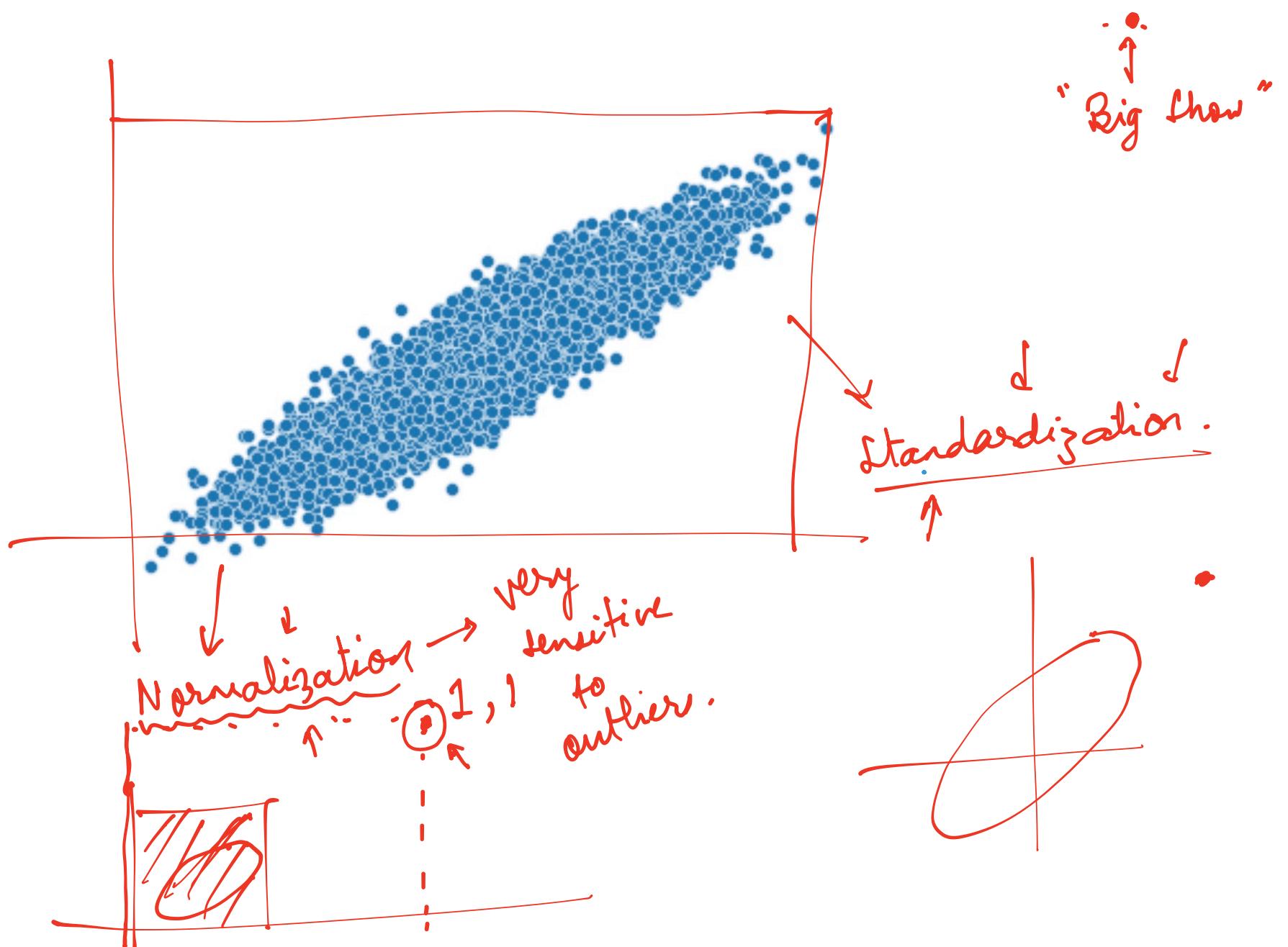
$$\bar{y}_i = \frac{\bar{x}_i - \mu}{\sigma} \rightarrow \text{standardization.}$$

## Normalization

→  $X = \{\bar{x}_i \in \mathbb{R}\}_{i=1}^n \rightarrow Y = \{\bar{y}_i \in \mathbb{R}\}_{i=1}^n$

$$y_i = \frac{x_i - \min(x_i)}{\max(x_i) - \min(x_i)} = 1$$



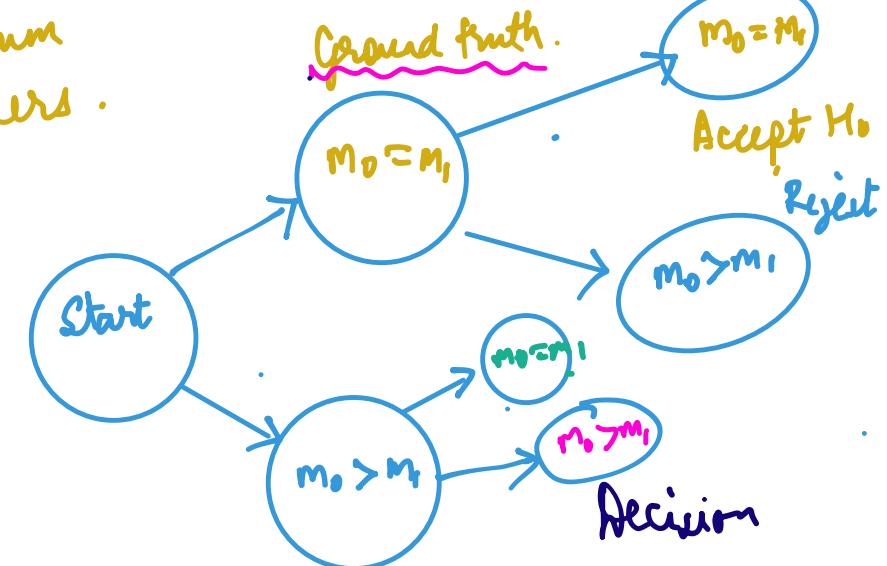
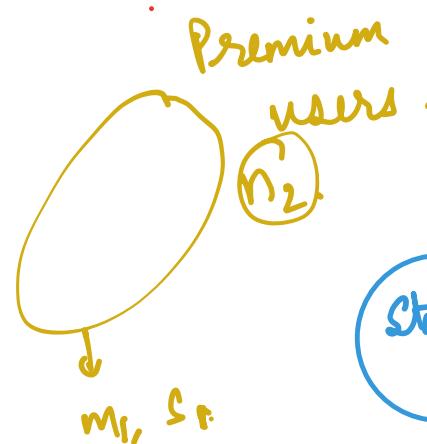
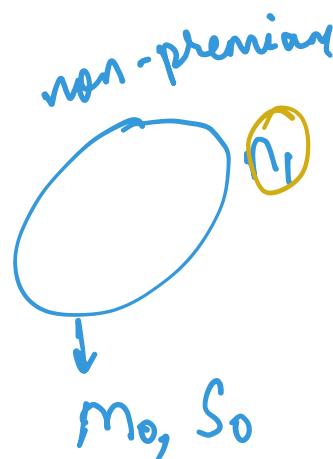


Experiment Design: YouTube data.

Watch time  $\rightarrow$  min.

Premium  $\rightarrow$  \* or 0

Decision



$$H_0 : m_0 = m_1$$

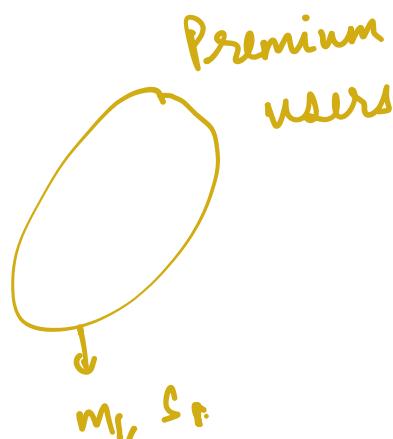
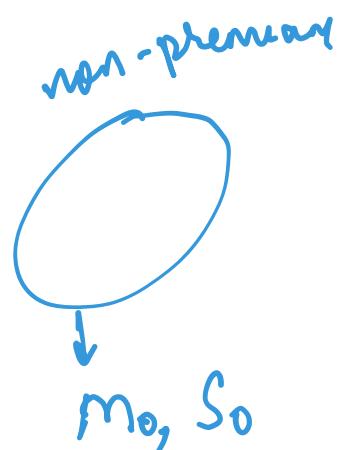
$$H_a : m_0 > m_1$$

$\alpha \rightarrow$  Significance level.  $\rightarrow$  Type I error.

$\beta \rightarrow$  Measure of false negatives  $\rightarrow$  Type 2 error.

Reality			Accept	Reject
	True	False	True-ve. $1 - \alpha$	False-ve. $\alpha$
True				
False			False-ve. $\beta$	True-ve. $1 - \beta$

"Power of the test"



Reality	Accept		Reject
	True	False	
True	True - ve. $1 - \alpha$	False - ve. $\beta$	False + ve. $\alpha$
False	False - ve. $\beta$	True + ve. $1 - \beta$	True + ve. $1 - \beta$

$$H_0 : m_0 = m_1$$

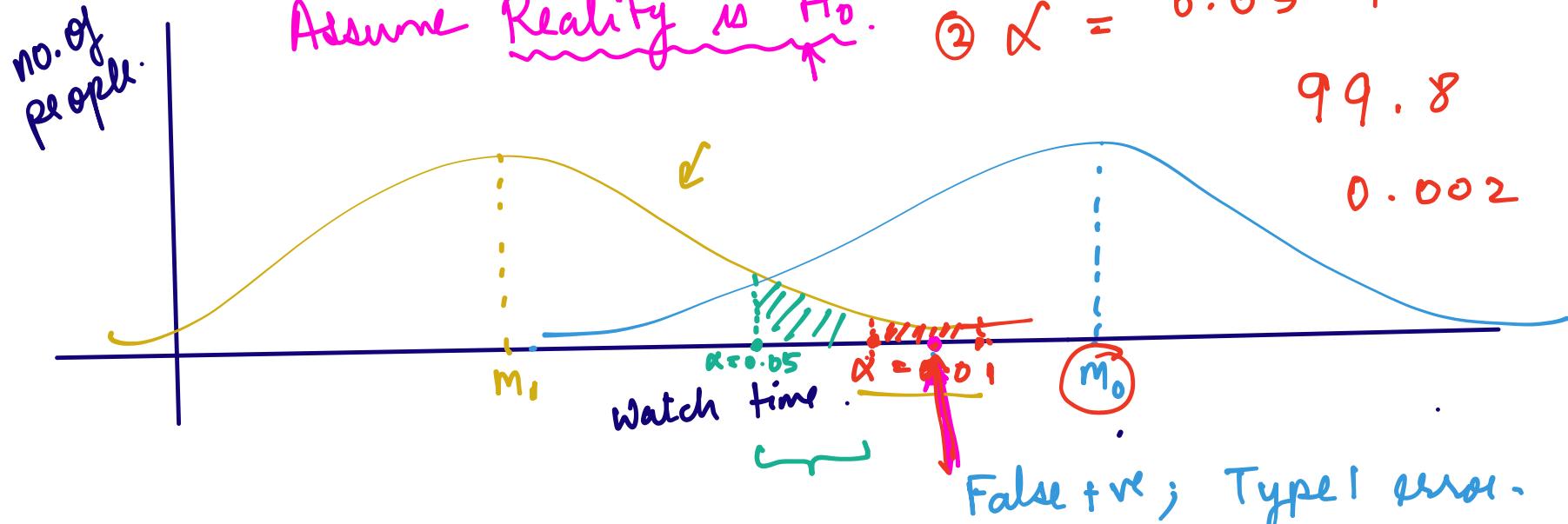
$$H_a : m_0 > m_1$$

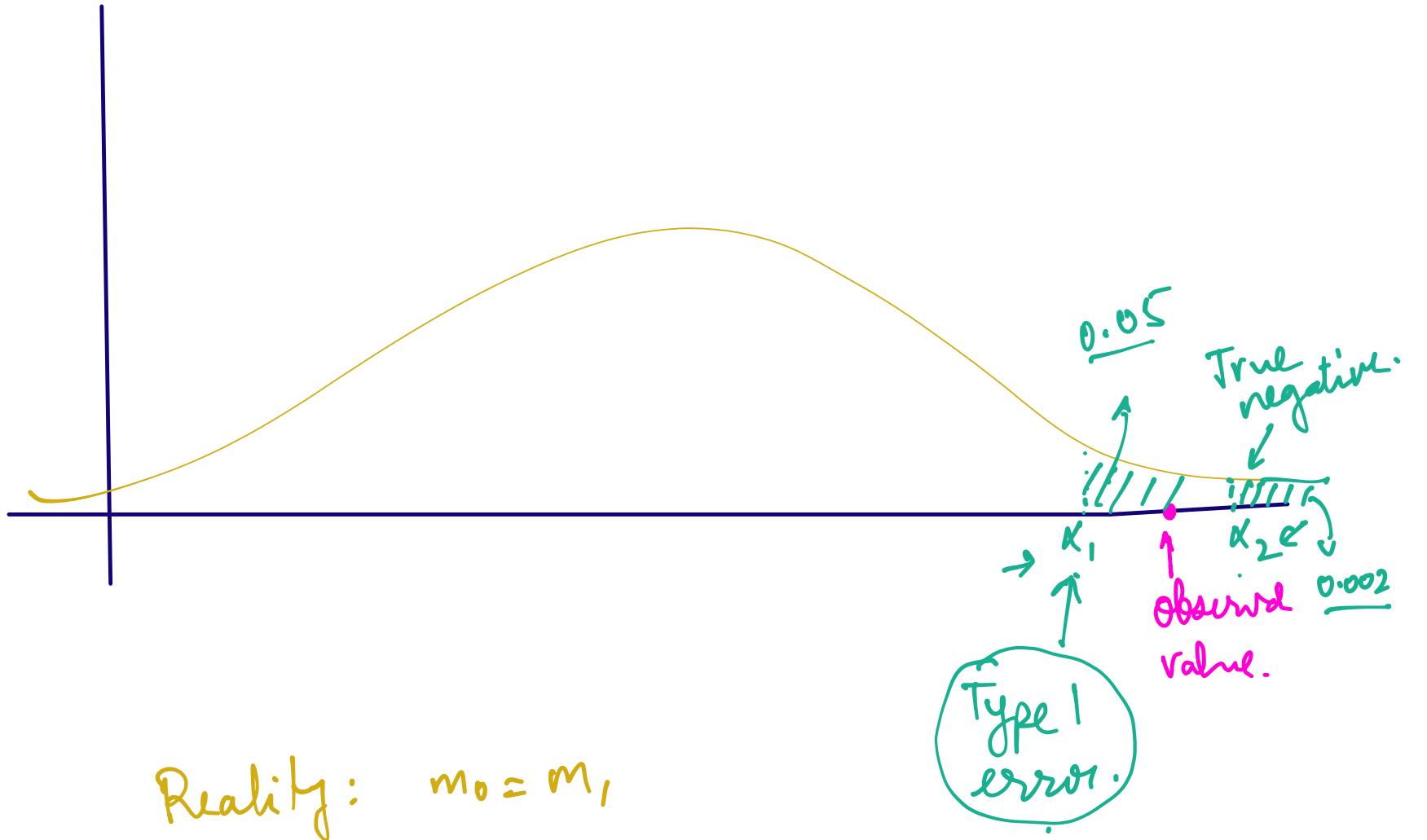
Right-tailed test

$$\textcircled{1} \quad \alpha = 0.01$$

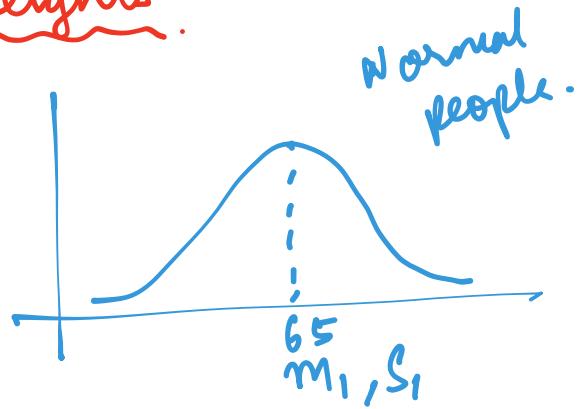
more strict.

Assume Reality is  $H_0$ .  $\textcircled{2} \quad \alpha = 0.05$





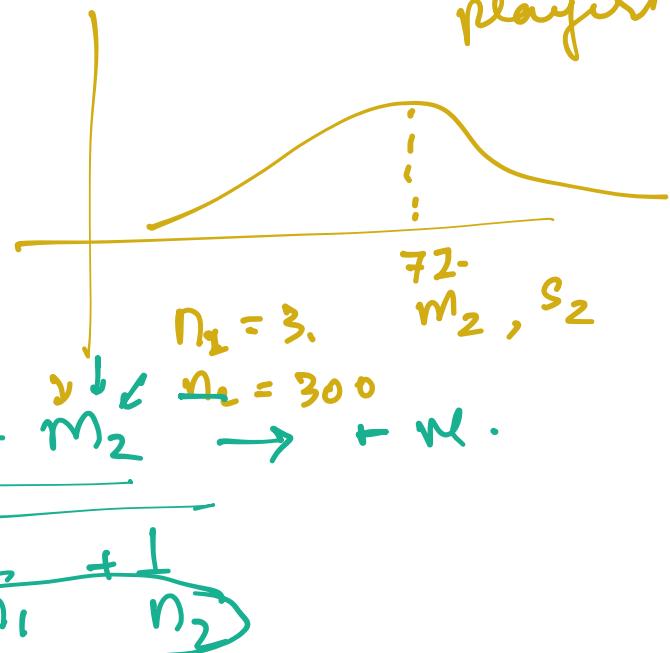
Heights:



$$n_1 = 2$$

$$n_1 = 200 \quad t\text{-statistic} : \frac{m_1 - m_2}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \rightarrow t \sim N.$$

USA Basketball players.

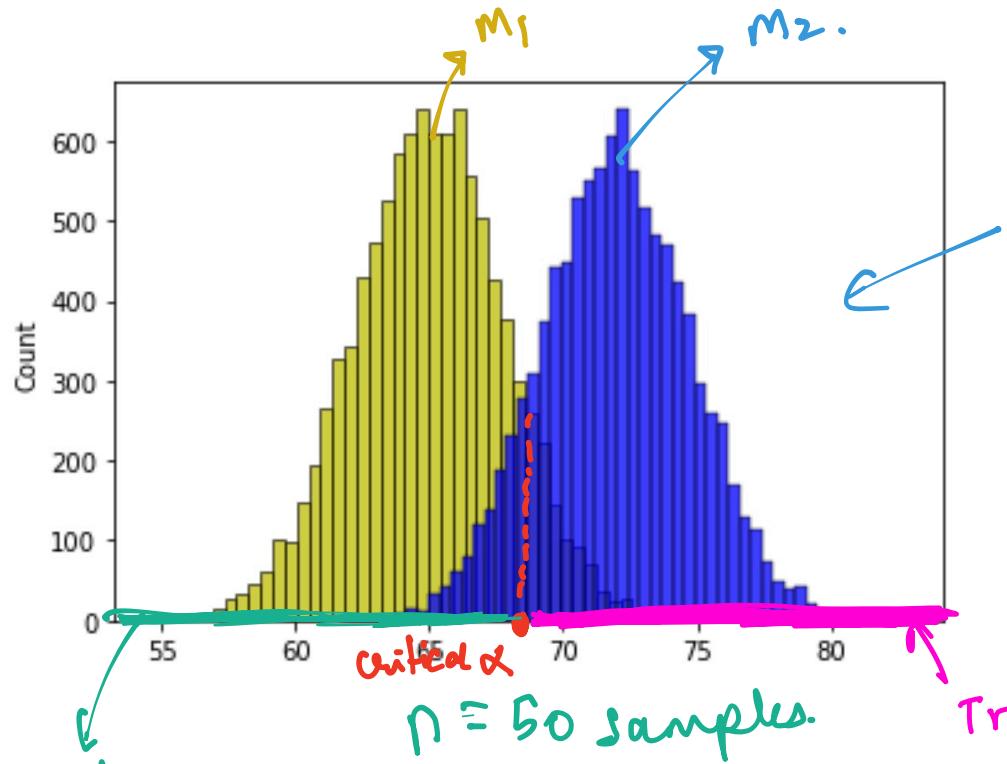


$S$  = Pooled sample standard deviation.

$$\sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}.$$

$$t_1 \rightarrow \text{where } n_1 = 2, n_2 = 3 = 2.0.$$

$$t_2 \rightarrow -$$

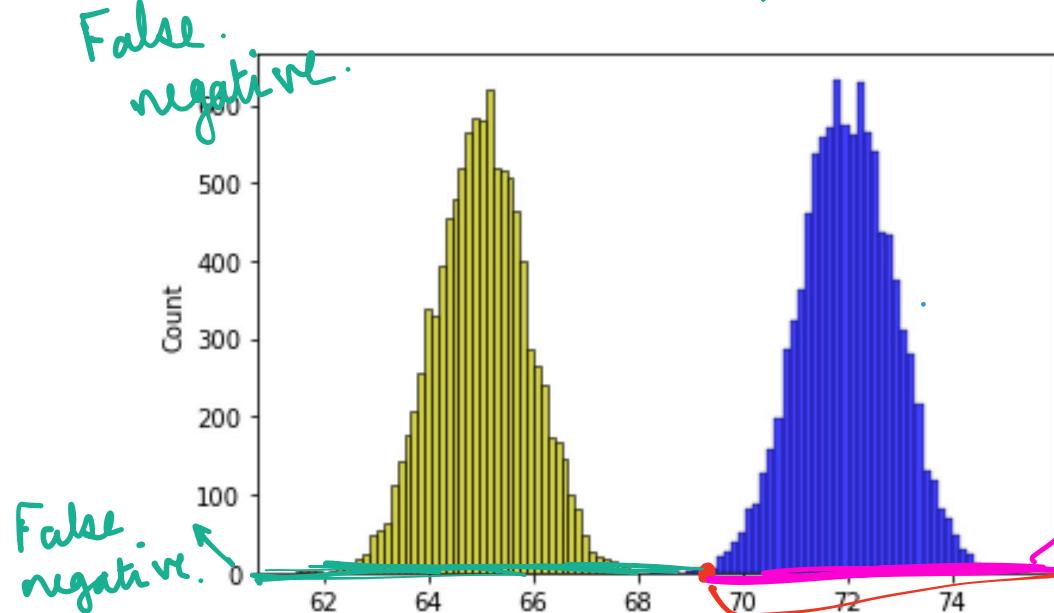


Reality :  $m_1 < m_2$ .  
 $n = 50$  samples.

• → observed  
 $m_2$ .

$$\sigma / \sqrt{n}$$

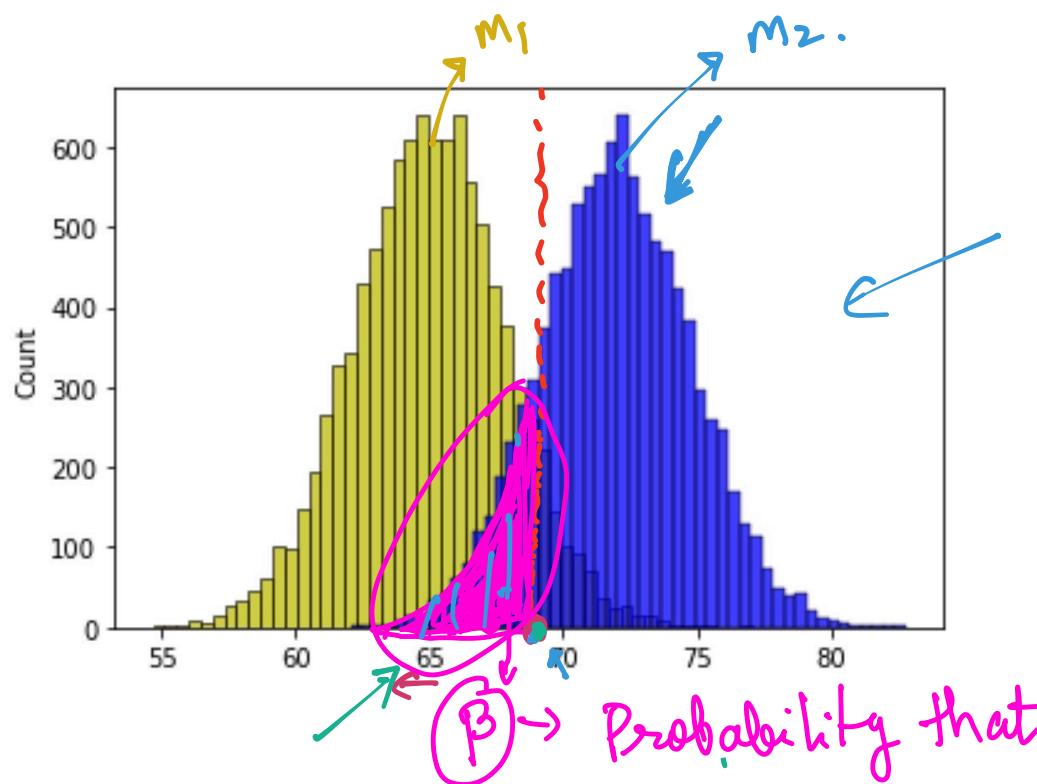
True positive.



$n \uparrow$  i.e.  $n = 500$   
Reduces False negatives.

True positive.

$\alpha$ .



Reality :  $M_1 < M_2$

$n = 50$  samples.

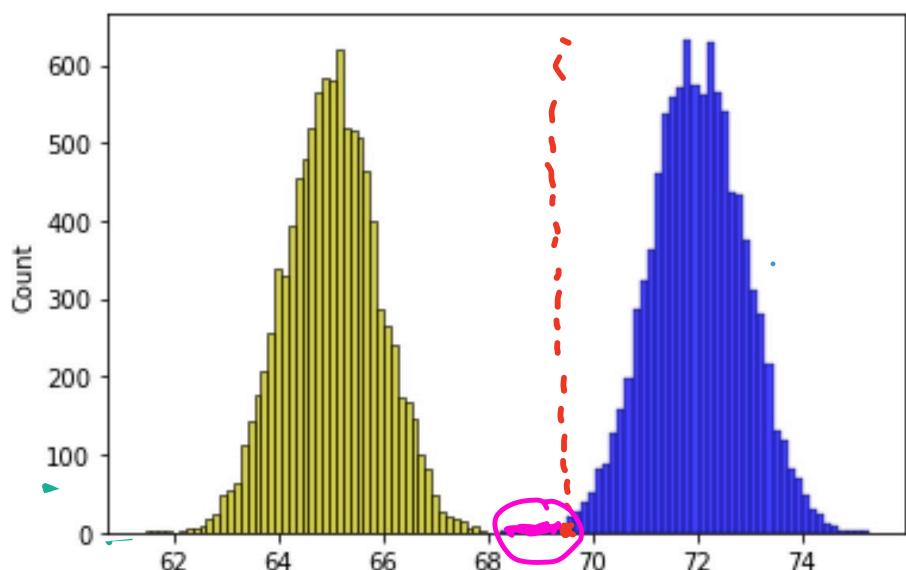
○ → observed  $M_2$ .

$$\sigma / \sqrt{n}$$

$M_1 < M_2$  but we fail to reject.

$n \uparrow$  i.e.  $n = 500$   
Reduces False negatives.

False - re.  
 $\downarrow$



$n \rightarrow$  Sample size is the  
most critical parameter.

Choosing  $n$  as large as possible  
is the best case scenario!

\* Core question: What  $n$  is just  
enough??

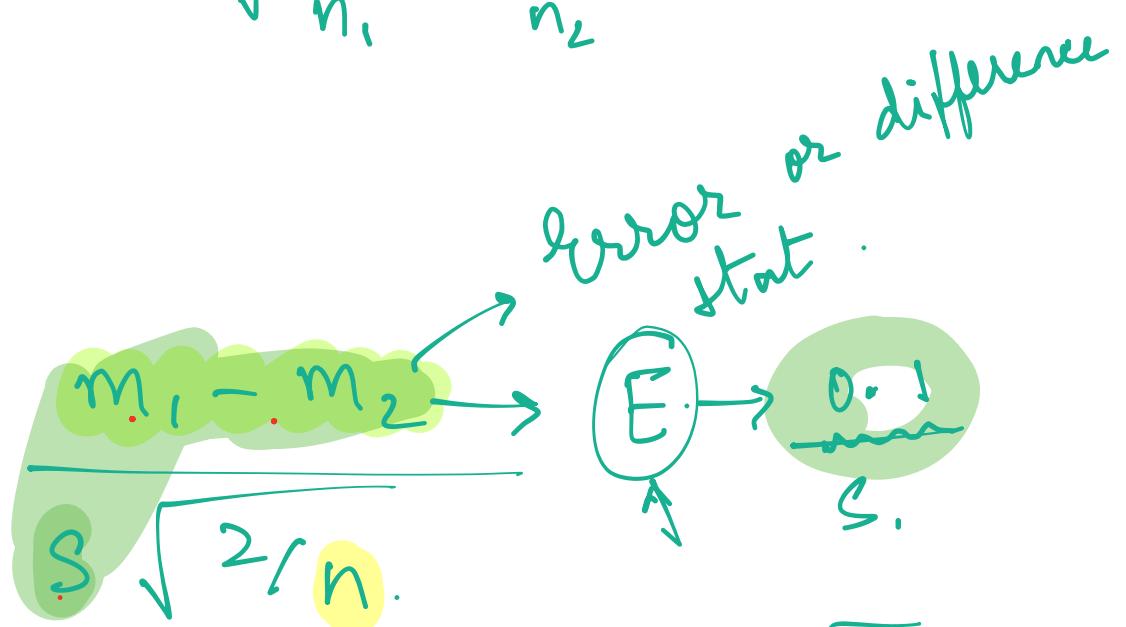
$$t\text{-stat} = \frac{m_1 - m_2}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$n_1 = n_2$$

$t_{critical} =$   
decided by  $\alpha$  &  $1-\beta$ .

$$t_c = \frac{E}{S \cdot \sqrt{2/n}}$$

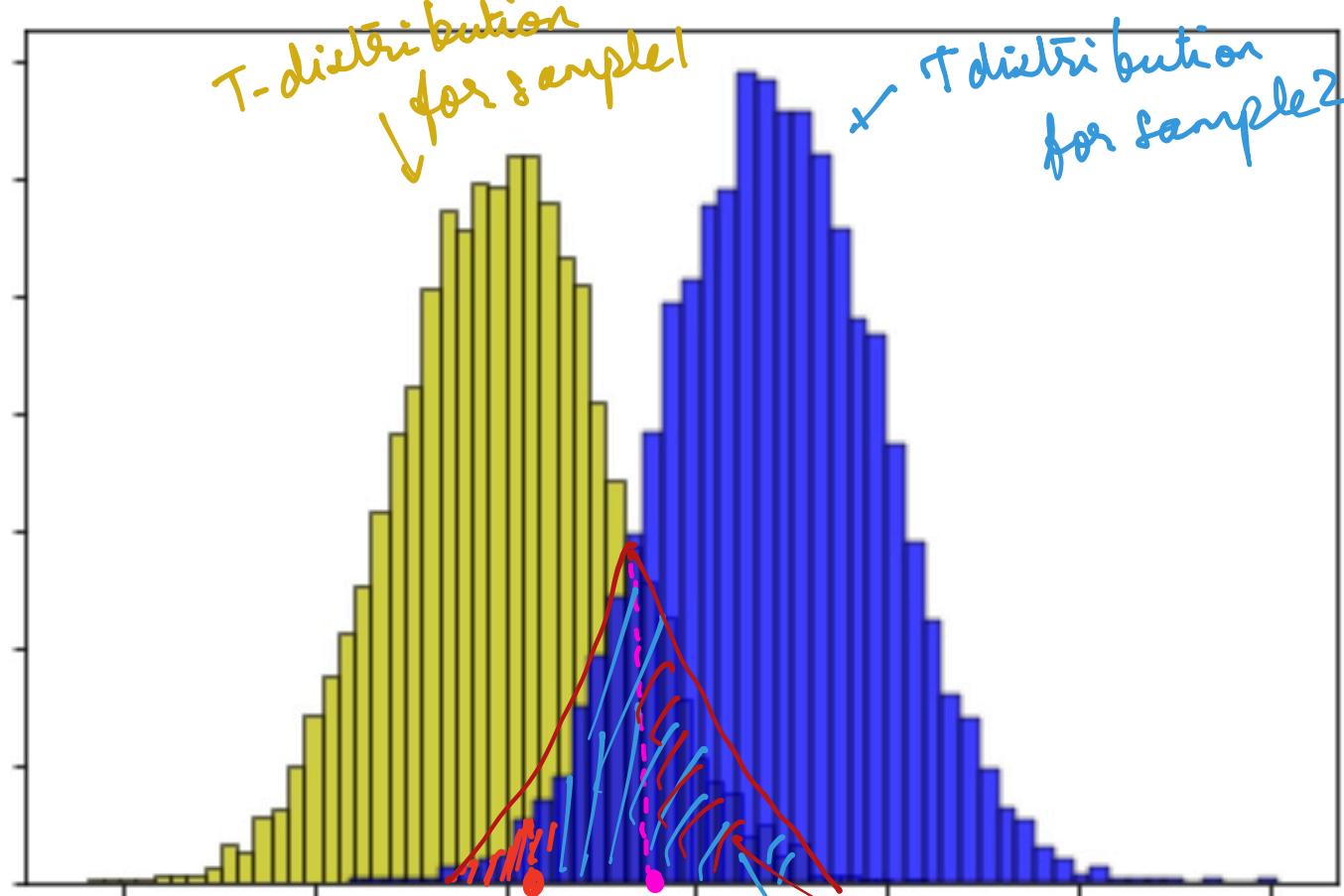
Cohen's effect size.



$$\left( \frac{t_c \times S \times \sqrt{\frac{2}{n}}}{E} \right)^2 = n_{0.1}$$

n that is just enough!

$$E = \frac{\mu_{\text{difference}}}{\sigma_{\text{sample std.}}}$$



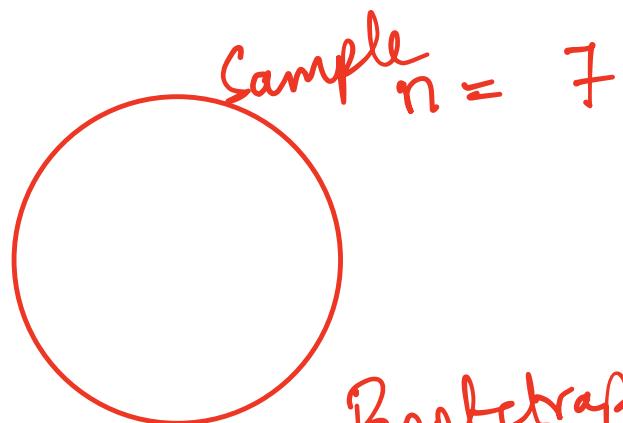
$$\rightarrow \beta_1 < \beta_2.$$

$$\rightarrow \alpha_1 > \alpha_2$$

$$\checkmark \beta = \underline{\underline{0.9}} \checkmark$$

$$\checkmark \alpha = \underline{\underline{0.01}} \checkmark$$

$S.E.$  → defined for a sample.



Sample  $n = 7$

$\sigma \rightarrow$  Don't know.

Bootstrapping → Estimate  $\bar{S} \cdot \leftarrow$

$$S.E. = \frac{\bar{S} \cdot \leftarrow}{\sqrt{n}}$$

We believe that the loan.csv data shows that graduate unmarried men are more likely to get a loan than graduate women.

To prove this, would a t-test be more appropriate or a chi-square test? Carry out the appropriate test on the 'Loan\_Status' column for the two groups and report the p-value. Also report your interpretation.

Note: Assume a confidence level of 5% and round off the p-value to 2 decimal places.

category 1: Graduate unmarried men .

Observations

about whether they got loan or not -

Category 2: Graduate women



Observations about whether they got a loan .

