

Wrap up: CLT + CI

$\xrightarrow{\text{--}}$ Gaussian, know σ
 $\xrightarrow{\text{--}}$ dataset "Bootstrap"

Hypothesis testing (Q to 3)

$\xrightarrow{\text{Tricky}}$ $\xrightarrow{\text{transf}}$ Normal
 $\xrightarrow{\text{intervew}}$ "p-value" \rightarrow later

~~Recall:~~

$$X \rightarrow E[X] = \mu, \quad \text{Var}[X] = \sigma^2$$

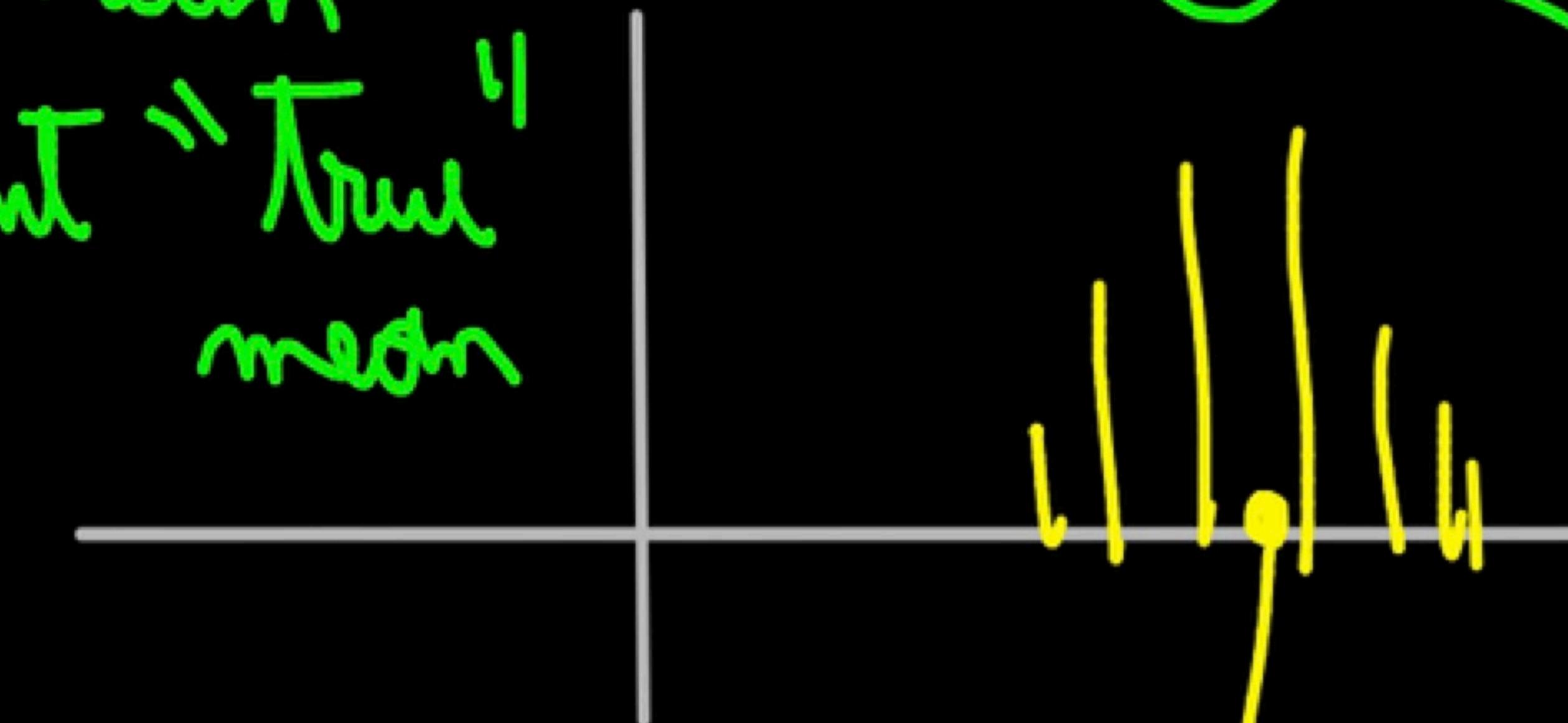
→ CLT: Sum of $X_1 + X_2 + \dots + X_n \rightarrow \text{Behaviour}$

Mean $\bar{X} = X_1 + X_2 + \dots + X_n \rightarrow \text{Behaviour}$

sample mean

(constant "true" mean)

"Stand Gaussian"



"Gaussian"

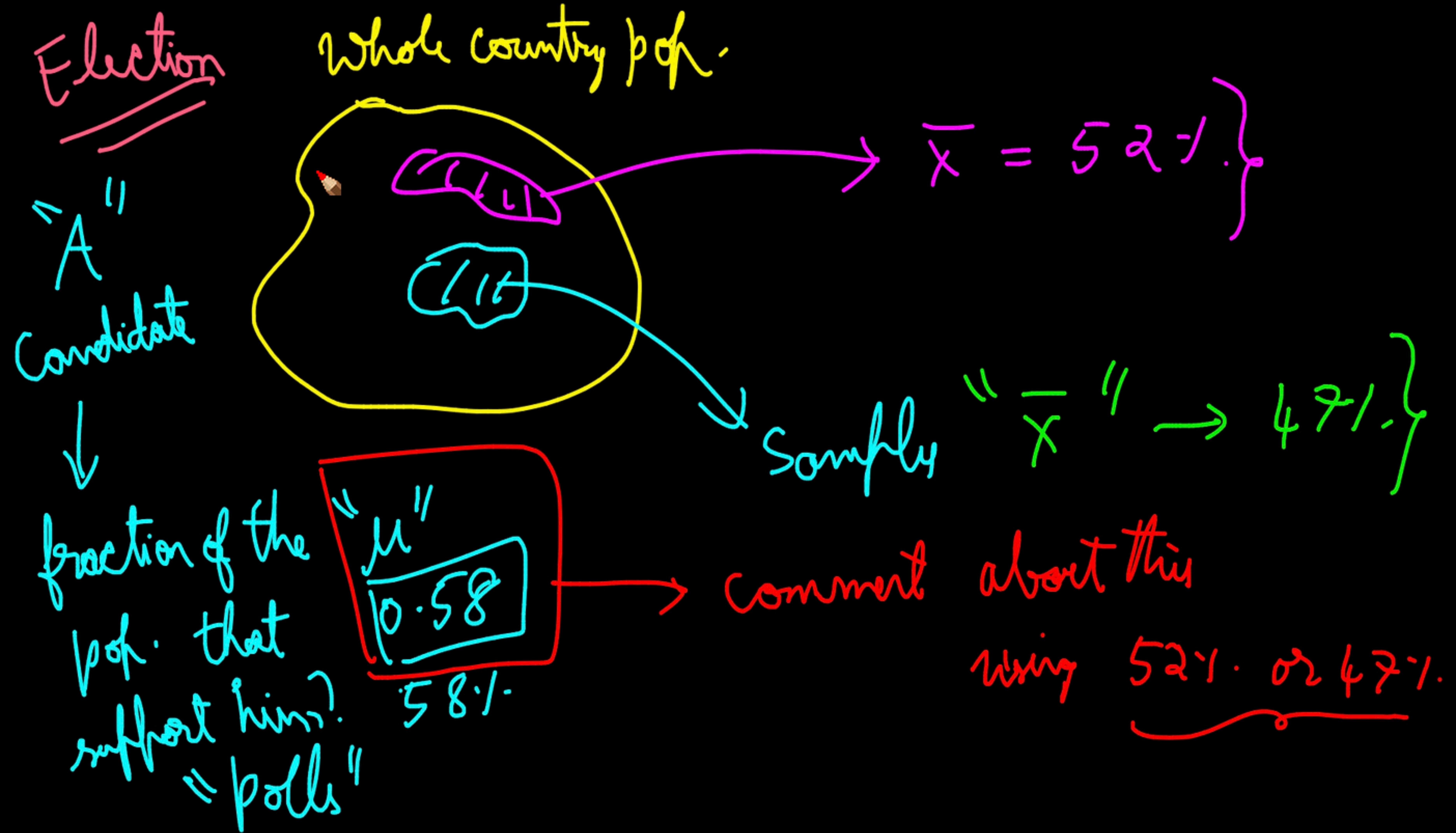
sample

$$Z = \frac{\bar{X} - E[\bar{X}]}{\sqrt{\text{Var}[\bar{X}]}} \rightarrow E[Z] = 0$$

$$\text{Var}[Z] = 1$$

$$E[\bar{X}] = E[X] = \mu$$

$$\text{Var}[\bar{X}] = \frac{\text{Var}[X]}{n} = \frac{\sigma^2}{n}$$



We got sample: \bar{x}

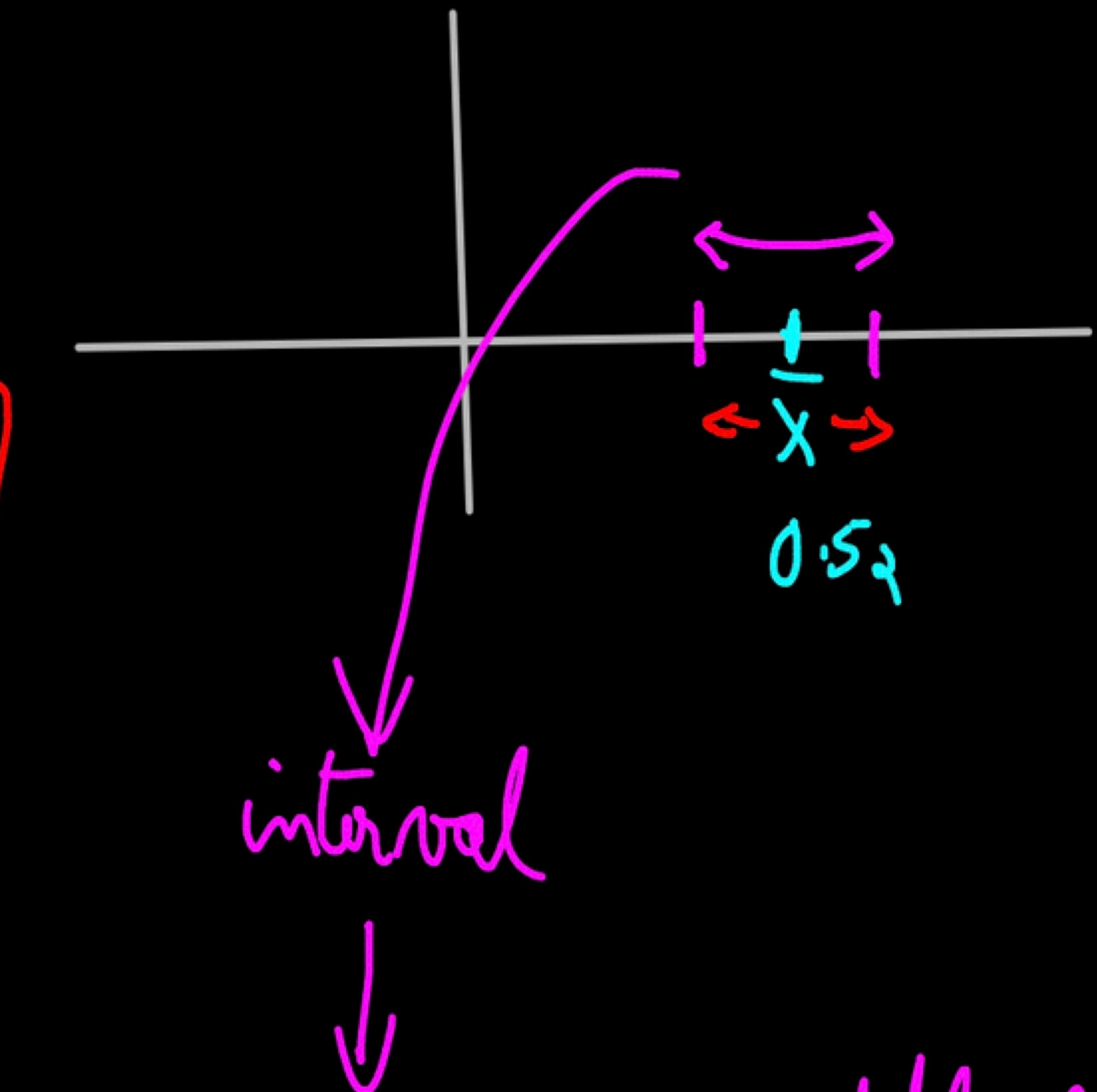
"Population" mean: μ

Goal is to say: $\mu \in \text{interval}$

95%

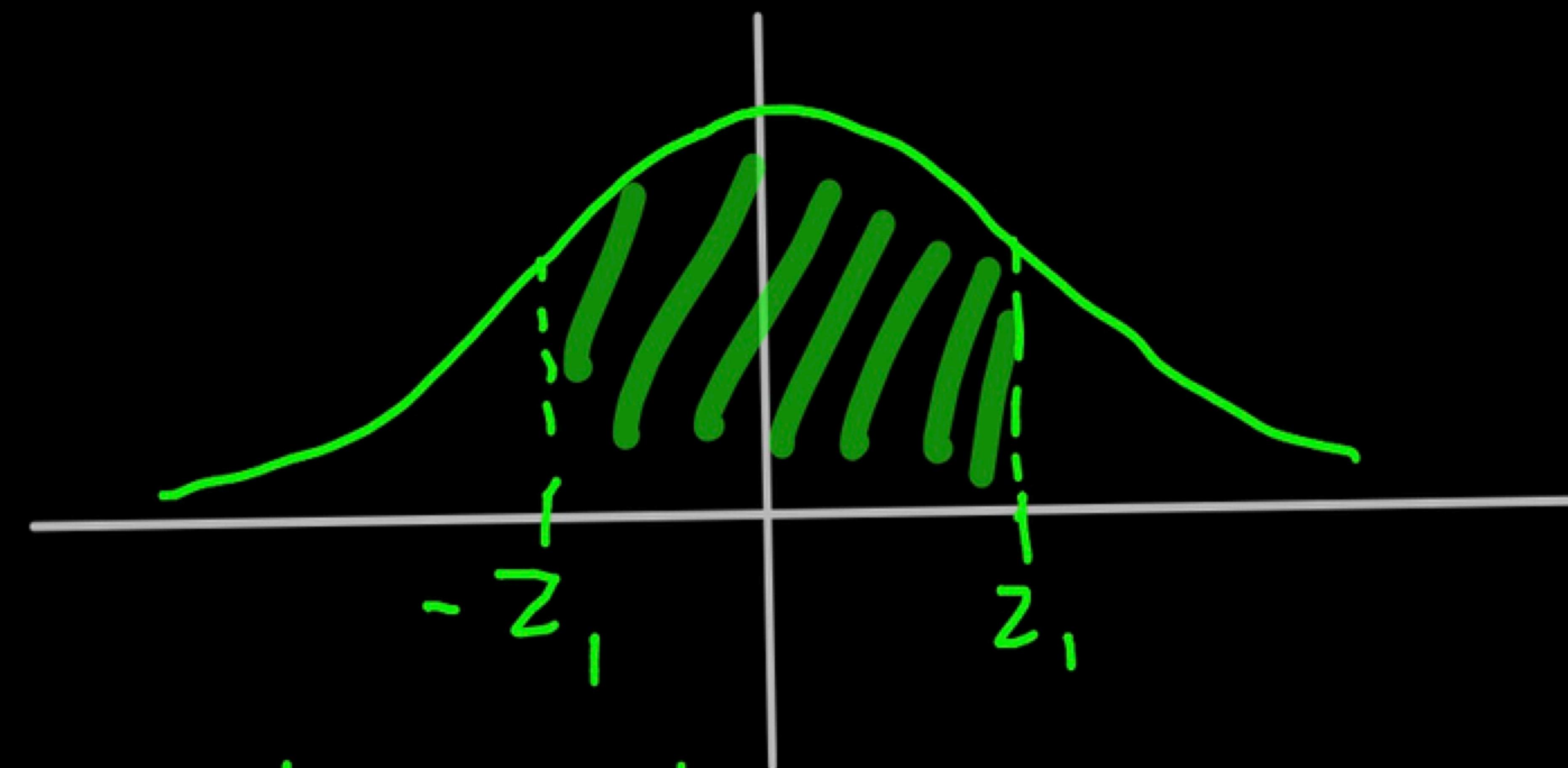
$$\mu \in \left[\bar{x} - \frac{1.96\sigma}{\sqrt{n}}, \bar{x} + \frac{1.96\sigma}{\sqrt{n}} \right]$$

If " σ " is known

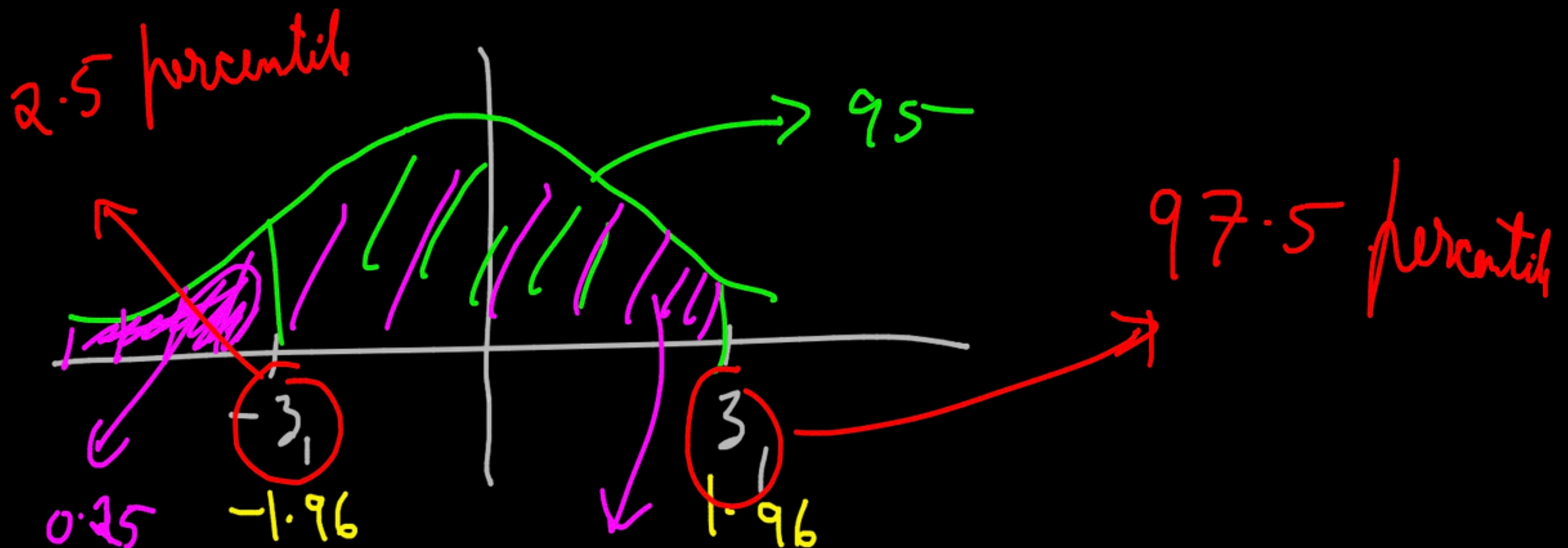


\bar{x} lies in the middle of the interval

What if σ is not known: "Bootstrap"



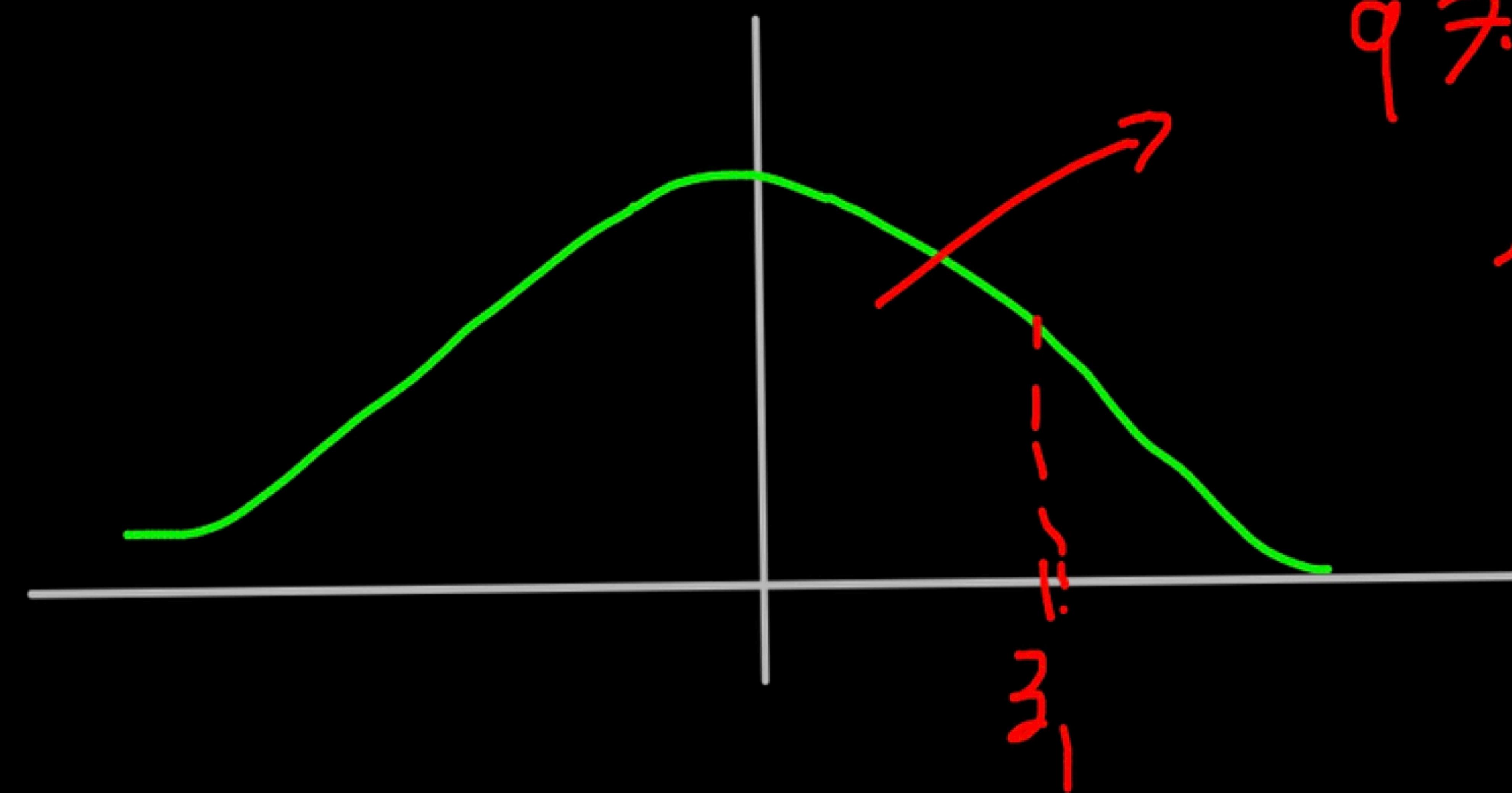
z_1 : That value such that 95%. area is in
 $[-z_1, z_1]$.
norm.ppf(0.975)



$$95 + 0.25 = 0.975$$

z_1 is that value such that $\text{norm.cdf}(z_1) = 0.975$

$$\underline{\underline{z_1 = \text{norm.ppf}(0.975)}}$$



97.5% of values are
below z_1

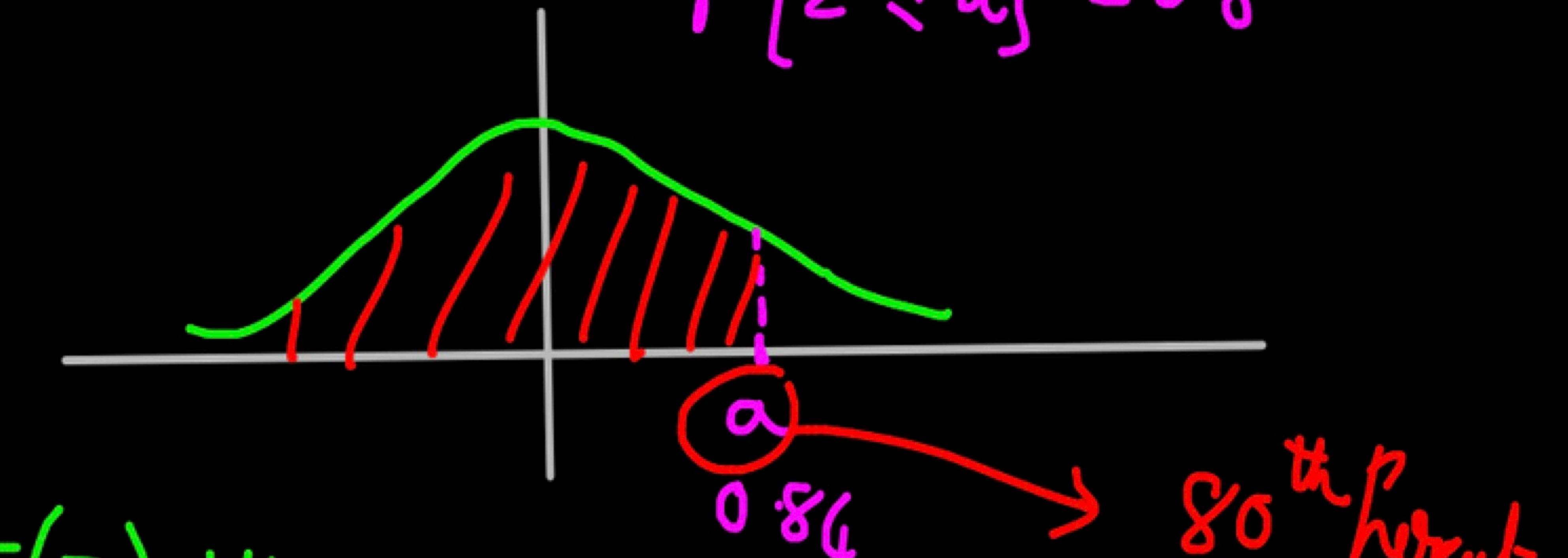
IQ Scores

$$\mu = 100, \sigma = 15$$

$$\text{Var} \rightarrow 15^2$$

what should a person IQ be such that he is in the top 20%.

$$X \sim N(\mu=100, \sigma=15)$$



$$Z = \frac{X - 100}{15} \rightarrow X = 15(Z) + 100$$

How do we find a ? $\text{norm_cdf}(a) = 0.8 \rightarrow \text{norm_ppf}(0.8) = a$

$$a = 0.84$$

$$15(0.84) + 100 = 112.64 \rightarrow \begin{matrix} \text{percentile} \\ \text{cut-off} \end{matrix}$$

80th percentile has two meanings

Distribution

norm.ppf

Sample

Use the dataset
np.percentile

95%). Confidence

→ 2.5 percentile "Sample"

→ 97.5 percentile

$$\text{norm.} \cdot \text{cuff}(0.975) \\ = 1.96$$

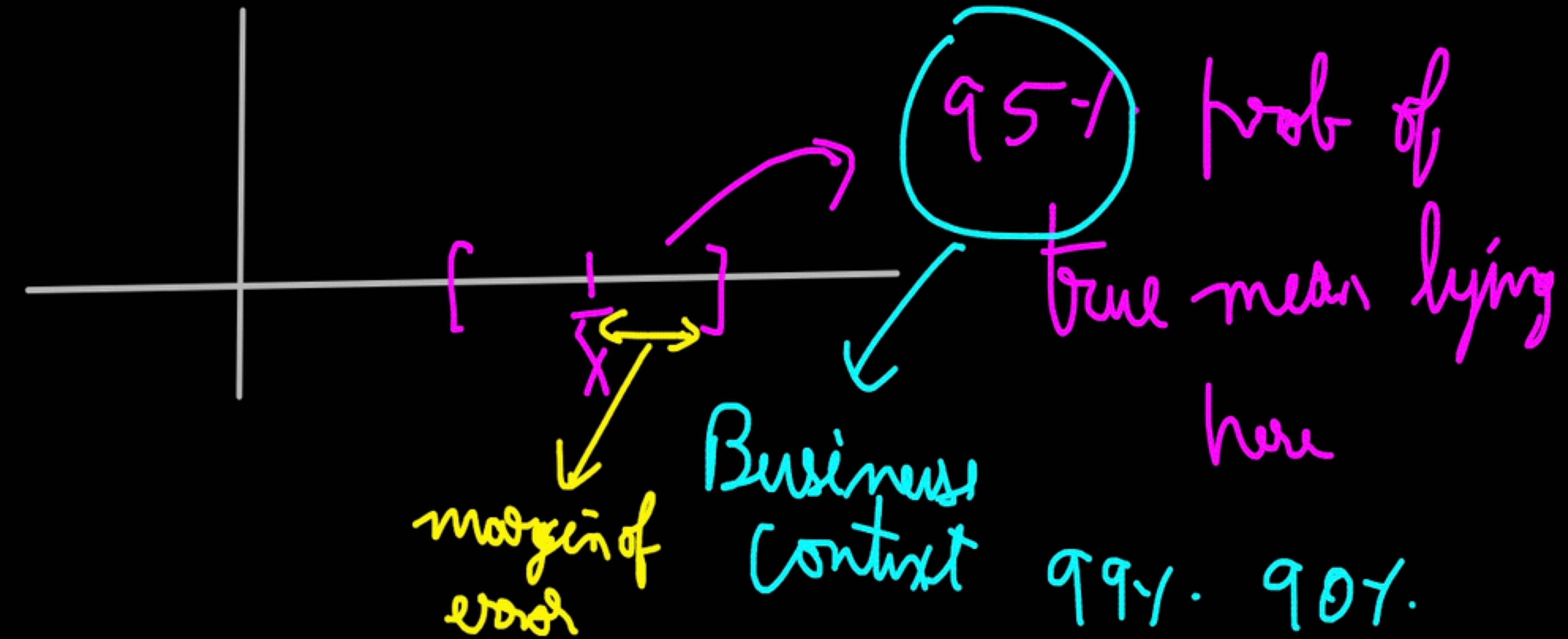
$n = 10000$

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$\rightarrow n \text{ large} \sim \text{Gaussian}$
(CLT)

Sample mean

How to compute Conf. interval? for true mean?



10,000 samples : "row"

97.5
2.5

percentile } of the sample
percentile } mean

→ "Bootstrap"

Use the 10,000 samples → sample with replacement

One iteration

- 1) pick one row randomly → 17th row
- 2) pick one row randomly → 585th row

in a loop $\rightarrow k \leq 10,000$

some will repeat

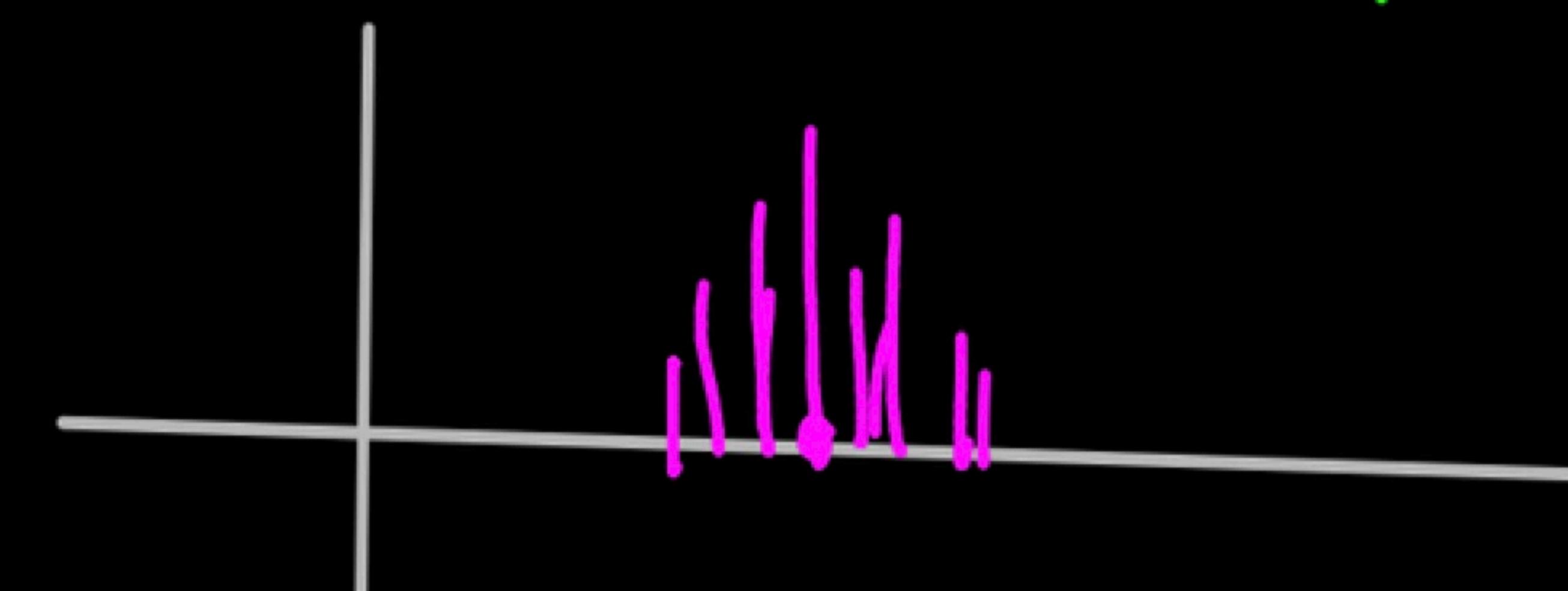
Every sample \rightarrow "new" dataset

| 1st set of sample

2nd set of sample

.

100th



$$\bar{x}_1 \quad \bar{x}_2 \quad \bar{x}_{100}$$

histogram

looks like a
gaussian

Dataset 1

$\mu = 111$

for \rightarrow Sample with replacement
sample → bootstrapped means $\rightarrow \{\bar{x}_1, \bar{x}_2\}$
less variance

Sample with replacement → "smaller CI"

Dataset 2

$\mu = 100$

$\{102, 97, 111, 98, \dots\}$

$\{\bar{x}_1, \bar{x}_2\}$ more variance \rightarrow "larger CI"

$$\{1, 2, 3, 4, 5, 6\}$$

→ sample with replace 6 times

chosen
[3, 4, 3, 5, 1, 2]

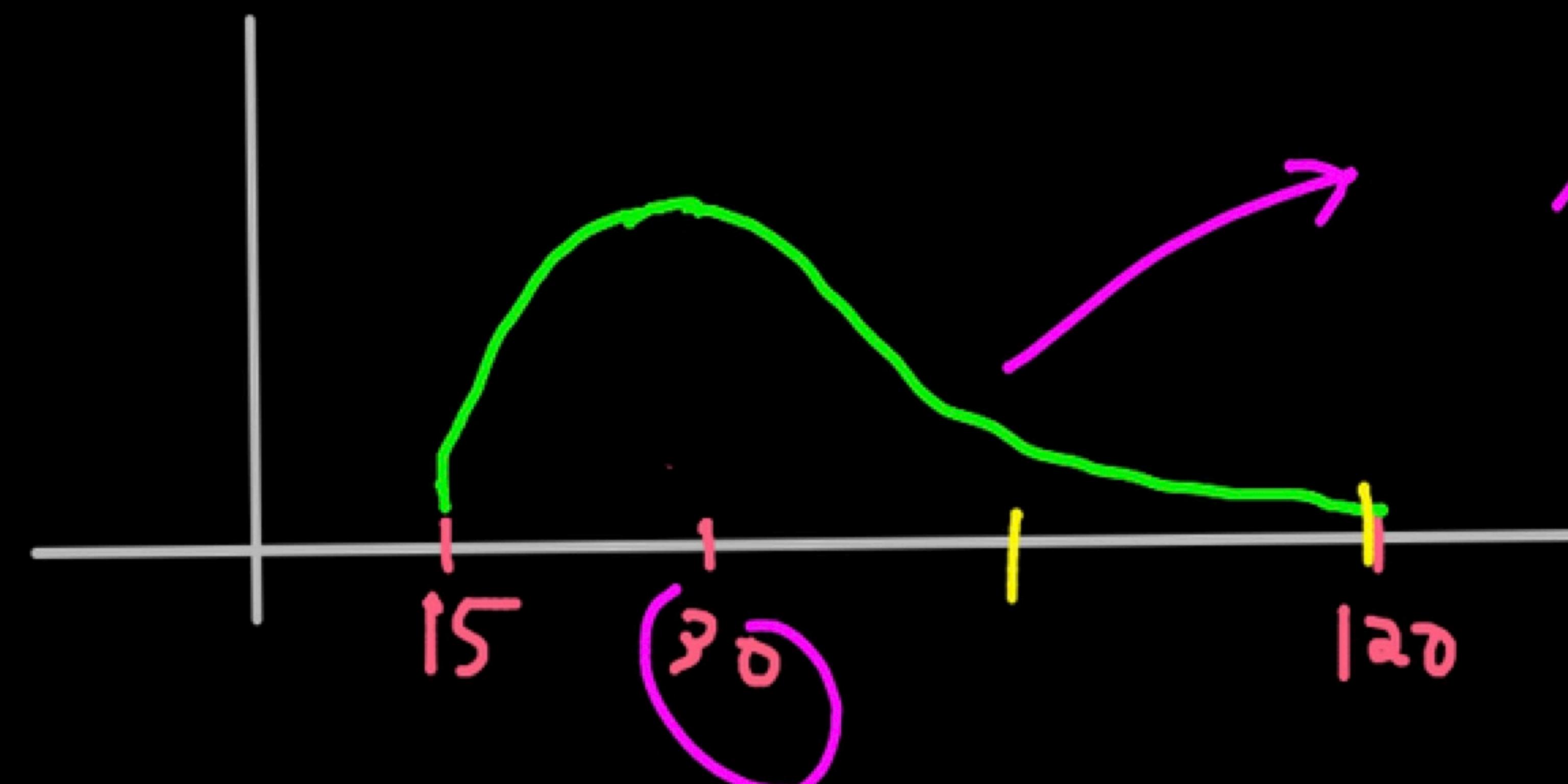
→ sample without replacement

chosen
[3, 1, 4, 2, 5, 6]
same sample now
everytime

Swiggy delivery

30 min → 100 mixy
skewed

$X_1, X_2, X_3 \rightarrow$ data



$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

Closer to μ

$$\mu \in [27, 33] \rightarrow 95\%-CI$$

↳ Publicize

- μ) Conf. Int for max
2) Conf. Int for 90th percentile
 $[45 - 55]$ CI 90th perc.

[45 - 55] → 90th percentile

Only 10% of times will you have to wait
more than 55 mins.

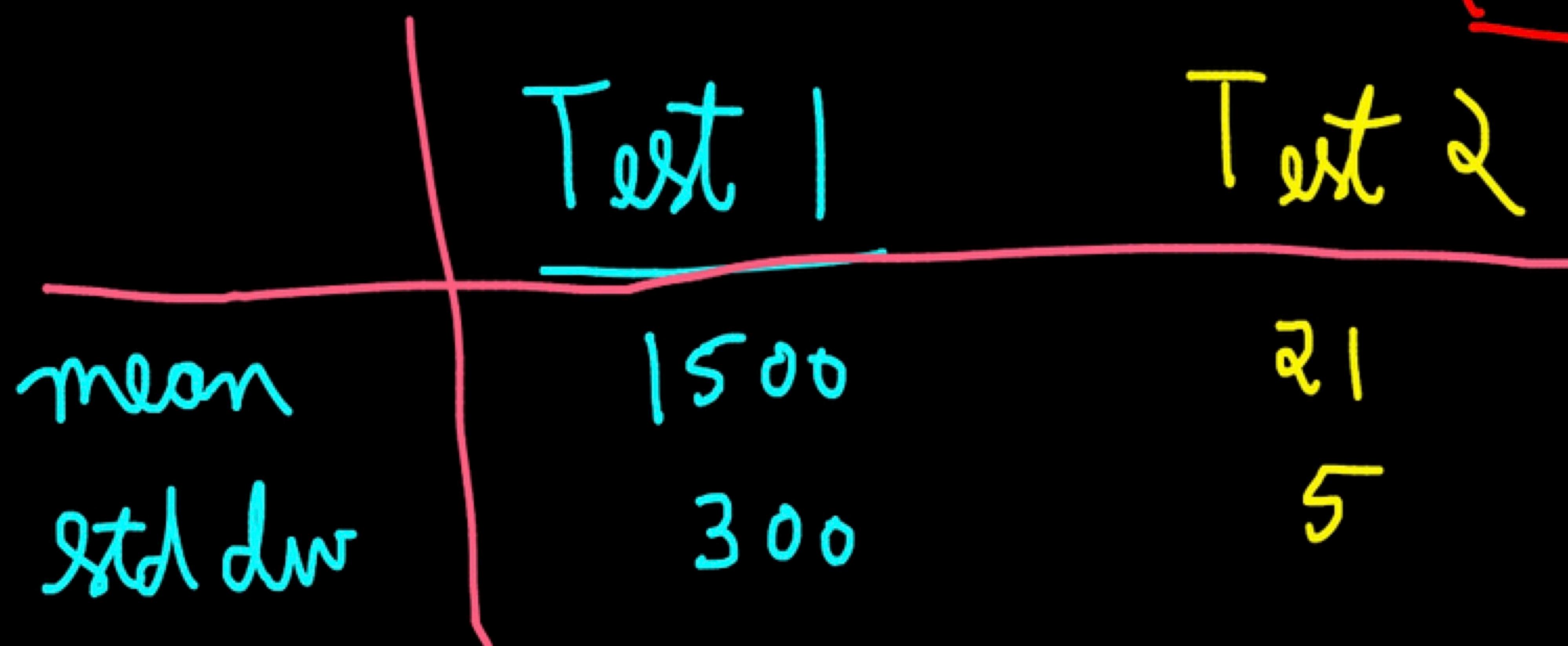
Earlier [] → mean

Now [] → 90th percentile

median

Two types of Exams - Range of marks

different range

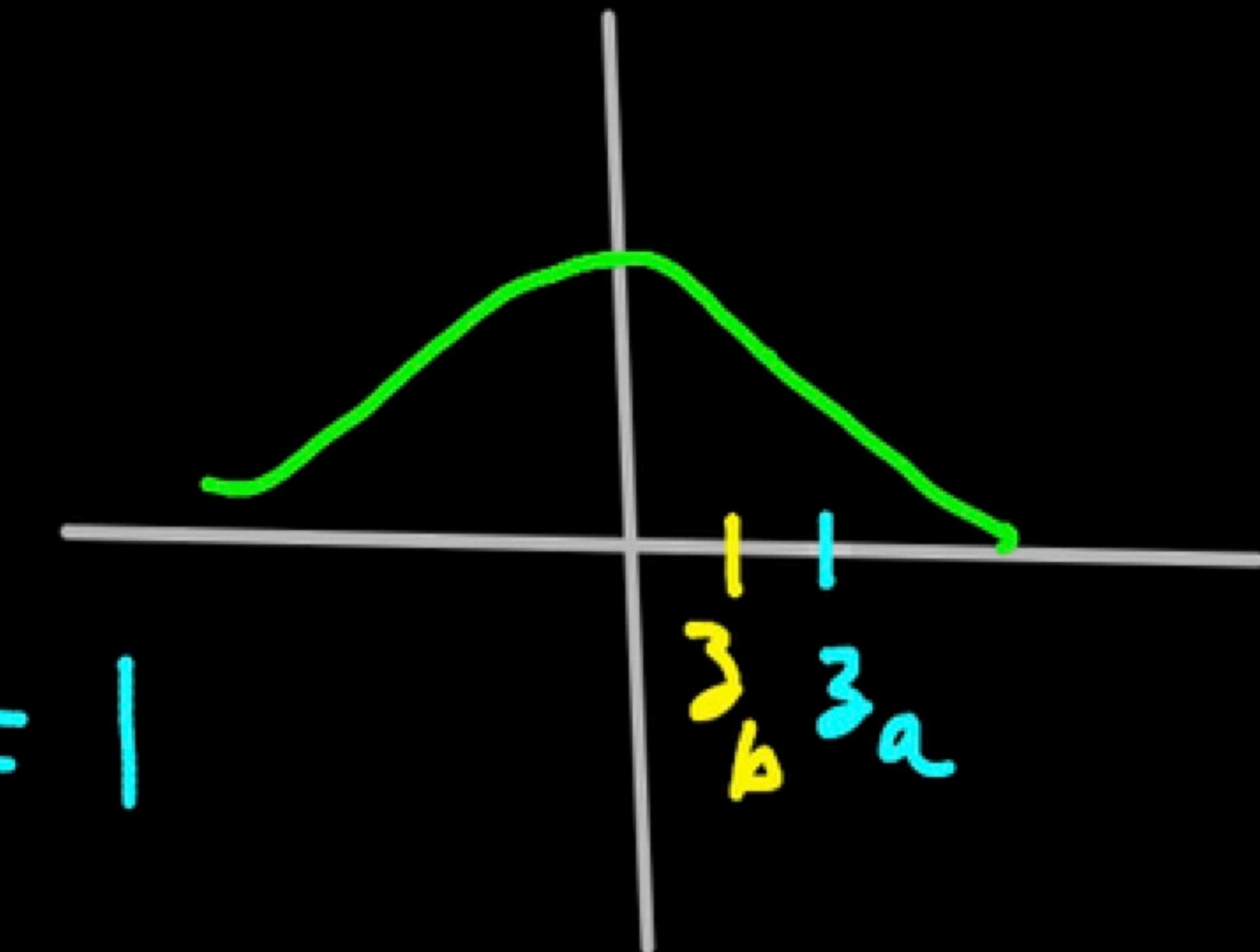


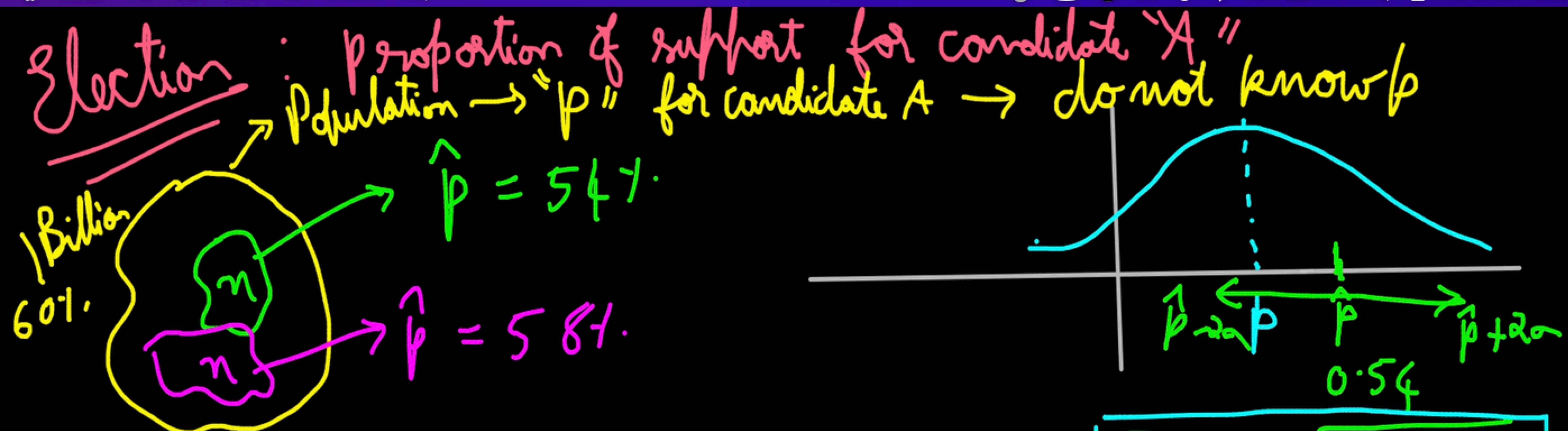
$$A: \text{test 1} \rightarrow 1800 \rightarrow z_a = \frac{1800 - 1500}{300} = 1$$

$$B: \text{test 2} \rightarrow 24 \rightarrow z_b = \frac{24 - 21}{5} = \frac{3}{5}$$

z -score

$$z = \frac{x - \mu}{\sigma}$$





Standard Error:

$$\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$[\hat{p} - 2SE, \hat{p} + 2SE]$$

Eg: $[52, 56]$

$$\hat{p} = \sqrt{\frac{p(1-p)}{n}}$$

↳ Skip the proof

p lies in CI around \hat{p} with 95% Conf.

Hypothesis testing → Tricky

(coin: Cricket captain "7 heads" always

- 1) 10-match series
Won 7 times

not suspicious

Is there
math to
guide this?

- 2) 100-match series
Won 70 times

not so sure?

- 3) 1000-match series
Won 700 times

very suspicious

i) (10-match 7 wins)

I invite 1000 people to toss 10 times

→ How many got ≥ 7 heads
≈ 16%

→ p-value : $\frac{m}{1000}$

→ ~160 out of 1000 0.16 16%

2) (100 matches \rightarrow 70 heads)

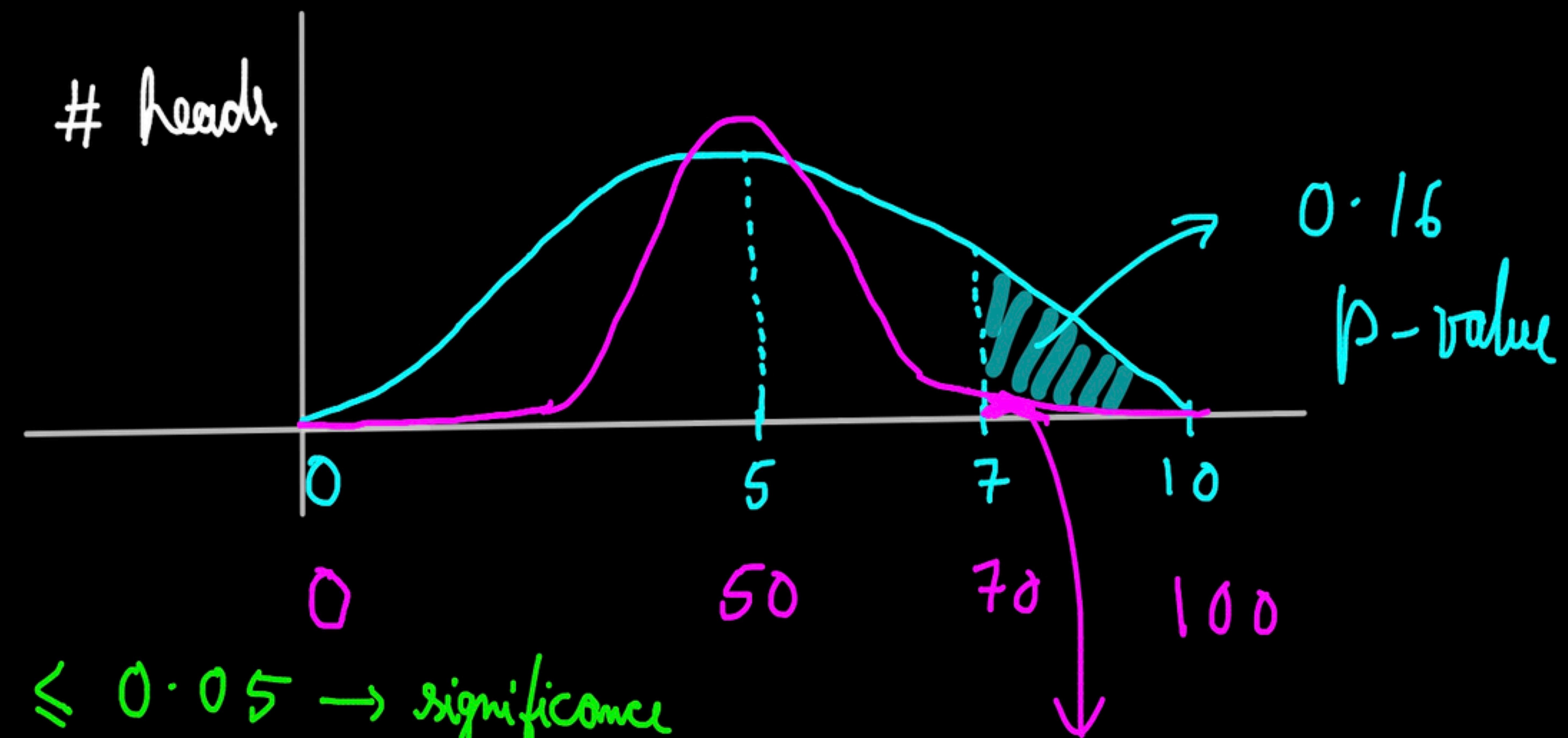
1000 \rightarrow tails 100 times

\rightarrow # people \geq 70 heads $\rightarrow n$

$$\left[\frac{n}{1000} \right]$$

p-value

"Under Null hypothesis" \rightarrow fair coin



strong evidence to
reject the H_0

very very small
p-value

"Coin is fair"

→ Null Hypothesis (H_0)

"Coin is biased towards Head"

→ alternate
hypothesis

"Reject" null hypothesis?

We reject only if we have ^{very} strong evidence (H_a)

Court Judge: Crime has happened. There is a suspect.

What is the null hypothesis made by the judge regarding the suspect?

H_0 : "Innocent until strongly proven guilty"

Reject only if you have strong evidence

Gicket: 3rd umpire DRS

H_0 : On-field umpire is correct

Reject if we have strong evidence

$$\text{CLT } P\left[\frac{\bar{X} - \left(\text{approx. } \mu\right)}{\sqrt{n}}, \bar{X} + \left(\text{approx. } \sigma\right)\right]$$

sample mean only

median , any percentile \rightarrow bootstrap