

Confidence Interval

dist $\leftarrow X$: height

$\{x_1, x_2, x_3, \dots, x_{10}\}$ random sample from X of size 10

Ques estimate population mean of $X = \mu$

$$\mu \approx \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{simple avg}$$

pepl mean } sample mean

as $n \uparrow$, $(\bar{x} \rightarrow \mu)$

$\{180, 162, 158, 172, 168\} \rightarrow \text{height}$
 $\{150, 171, 183, 165, 176\}$

Point estimate of $\mu = \frac{1}{10} \sum_{i=1}^{10} x_i = \underline{168.5 \text{ cm}}$

$\rightarrow \mu \in [162.1, 174.9]$ with 95% probability

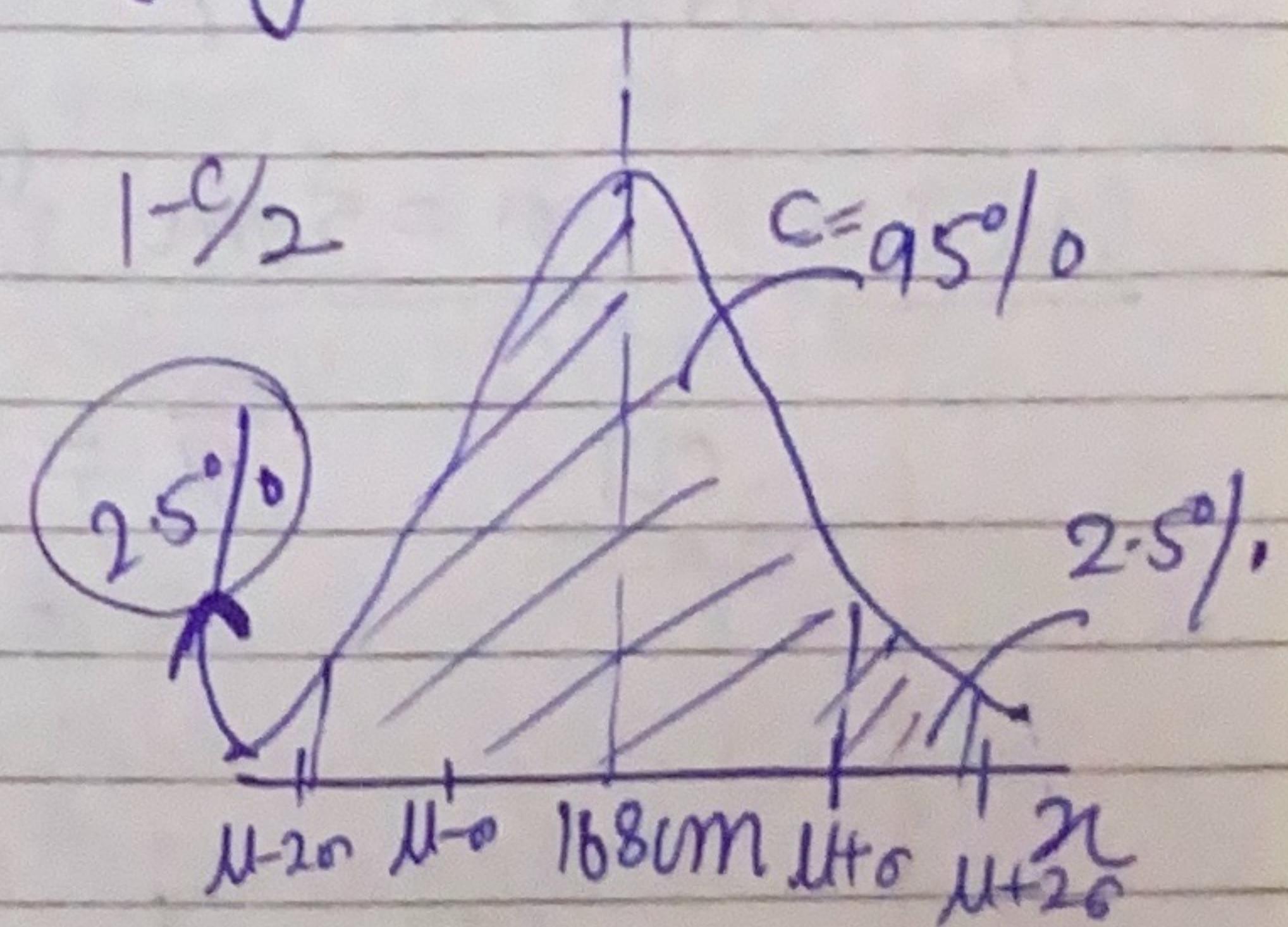
↓ ↓

interval confidence

Computing C.I given underlying dist

heights $\leftarrow X \sim N(\mu, \sigma^2)$ $1 - \frac{\alpha}{2}$

$$\text{dist} \quad (\mu = 168\text{cm}, \sigma = 5\text{cm})$$



Ques C.I of heights

$(\mu - 2\sigma, \mu + 2\sigma) \leftarrow$
95% of obs lies in

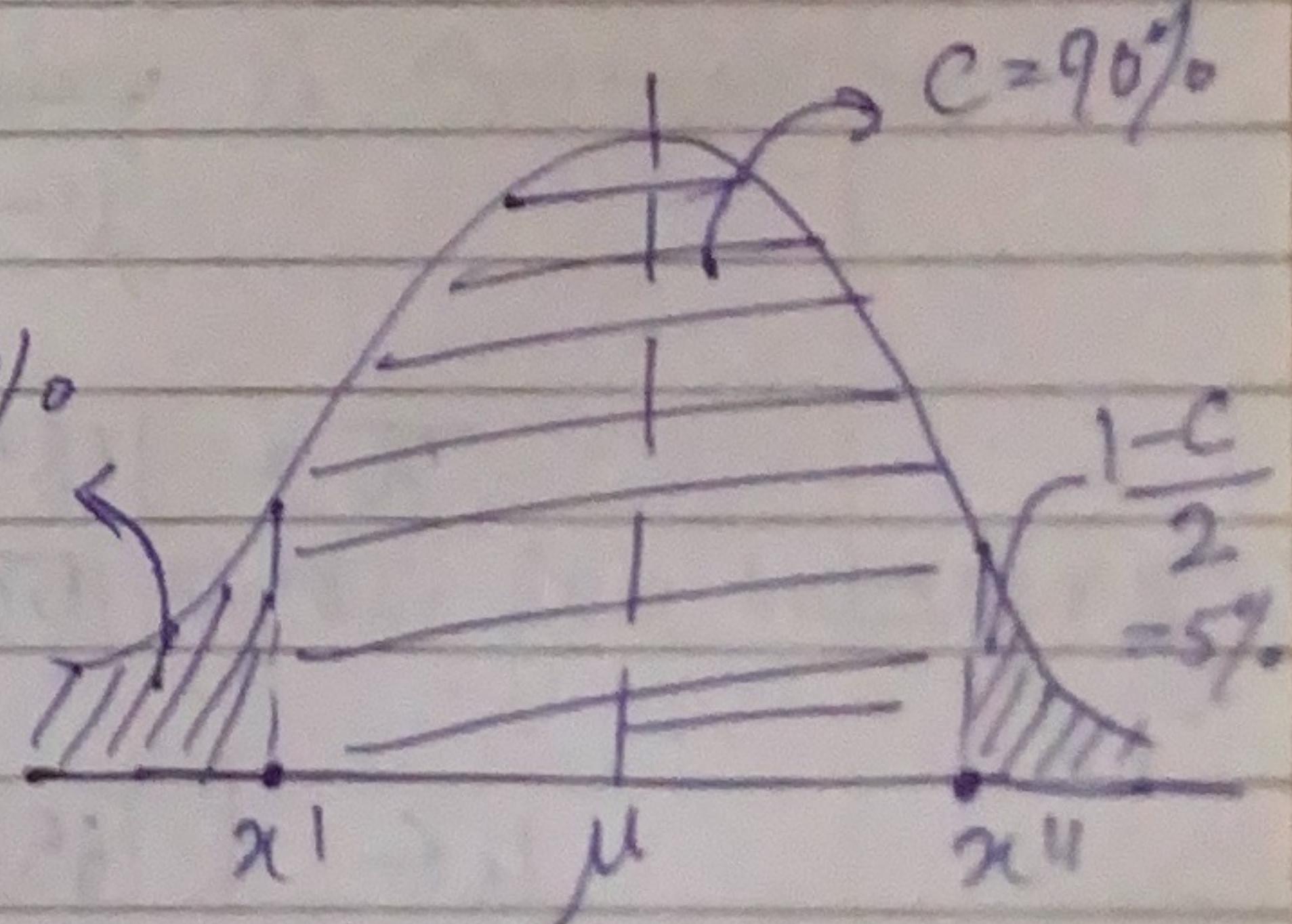
$$\frac{1 - 0.95}{2} = 2.5$$

$(158, 178)$ with 95% prob

C.I

lie in $[x_1, x_{11}]$ with 90% confidence 5%

lower bound upper bound



C.I for mean (μ) of a size

$X \sim N(\mu, \sigma^2)$ with pop mean of μ & std dev σ

$\{x_1, x_2, \dots, x_{10}\} \rightarrow$ sample of size $n=10$
 $\{180, 162, 158, 172, 168, 150, 171, 183, 165, 176\}$

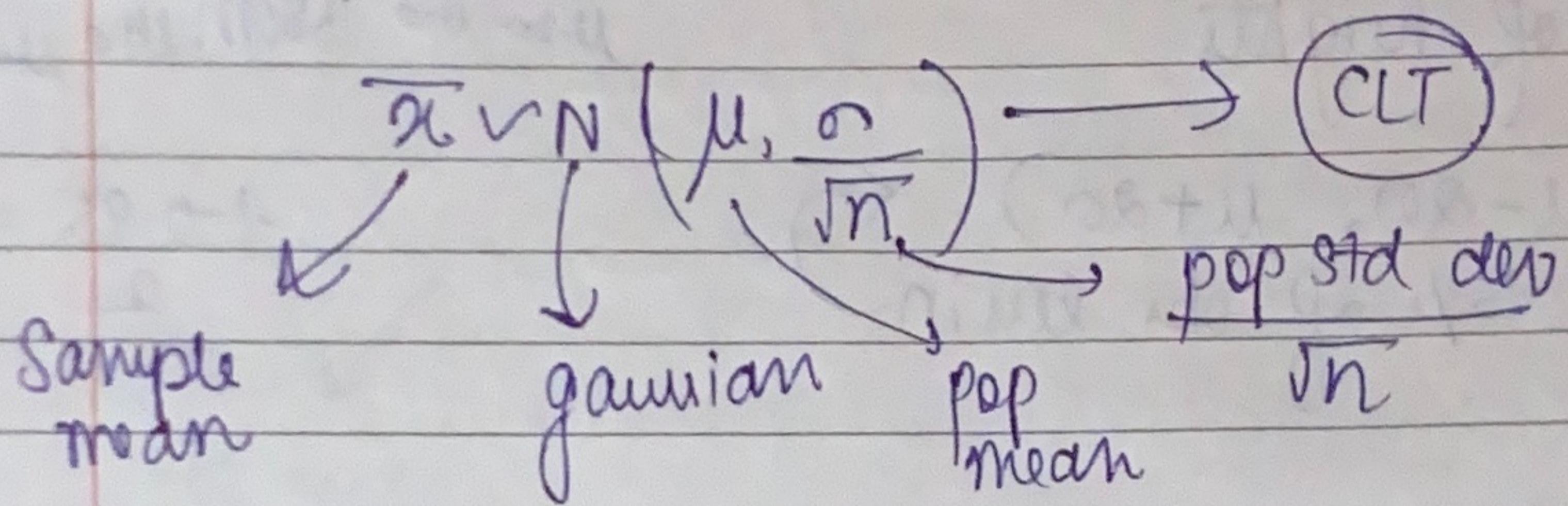
$$\bar{x} = 168.5\text{cm}$$

Q What is 95% C.I. of μ .

Case I $\sigma = 5\text{cm}$ {We know population std dev}

CLT $\bar{x} = \text{Sample mean}$

$$= \frac{1}{10} \sum_{i=1}^{10} x_i$$



$$\mu \in \left[\bar{x} - \frac{2\sigma}{\sqrt{n}}, \bar{x} + \frac{2\sigma}{\sqrt{n}} \right] \text{ with } 95\% \text{ confidence}$$

$\mu - 2\sigma \quad \mu + 2\sigma$

$$\bar{x} = 168.5 \quad \sqrt{n} = \sqrt{10}$$

$$\sigma = 5\text{cm}$$

$$\mu \in [165.34, 171.66] \text{ with } 95\% \text{ confidence}$$

Case II If we don't know σ (pop std dev)

Sample (n)

t -dist^b \rightarrow Student's t -dist

$\checkmark \bar{x} \sim t(n-1)$
sample mean \hookrightarrow degree of freedom

Estimate C.I. of μ of a.r.v.

→ Case I: σ is known CLT; $N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$

→ Case II: σ is unknown t -dist $(n-1)$

(Q)

estimate C.I. of σ of a.r.v

median

90th percentile

(Computation) → (bootstrap C.I.)

C.I. using empirical bootstrap

computational

→ C.I. for median, var, std dev, 90th percentile

→ computer → programming & simulation

60-70 yrs.

→ XNF task: estimate 95% CI for median of X

X:

$S \supset$

Sample of size n: $\{x_1, x_2, \dots, x_n\} \quad n=10$

using only this

↳ C.S. of median of

$s_1: x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots, x_m^{(1)} \quad m \leq n$

↳ random sample of size m

Sampling with replacement from S

$$S = \{x_1, x_2, \dots, x_n\}$$

↓ using sampling with replacement

bootstrap

samples

$$S_1: (x_1^{(1)}, x_2^{(1)}, \dots, x_m^{(1)}) \rightarrow m_1 \rightarrow \text{median of } S_1$$

$$S_2: x_1^{(2)}, x_2^{(2)}, \dots, x_m^{(2)} \rightarrow m_2$$

$$S_3:$$

⋮

$$S_k: x_1^{(k)}, x_2^{(k)}, \dots, x_m^{(k)} \rightarrow m_k$$

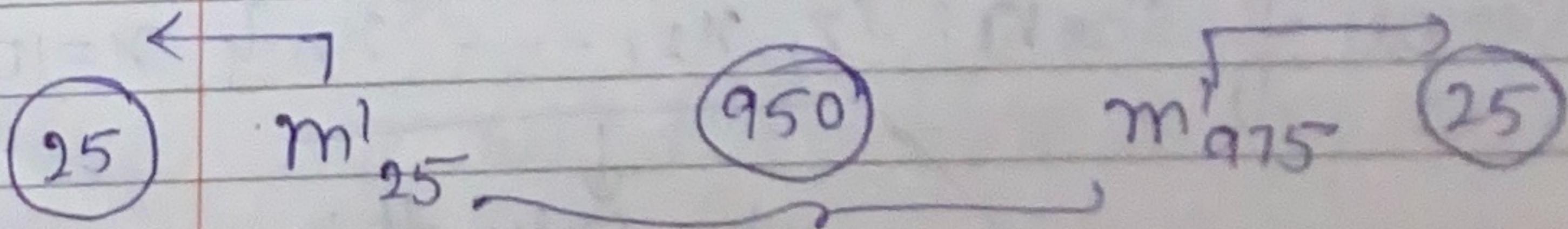
$m_1, m_2, \dots, m_{1000} \leftarrow$ 1000 medians

generated using
bootstrap samples

↓ Sort

$$m'_1 \leq m'_2 \leq m'_3 \leq \dots \leq m'_{1000}$$

↓ CI (95%)



95% CI of median of x is

$$[m'_{25}, m'_{975}]$$

Non-parametric Technique

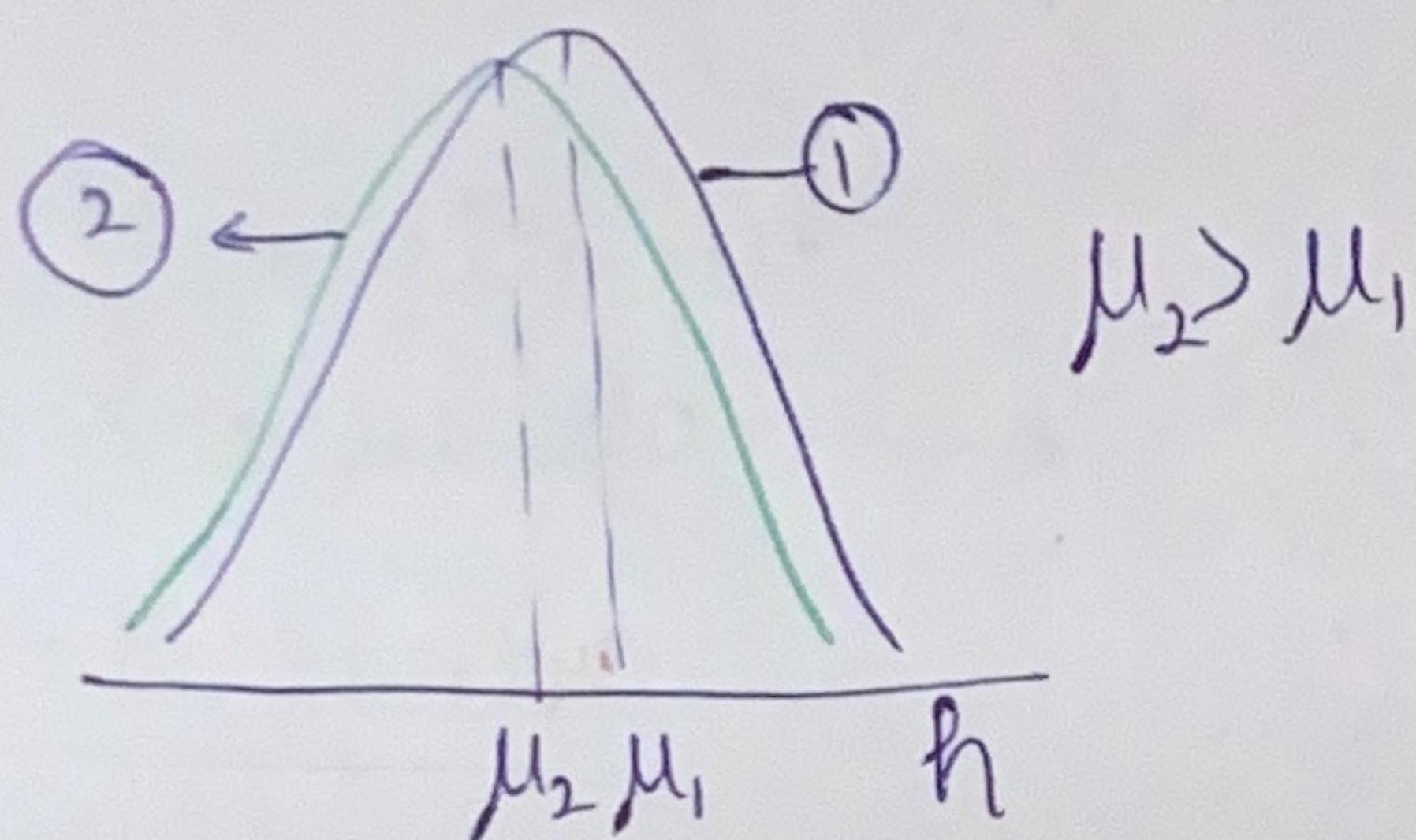
↓ doesn't make any assumptions about the distribution

Probability And Statistics

①

Hypothesis Testing

(Q) Is there a difference in heights of students in Cl₁ and Cl₂



	Cl ₁	Cl ₂
1	160	162
2	152	156
	⋮	⋮
50	148	182

① Choosing a test-statistic

$$n = (\mu_2 - \mu_1)$$

$$n = 10\text{cm}$$

② Null hypothesis (H_0) (Proof By Contradiction)

H_0 : no difference in μ_1 & μ_2

Alternative hypothesis (H_1): diff in μ_1 & μ_2

③ P-value: probability of observing $(\mu_2 - \mu_1)$ if null hypothesis is true

assume H_0 is true.

Cl₁ Cl₂
⑤₀ ⑤₀

if P-value = 0.9

\Rightarrow prob of 10cm is 0.9 if H_0 is true

H_0 \leftarrow if P-value is 0.05
 accept H_1 \rightarrow 5% chance that diff is 10cm
 if H_0 is true

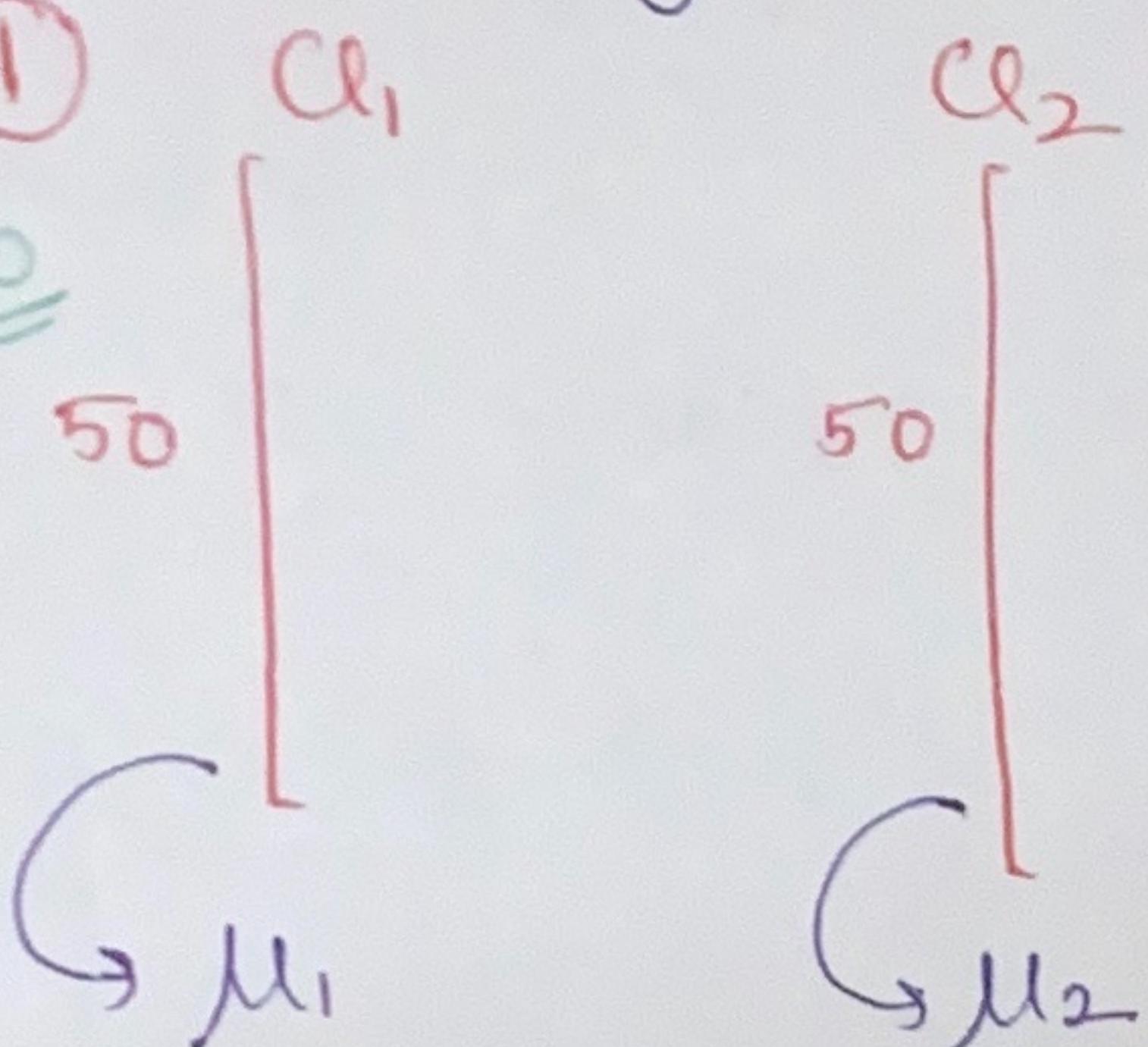
Hypothesis Testing & Permutation Test / Resampling

↳ Same example

→ focusing on important points

①

Step 0



we already made
← observation
actual data

Sample size = 50

$$\Delta/x: \mu_2 - \mu_1 = 10\text{cm} \quad \text{observation}$$

② What is prob of observing a value of $x \geq 10\text{cm}$ if there was no difference in class heights?

$$\text{prob}(x \geq 10\text{cm} \mid H_0)$$

H_0

→ null hypothesis

p-value

if p-value is small = 0.01 or 1%

\downarrow $\rightarrow < 5\% \rightarrow$ customary

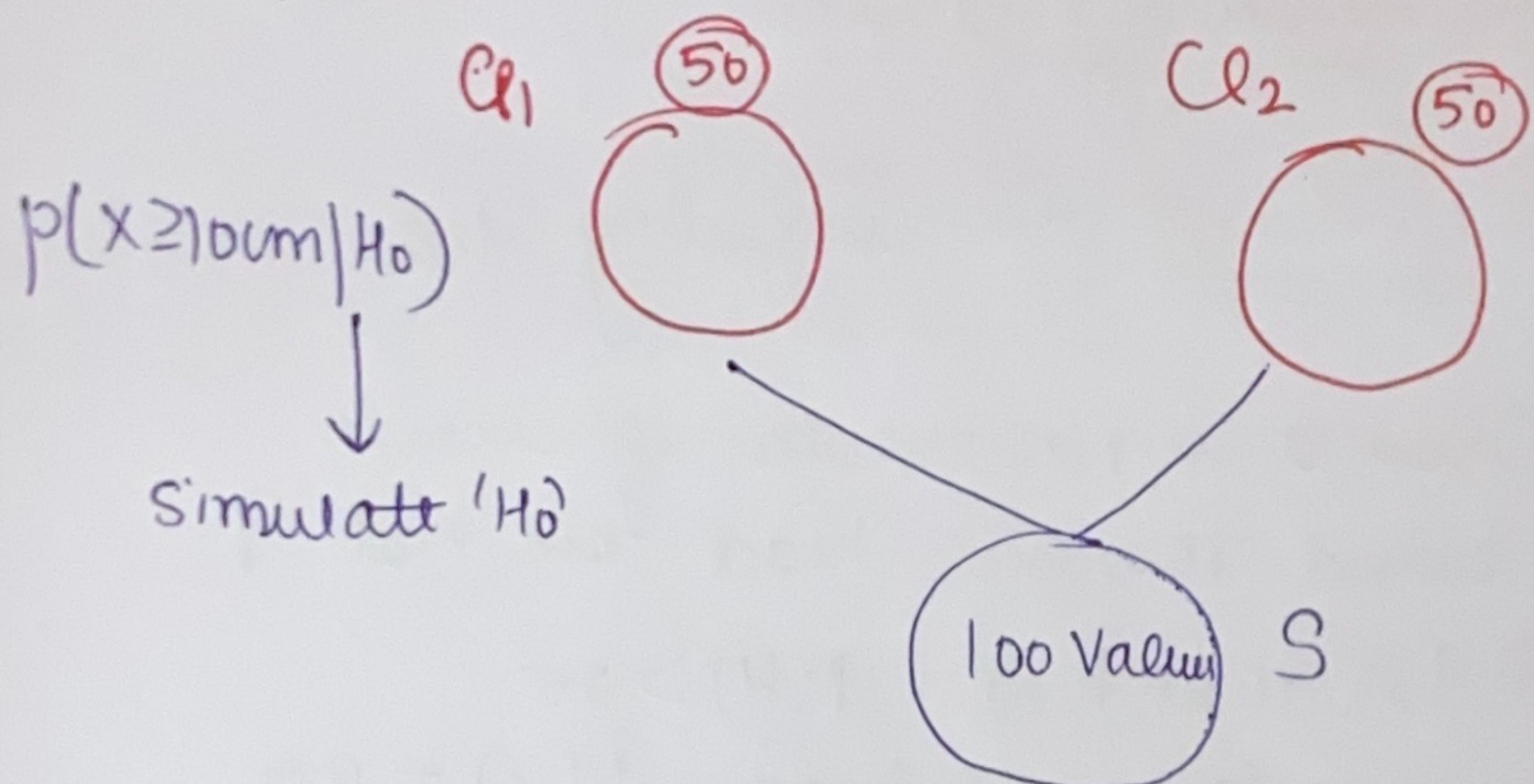
$$P(x \geq 10\text{cm} \mid H_0) = 0.01/1\%$$

(2)

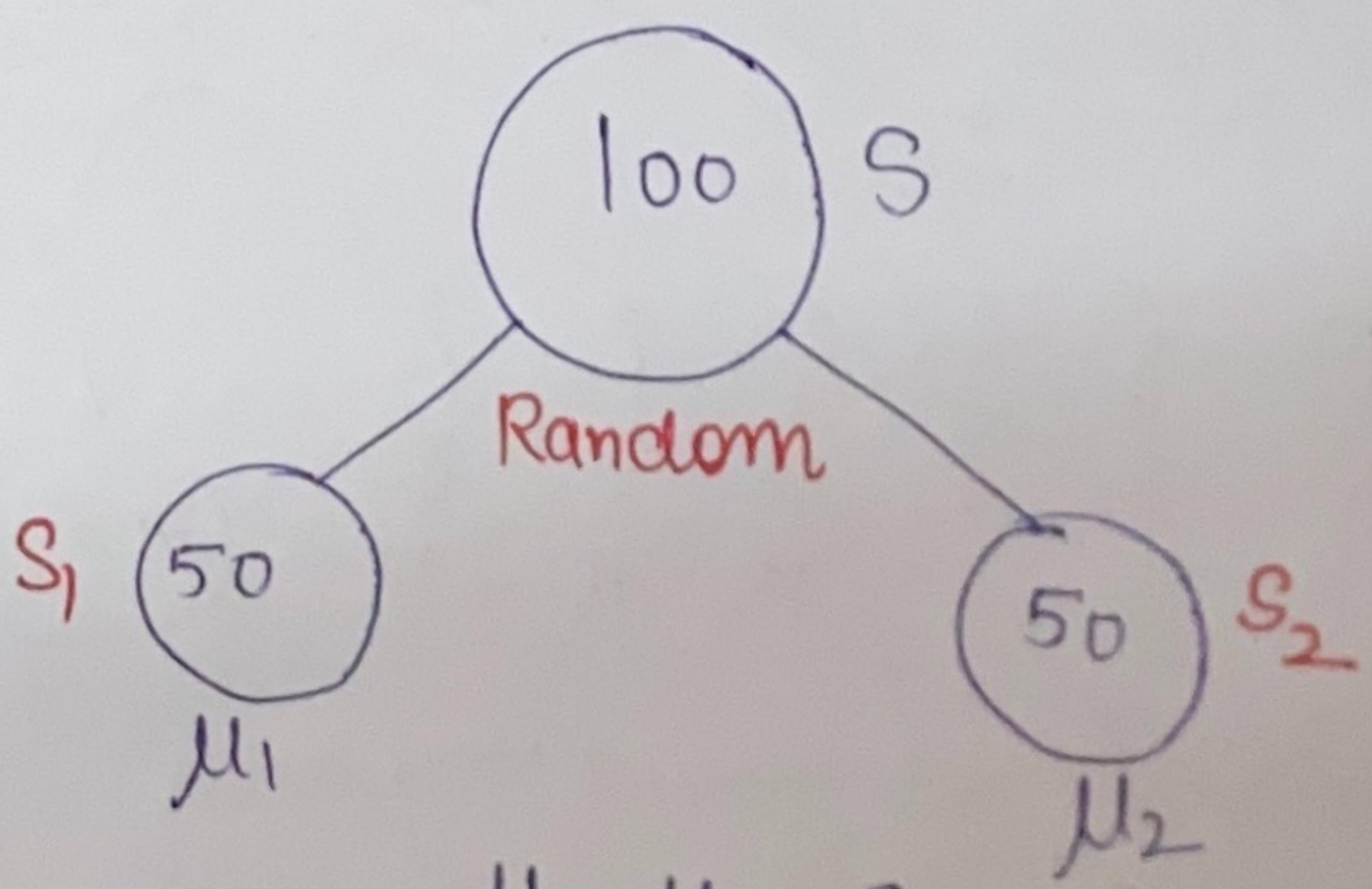
$x=10\text{cm}$ or $x \geq 10\text{cm} \leftarrow \text{True}$

H_0 is less probable or H_0 may not be true

P Value



if H_0 was true \rightarrow no diff in class heights



$\mu_2 - \mu_1 = \delta_1 =$ Simulated diff in class heights with sample size of 50 & H_0 true

$s_1, s_2, s_3, \dots, s_{10K} \leftarrow$ Simulated $s \times 5$ when H_0 being true

\downarrow Sent
 $(s'_1, s'_2, s'_3 - 10\text{cm}, \frac{s'_1}{s_{10K}})$

$p(x \geq 10\text{cm} | H_0) = 2000 = 5.0\%$

$$S' \subseteq S_2' \stackrel{9900}{\subseteq} S_3' \subseteq \dots \stackrel{100}{\subseteq} 10\text{cm} \subseteq S_{10K}'$$

$$P(X \geq 10\text{cm} | H_0) = \frac{100}{10K} = 1\%$$

H_0 ← easier to simulate

Hypothesis Testing ← confusing idea

example: Given a coin determine if coin is biased towards head or not?

biased towards head: - $P(H) > 0.5$

not biased towards head: - $P(H) = 0.5$

experiment: flip a coin 5 times and count no of heads = X ← Test Statistic
H.V

Perform exp f, f, f, f, f
 ↓ ↓ ↓ ↓ ↓
 H H H H H \downarrow
 $X = 5$ ← obs

$P(X=5) | \text{coin not biased towards heads} = P(\text{obs} | H_0)$

assumption
null hypothesis (H_0)

H_0 : coin is not biased towards head

$$P(X=5 | H_0) = \frac{1}{2^5} = \frac{1}{32} \approx 0.03 = 3\%$$

5 heads in
5 tosses

coin is
unbiased towards heads

$$P(X=5 | H_0) = 3\%$$

③

There is a 3% chance of getting 5 heads in 5 flips if coin is not biased towards heads

hyp testing

$$P(\text{obs by exp} | \text{assumption}) = 3\%$$

p value

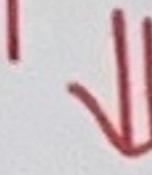
Small

if $P(\text{Obs} | H_0) < 5\%$

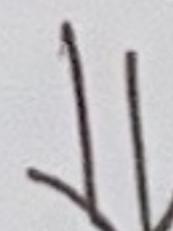
then H_0 may be incorrect



assumption / H_0 is not true



reject $H_0 \Rightarrow$ reject coin is not biased towards head



Coin is biased

Null hypo

H_0 : coin is not biased towards head

Alt hypot

H_1 : coin is biased towards head

{ rejecting $H_0 \Rightarrow$ accepting H_1 }
{ rejecting $H_1 \Rightarrow$ accepting H_0 }

(exp)

flip 3 times

Count no of heads = $X \leftarrow$ test statistic

(perform)

$\begin{matrix} F & F & F \\ H & H & H \end{matrix}$ $X=3 \leftarrow \text{Obs}$

$$P(\text{obs} | \text{assumption}) = \frac{1}{2^3} = \frac{1}{8} = 12.5\% > 5\%$$

↳ cannot biased

accept H_0

★ Imp key points:

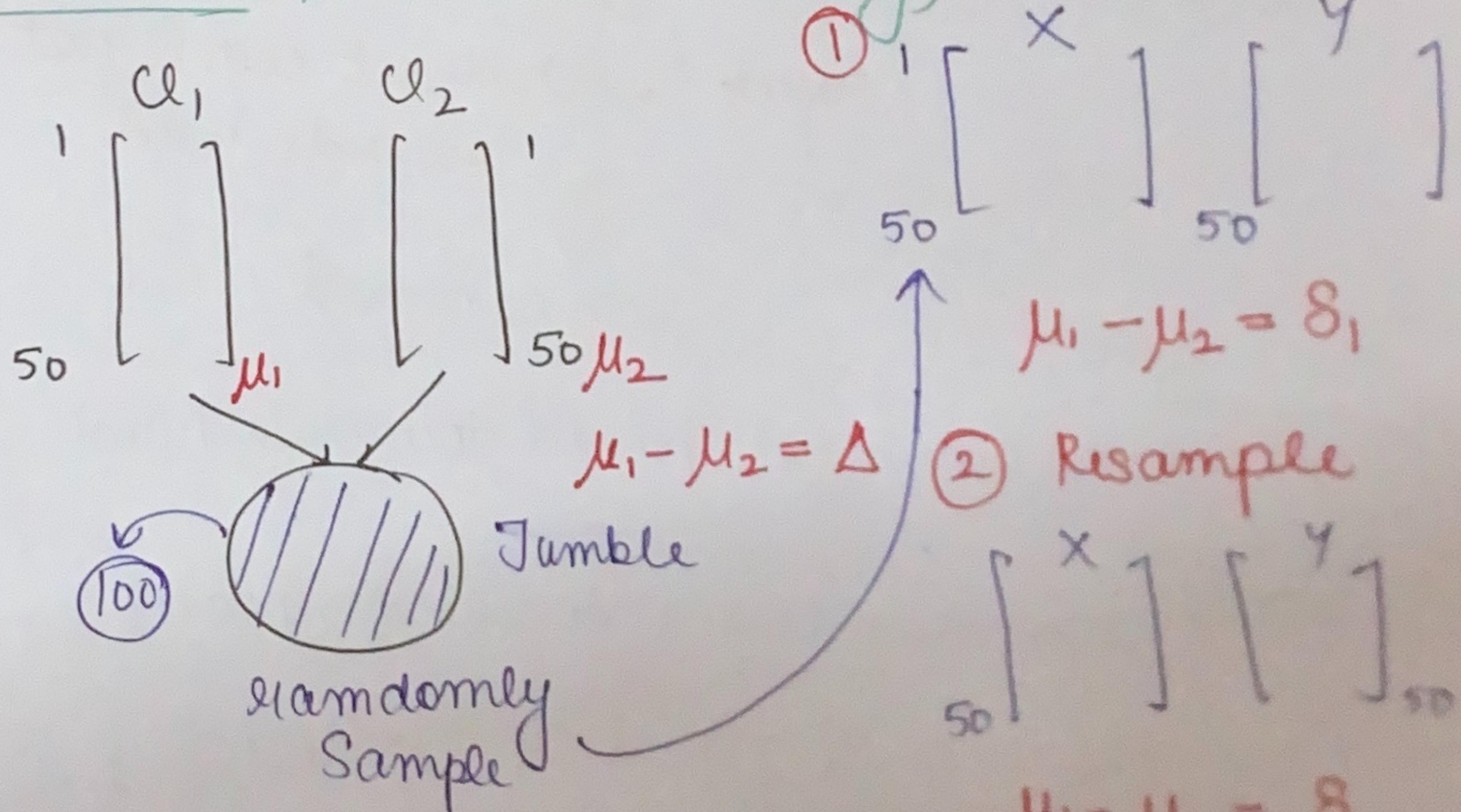
① design of experiment (Sample size)

② H_0 : null hypothesis

$P(\text{obs} | H_0)$ is easy & feasible

★ p value is probability of observation given assumption

pValue: (Permutation Testing)



③ Sort out Δ 's $\Delta = \mu_1 - \mu_2$

$$\Delta_1, \Delta_2, \Delta_3, \dots, \Delta_{10k} - \frac{\Delta}{5\%} \rightarrow \text{Sort out (mul)}$$

$$p \text{ value} = 0.05$$

$10k \rightarrow S_{p \text{ value}}$

p-value computing

(4)

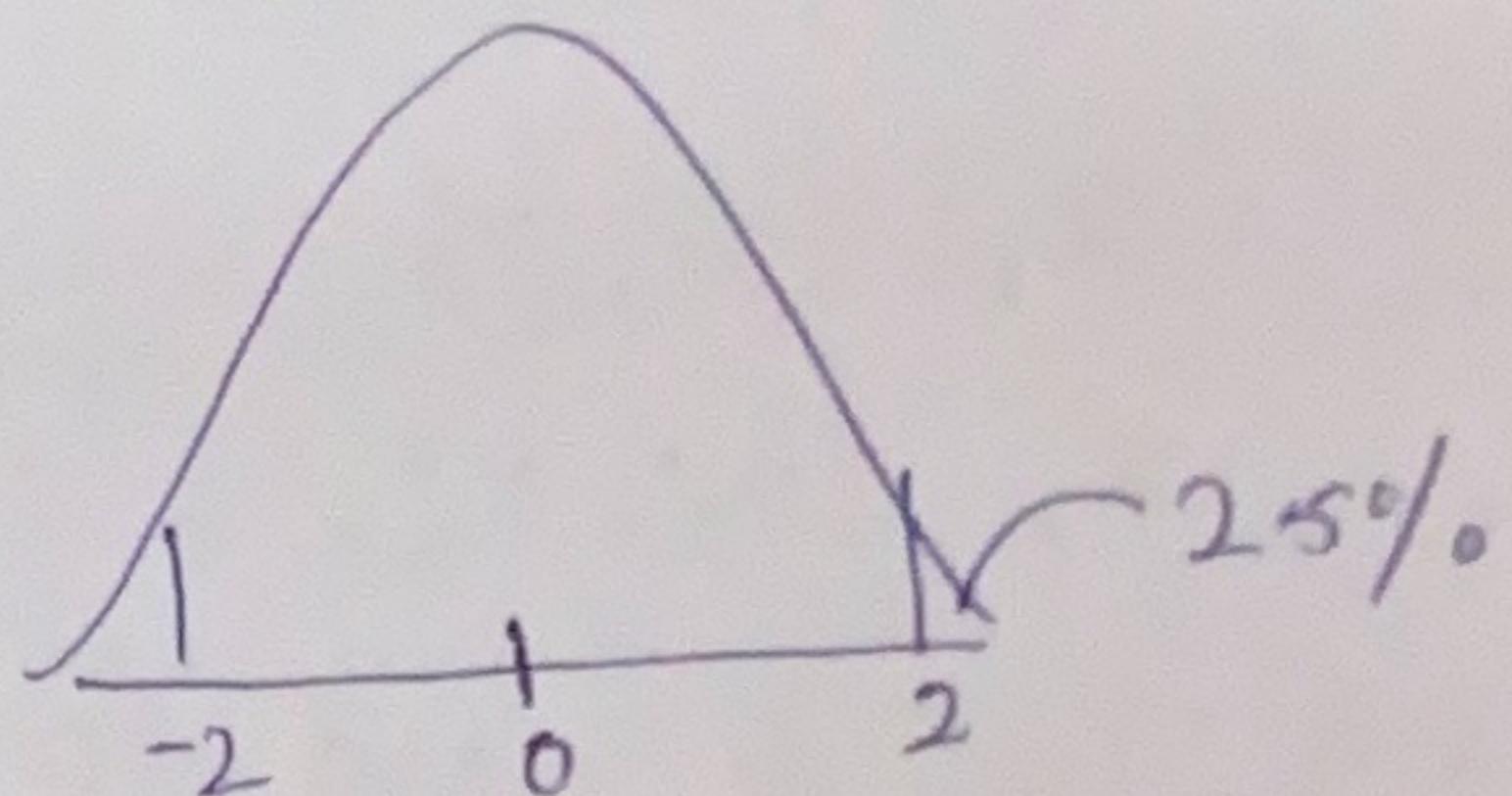
$S_i \rightarrow$

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

$$\Delta = 2$$

$$P(S_i \geq 2) = 0.025$$

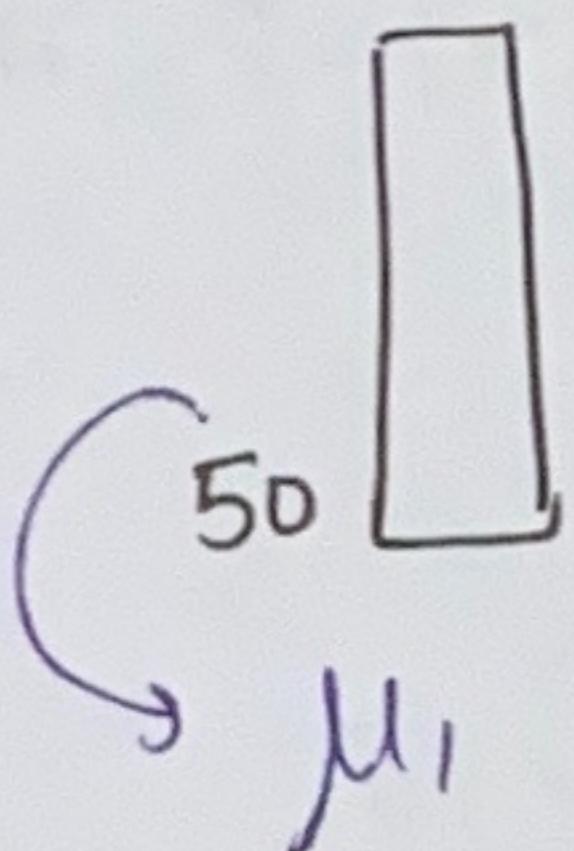
$$p\text{-value} = 0.025$$



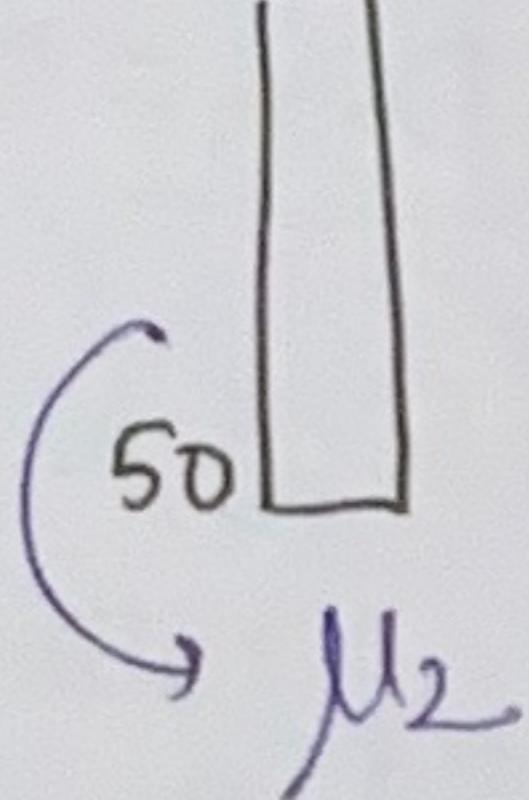
↳ Same example in more detail

①

C_{l_1}



C_{l_2}



we already made
obs (actual data)

$$\text{sample size} = 50$$

~~Step 0~~

$$\Delta \text{ on } x: \mu_2 - \mu_1 = 10 \text{ cm} \xrightarrow{\text{observation}} \text{truth}$$

② What is prob of observing a value of $x \geq 10 \text{ cm}$ if there was no diff in class height?

↳ p-value $P(x \geq 10 | H_0)$ null hypot

If p-value is small = 0.01 or 1%

$\hookrightarrow < 5\% \rightarrow$ significant

$$P(x \geq 10 | H_0) = 0.01 \text{ or } 1\%$$

H_0 is less probable on H_1 may not believe

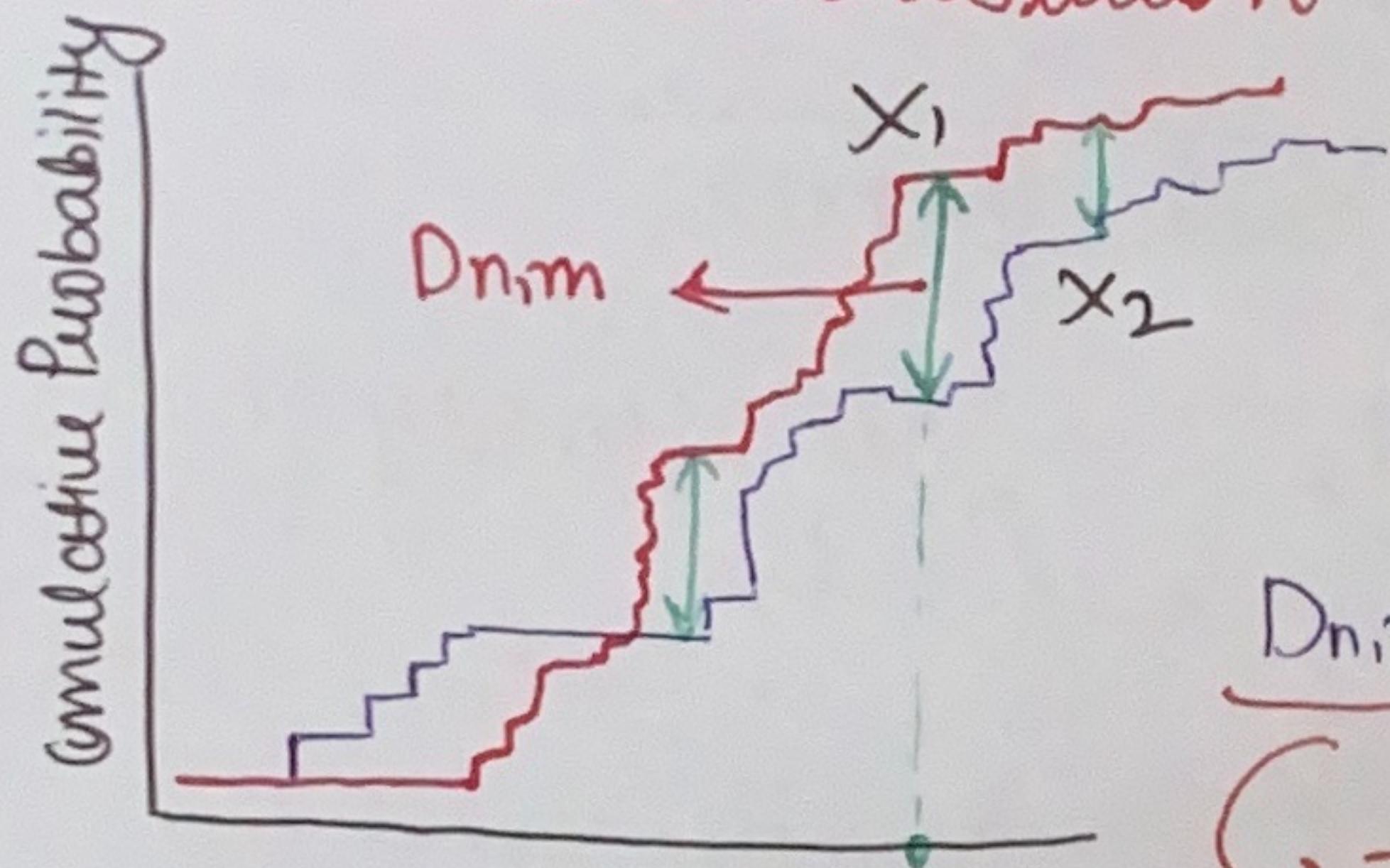
Reject H_0

Two Sample Kolmogorov-Smirnov Test (KS Test)

$$X_1: [x_1^1, x_2^1, \dots, x_n^1]$$

$$X_2: [x_1^2, x_2^2, \dots, x_m^2]$$

Q KS Test is used to determine whether X_1 & X_2 have same distribution



$H_0: X_1 \text{ & } X_2 \text{ have same same distribution}$

$$D_{n,m} = \sup_n |F_{1,n}(x) - F_{2,m}(x)|,$$

G Test Statistic

The null hypothesis is rejected at level α if

$$D_{n,m} > c(\alpha) \sqrt{\frac{n+m}{nm}}$$

$$c(\alpha) = \sqrt{-\frac{1}{2} \ln \alpha}$$

α	0.10	0.05	0.025	0.01	0.005	0.001	Look Up Table
$c(\alpha)$	1.22	1.36	1.48	1.63	1.73	1.95	

if $n=1000$ $m=5000$

$$\alpha=0.05$$

$$D_{n,m} > 0.047$$

then we reject null hypothesis at 0.05 level

② If $n=50$ $m=30$

$$\alpha=0.05$$

$$1.36 \times \sqrt{\frac{50+30}{50(30)}} = 0.31$$

$$D_{n,m} > 0.31$$

reject H_0

①
hypot
Testing

H_0 is rejected at a significance level α ⑤
if $p\text{-value} < \alpha$

②
KS Test

H_0 is rejected if $D_{n,m} > C(\alpha) \sqrt{\frac{n+m}{nm}}$

$$\text{let } D_{n,m} = D$$

$$\text{reject } H_0 \text{ if } D \sqrt{\frac{nm}{n+m}} > c\alpha = \sqrt{-\frac{1}{2} \ln \alpha}$$

$$\Rightarrow D \sqrt{\frac{nm}{n+m}} > \sqrt{-\frac{1}{2} \ln \alpha}$$

$$\Rightarrow D^2 \left(\frac{nm}{n+m} \right) > -\frac{1}{2} \ln \alpha$$

$$= -D^2 \left(\frac{nm}{n+m} \right) < \frac{1}{2} \ln(\alpha)$$

reject H_0 if

$$C = \exp \left[-2D^2 \left(\frac{nm}{n+m} \right) \right] < \alpha$$

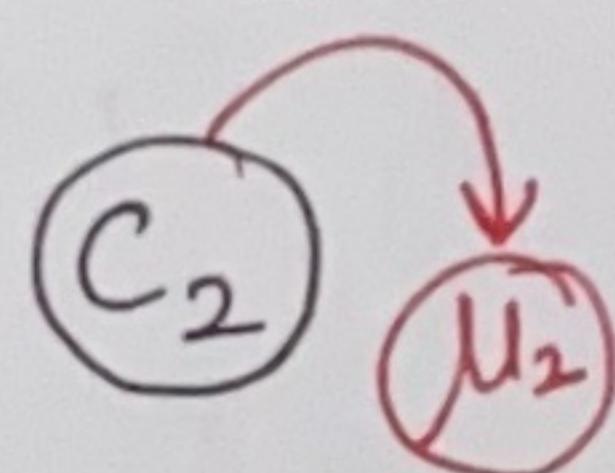
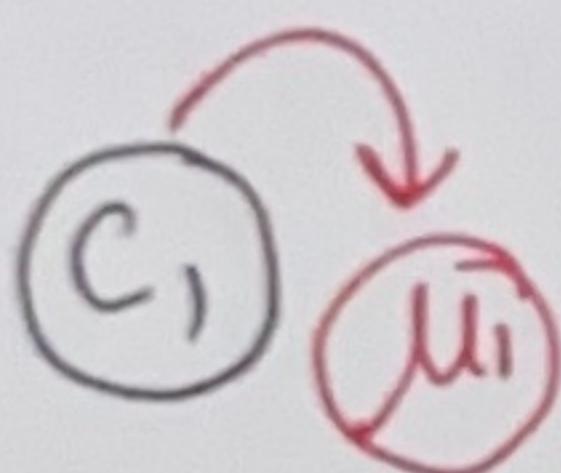
\hookrightarrow p-value for KS Test

Hypothesis Testing: Difference of means

e.g. ① Coin toss \rightarrow simple

$P(\text{obs} | H_0)$: very easy

e.g. ②



Determine if population means of heights of people in these two cities is same or not?

μ_1 / μ_2 are same or different

(exp):

$$\begin{array}{c} C_1 \\ h_1 \\ h_2 \\ \vdots \\ h_{50} \end{array} \quad \begin{array}{c} C_2 \\ h_1 \\ h_2 \\ \vdots \\ h_{50} \end{array}$$

randomly

sample heights of 50 people

$$\mu_1 = \frac{h_1 + h_2 + \dots + h_{50}}{50} \rightarrow 162 \text{ cm}$$

$$\mu_2 = \frac{h'_1 + h'_2 + \dots + h'_{50}}{50} \rightarrow 167 \text{ cm}$$

Test Statistic: $| \mu_1 - \mu_2 | = x = 162 - 167 = 5 \text{ cm}$

Null hypothesis: $H_0 = \text{no difference in population means}$

Compute: $P(x = 5 \text{ cm} | H_0)$

probability of obtaining a difference of sum in sample mean heights of sample size 50 between C_1 & C_2 if there is no difference in mean heights

Case 1: $p(n=5 | H_0) = 0.2 = 20\%$, ⑥

There is 20% chance of observing a diff of 5cm in sample mean heights of C₁, 8C₂ with sample size 50 if there is no pop mean diff.

- ⇒ assumption must be true
⇒ accept H₀

Case 2: $p(n=5 | H_0) = 0.03\% = 3\%$

$p(\text{obs} | \text{assumption}) = 3\%$ ↗ small
⇒ assumption incorrect
⇒ reject H₀ ⇒ accept H₁

Resampling & Permutation Testing

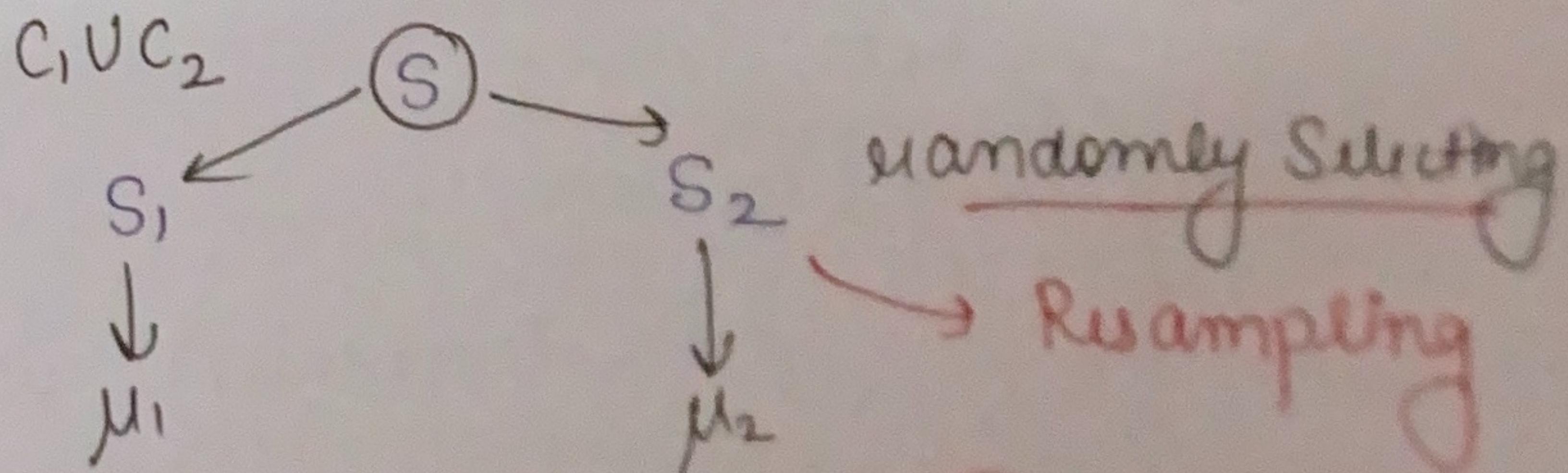
$n = \mu_1 - \mu_2 \leftarrow$ diff in sample mean with sample size of 50

H₀: no diff in population means

Step 1

$$S = \{h_1, h_2, h_3, \dots, h_1', h_2', h_3', \dots, h_{50}'\}$$

Step 2



① $\mu_2 - \mu_1 = 3\text{cm}$ ↗ ⑧₁

② $\mu_2 - \mu_1 = -2\text{cm}$ ↗ ⑧₂

③ $\mu_2 - \mu_1 = +1\text{cm}$ ↗ ⑧₃

④ $\mu_2 - \mu_1 = 6\text{cm}$ ↗ ⑧₄

repeat sampling

K=1000

Step 3

Sort S_i 's

$$S'_1 \leq S'_2 \leq S'_3 \leq S'_4 \dots \leq S'_k$$

inc order

case 1

$$\text{obs-difference} = n = 5\text{cm}$$

$$P(\text{diff} \geq 5\text{cm} | H_0) = ?$$

$$\underbrace{S'_1 \leq S'_2 \leq S'_3}_{80\% \leq 5\text{cm}} \quad - 5\text{cm} \quad \underbrace{S \leq S'_{1000}}_{20\% \geq 5\text{cm}}$$

pvalue $P(\text{obs diff} | \text{assumption}) = 20\% > 5\%$

Case 2

$$\underbrace{S'_1 \leq S'_2 \leq \dots}_{97\%} \quad - 5\text{cm} \quad \dots \leq \underbrace{S'_{1000}}_{3\%}$$

$$P(\text{obs diff} \geq 5\text{cm} | H_0) = 3\%$$

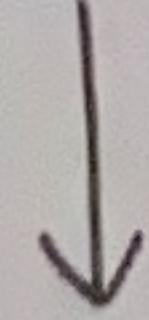
→ assumption incorrect
→ reject H_0

How to use hypothesis-testing

KS test → used if two r.v's have same distribution or not.

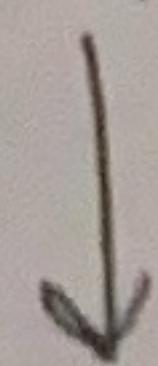
Eg designing drug / medicine

(in market) D_1



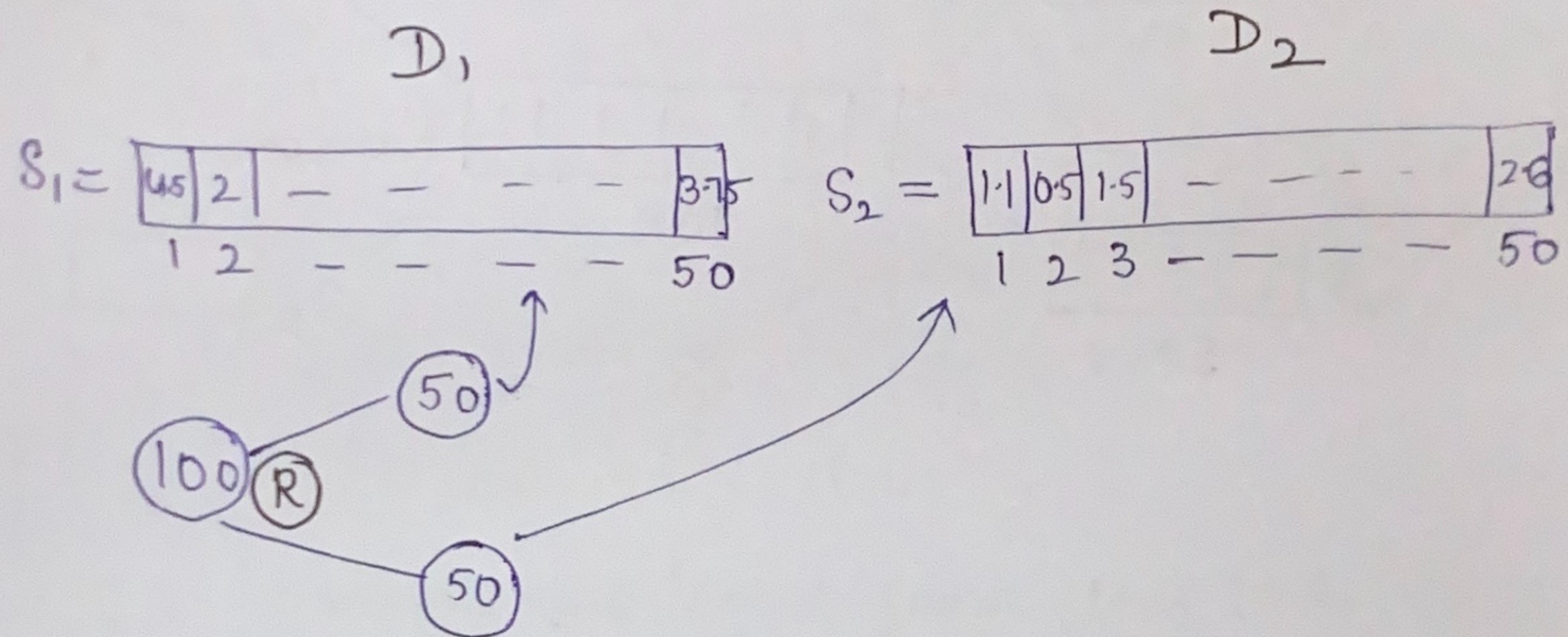
reduces fever
in 4 hours

D_2 (new drug)



claim: reduce fever faster
than drug 1

experiment → to determine if true or not ①
(claim)



$$\mu_1 = \text{mean time} = 4.0 \text{ hrs}$$

$$\mu_2 = \text{mean time} = 2 \text{ hrs}$$

Hyp testing: ① $H_0: D_2 \neq D_1$, are not different
have same time taken to
reduce fewer

② Test - Statistic: $X = \mu_2 - \mu_1 = 4 - 2 = 2 \text{ hrs}$

$P(X \geq 2 | H_0) = \text{very small } 1\%$ obs value

If there is no difference in D_1 and D_2 , the

probability of observing mean diff $\xrightarrow[H_0]{\text{obs}} 2.22$
is very small (1%)

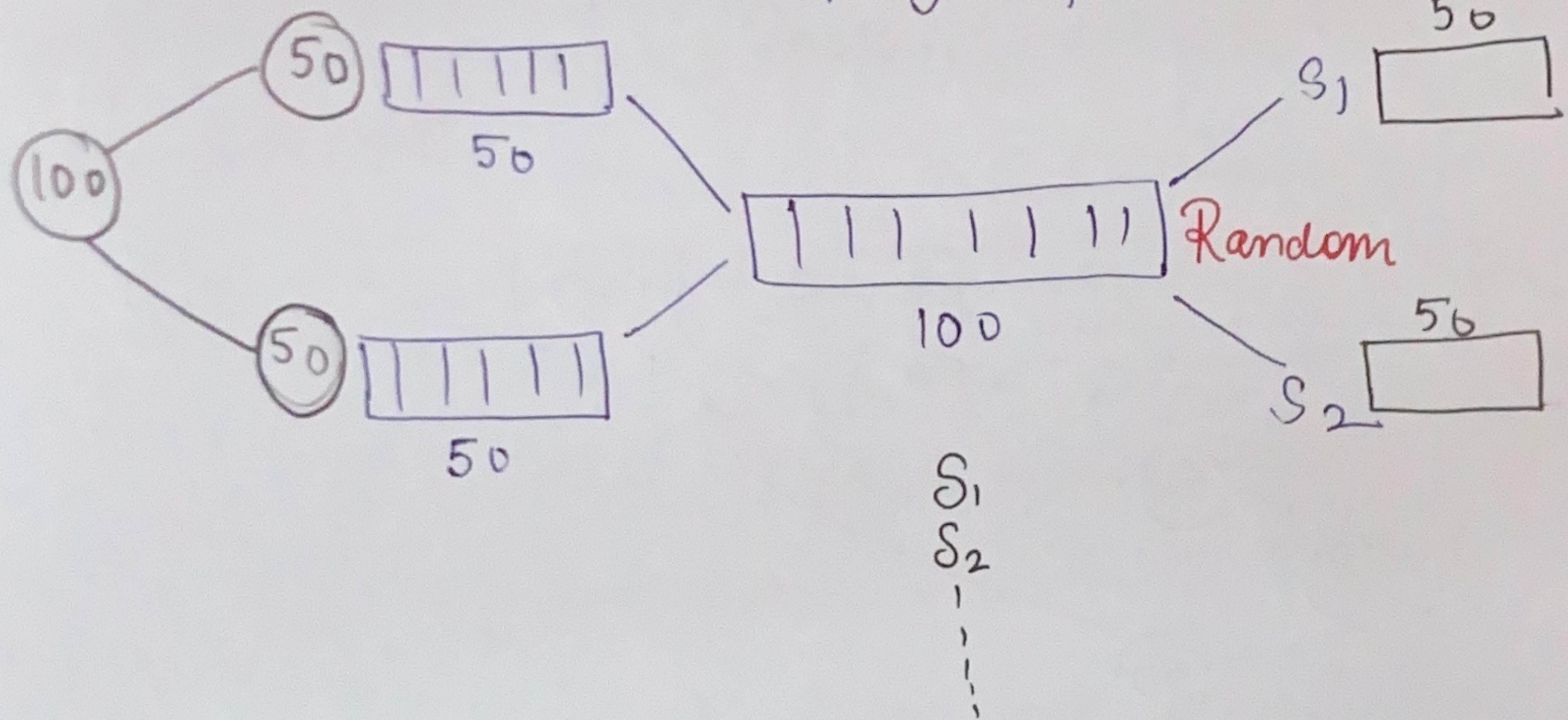


H_0 & obs do not agree with each other



H_0 must be wrong

$P(X \geq 2 | H_0) \rightarrow$ subsampling & permutation test



Proportional Sampling

$d = \frac{d_1}{2.0} \frac{d_2}{6.0} \frac{d_3}{1.2} \frac{d_4}{5.8} \frac{d_5}{20.0}$ randomly

Task: pick an element amongst n elements s.t
the prob of picking an element is proportional
to the d_i 's value

Step 1 ① $S = \sum_{i=1}^n d_i = 35 \leftarrow$ compute the sum

② $d'_i = d_i / S \quad \leftarrow$ normalizing values using sum

$$d'_1 = 0.0571$$

$$d'_2 = 0.171428$$

$$\sum d'_i = \sum \frac{d_i}{S} = 1 \quad d'_3 = 0.0343$$

$$d'_4 = 0.1657$$

$$d'_5 = 0.5714$$

} 0 to 1
Sum to 1

© cumulative normalized Sum

(8)

$$\text{d}_3 = \text{sum} \left\{ \begin{array}{l} d_1 = 0.0571 \\ d_2 = 0.171428 \\ d_3 = 0.0343 \\ d_4 = 0.1657 \\ d_5 = 0.5714 \end{array} \right\} \rightarrow \left| \begin{array}{l} \tilde{d}_1 = d_1 = 0.0571 \\ \tilde{d}_2 = \tilde{d}_1 + d_2 = 0.228528 \\ \tilde{d}_3 = \tilde{d}_1 + \tilde{d}_2 + d_3 = 0.262828 \\ \tilde{d}_4 = 0.428528 \\ \tilde{d}_5 = 1.00 \end{array} \right.$$

\tilde{d}_i = cum-nor-values

Step 2 Sample one value from $\text{unif}(0.0, 1.0)$

$\gamma = \text{numpy.random.uniform}(0.0, 1.0, 1)$

det

$\gamma = 0.6$

range
↓
no of values

Step 3 proper-sampling

if $\gamma \leq d_1$

return 1

else if $\gamma \leq d_2$

return 2

else if $\gamma \leq d_3$

return 3

⋮

⋮

\tilde{d}_1

\tilde{d}_2

\tilde{d}_3

\tilde{d}_4

\tilde{d}_5

$\left. \begin{array}{l} \tilde{d}_1 \\ \tilde{d}_2 \\ \tilde{d}_3 \\ \tilde{d}_4 \\ \tilde{d}_5 \end{array} \right] \leftarrow \gamma$

prob of γ lying b/w \tilde{d}_3 & $\tilde{d}_4 = d_4 \propto d_4$

$d_4 = d_4/S$