

Agenda

- * ~~***~~ Central Limit theorem
- ✓ * Application of CLT on real life dataset → height - weight
- ✓ * Confidence Intervals
 - using CLT ✓
 - using Bootstrapping ✓
 - Assignments

Recap

* Normal / Gaussian distribution

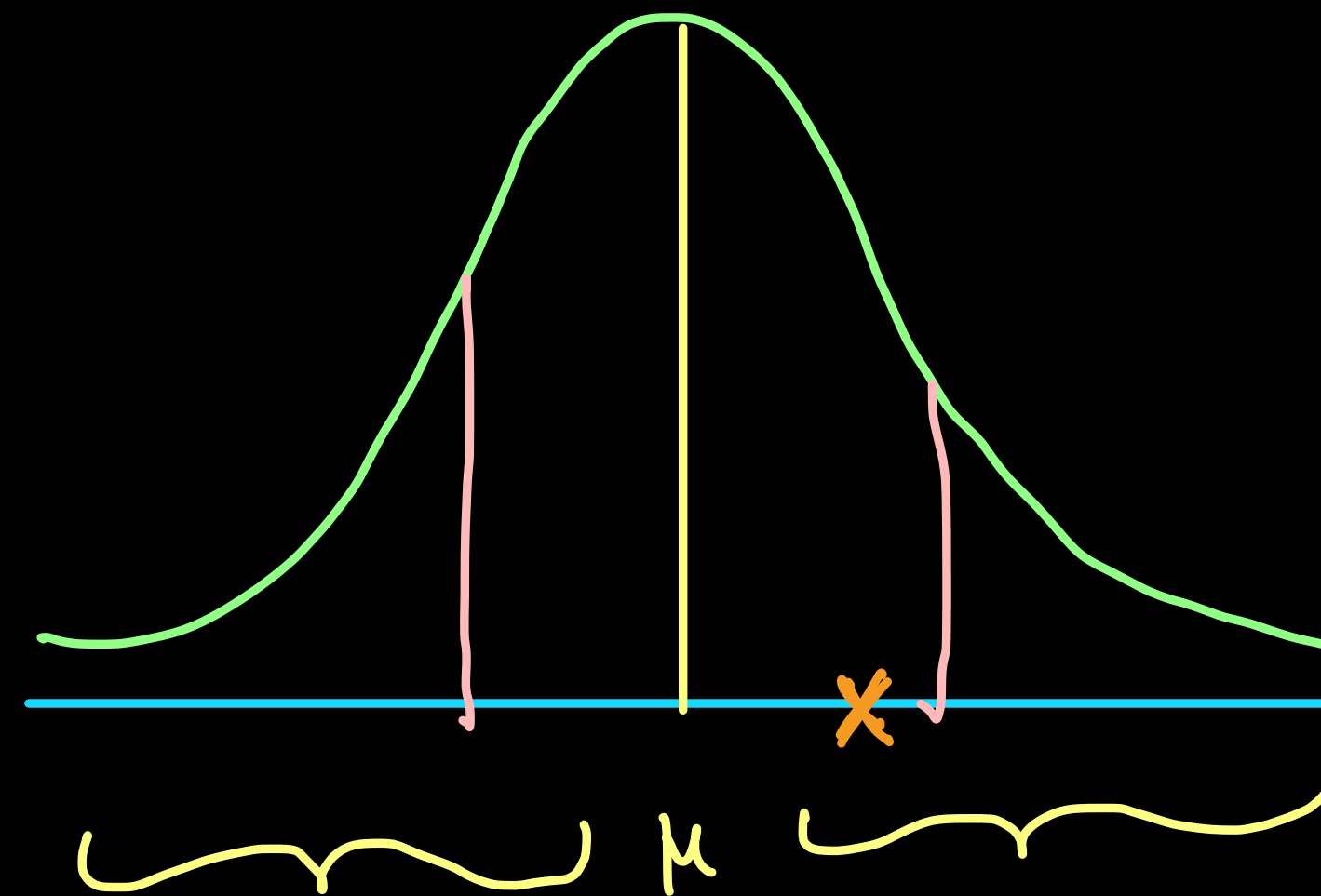
* Empirical rule

$$\left\{ \begin{array}{l} \mu - \sigma \leq \textcircled{x} \leq \mu + \sigma \rightarrow 0.68 \\ \mu - 2\sigma \leq x \leq \mu + 2\sigma \rightarrow 0.95 \\ \mu - 3\sigma \leq x \leq \mu + 3\sigma \rightarrow 0.997 \end{array} \right.$$

z-score

z-table

python \rightarrow scipy



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

*
✓
✓

standard Normal Distribution / z-distribution

The special case, when normal distribution has

$$\mu = 0$$

$$\sigma = 1$$

standard

✓ Example

Imagine you have two friends, ✓

Alex and Taylor, who are both great at math but have different grading systems. Alex's math scores range from 0 to 100, while Taylor's scores range from 0 to 50. You want to know who is performing with more consistency

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

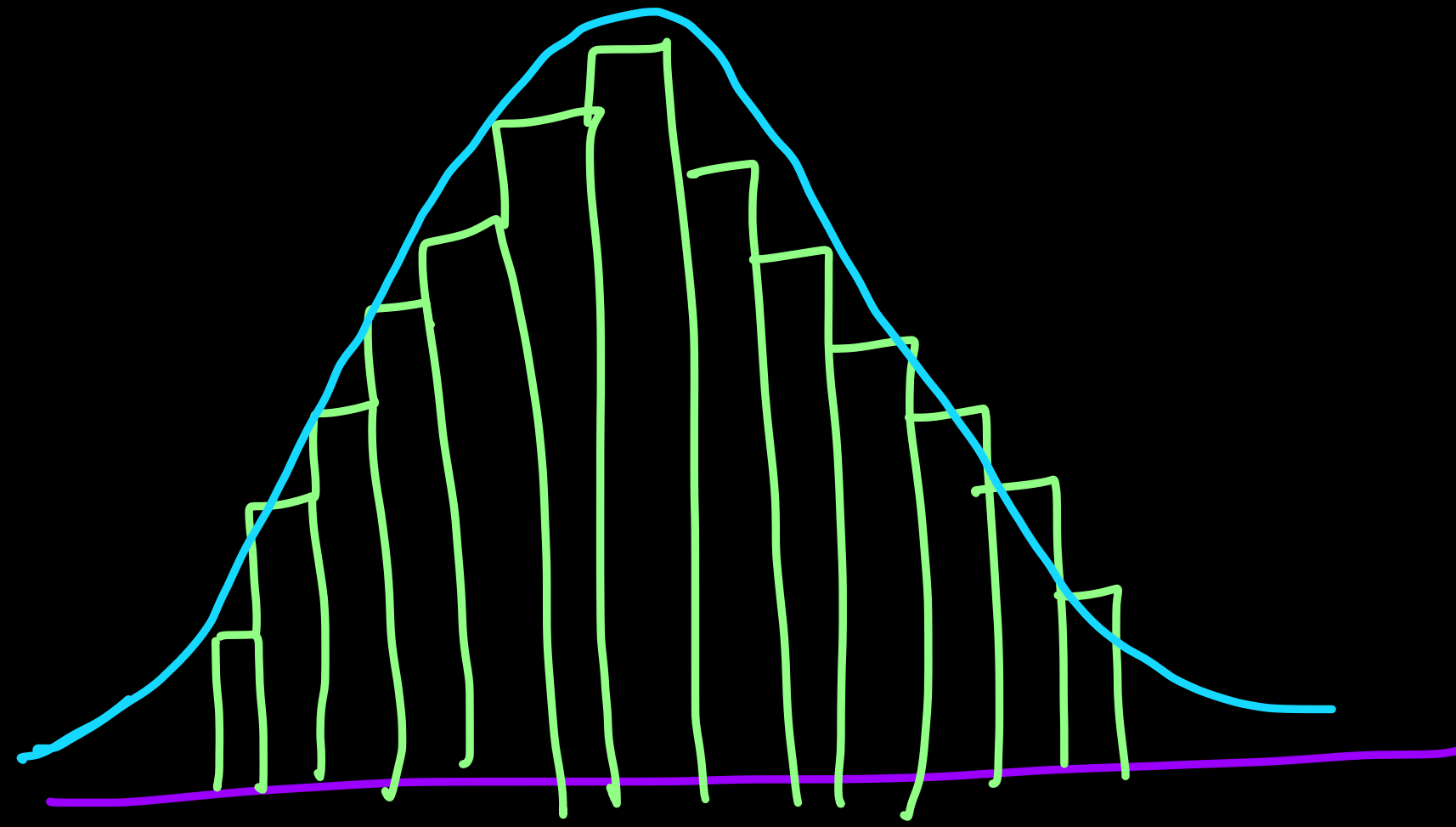
Central Limit Theorem



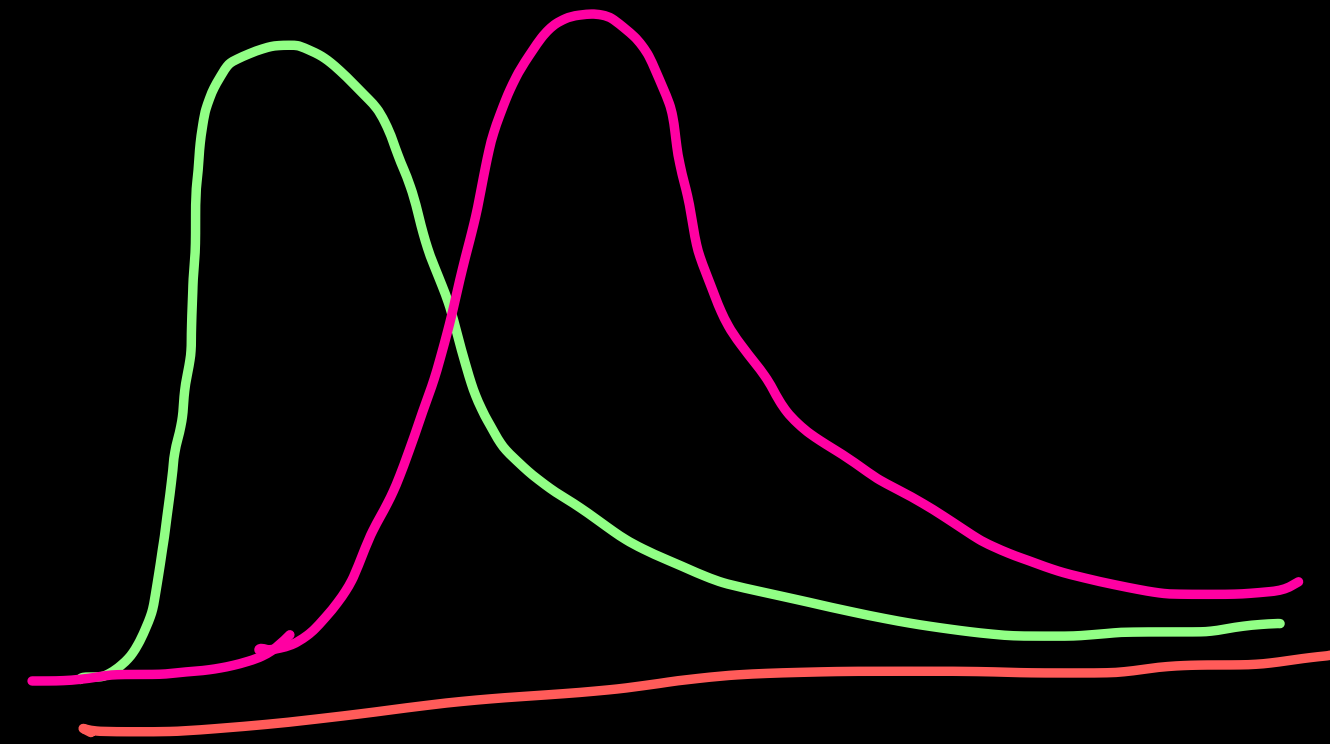
Q) what is the avg weight of the JB Basket?

$$\begin{aligned}
 s_1 &= [\underline{w_1}, w_2, \dots, w_{30}] = \bar{x}_1 \checkmark \\
 s_2 &= [w_1, w_2, \dots, w_{30}] = \bar{x}_2 \checkmark \\
 &\vdots \\
 s_{1000} &= [w_1, w_2, \dots, w_{30}] = \bar{x}_{1000} \checkmark
 \end{aligned}$$

mean of sample means \rightarrow $\frac{\bar{x}_1 + \bar{x}_2 + \dots + \bar{x}_{1000}}{1000} \approx$ population mean



\Rightarrow { Sample means always
✓ follows the Normal Distribution



Central Limit Theorem

CLT states that " the ^{mean} mean of a random sample will resemble even closer to the population mean as the sample size increases and it will approximate a normal distribution regardless of the shape of the population distribution ✓ ✓

mean of Sample means \approx population mean

$$n = 4$$

Q) What is the right sample size to take?

jelly beans samples taken

4 times

40 times

✓
trails

$n=12$

$n=40$

Sample
Size

→ ✓ Larger sample size → A sampling distribution that closely resembles a normal distribution

CLT tends to work well when n is sufficiently large. Typically consider as $n \geq 30$. ✓

sample size



3

>

66.
=

3

<

66.

= 5 \checkmark \nwarrow avg

\Rightarrow

66.3
=

> 66

$\mu = 66.36$
 $\sigma = 3.84$
 population $\frac{10000}{=}$
 mean of sample means
 $s_1 = 1.735$
 Sample size $\checkmark = 5$
 $\underline{\underline{66.35}}$

Sample size = 20
 $\bar{x}_{20} = 66.36$
 $s_2 = 0.86$
 $\frac{\sigma}{\sqrt{n}}$

Sample size $\overset{n}{=} \underline{\underline{100}}$
 $= 66.36$
 mean of sample mean
 $s_3 = \underline{\underline{0.38}}$

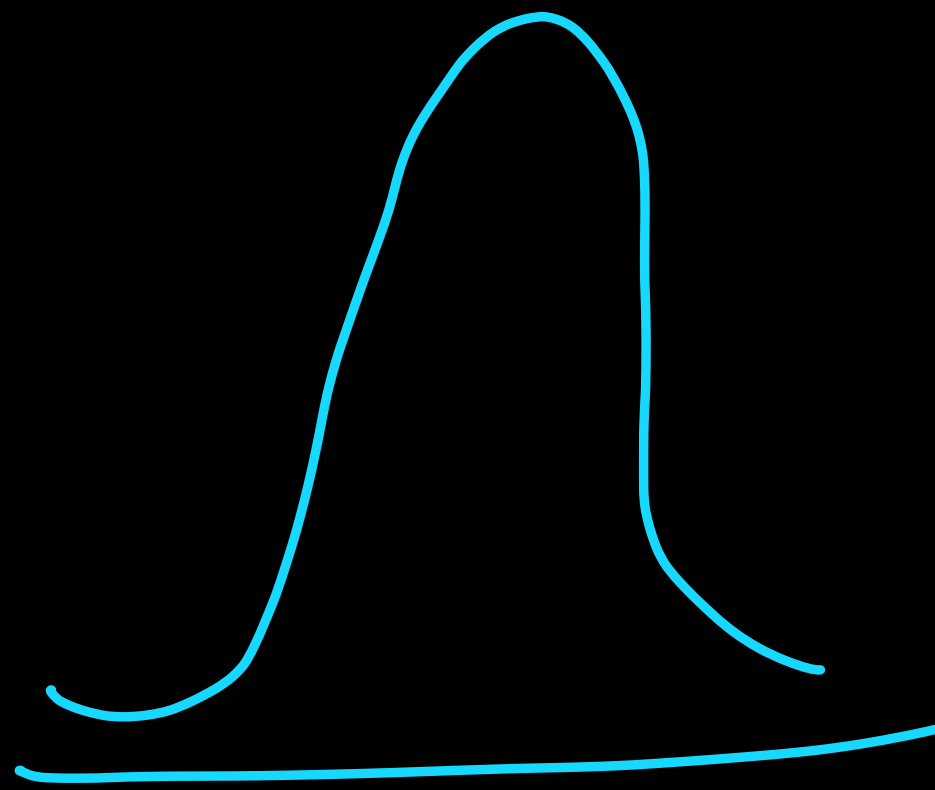
(1) mean of sample mean \approx population mean

(2) $\checkmark \frac{\sigma}{\sqrt{n}} \Rightarrow \frac{3.84}{\sqrt{100}} = \frac{3.84}{10} = \underline{\underline{0.384}}$

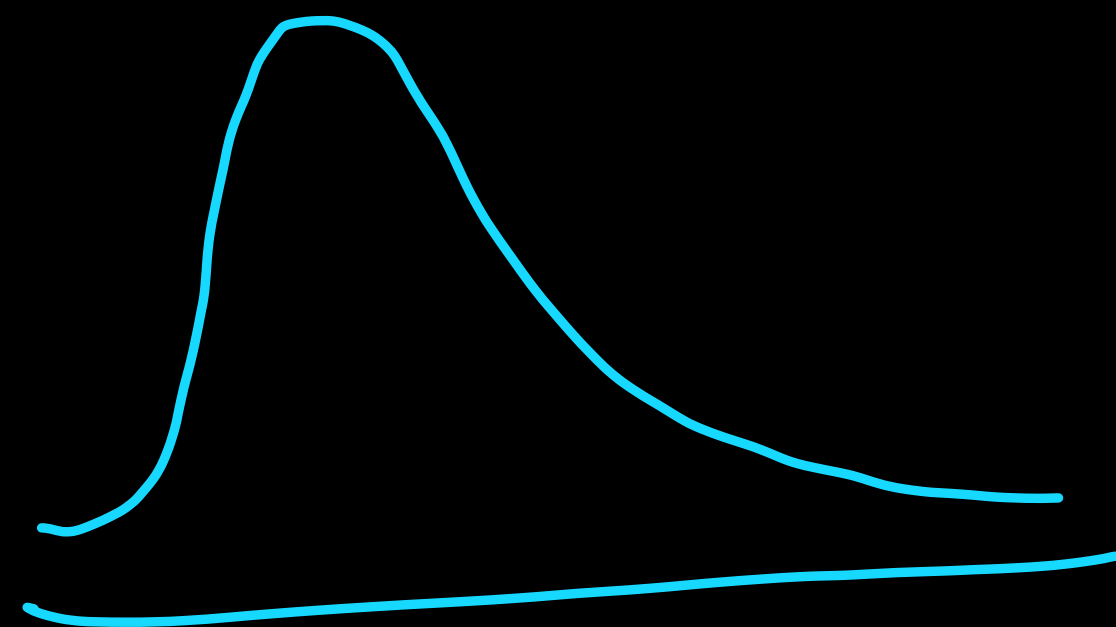
$$\sigma / \sqrt{n}$$



SD of sample means

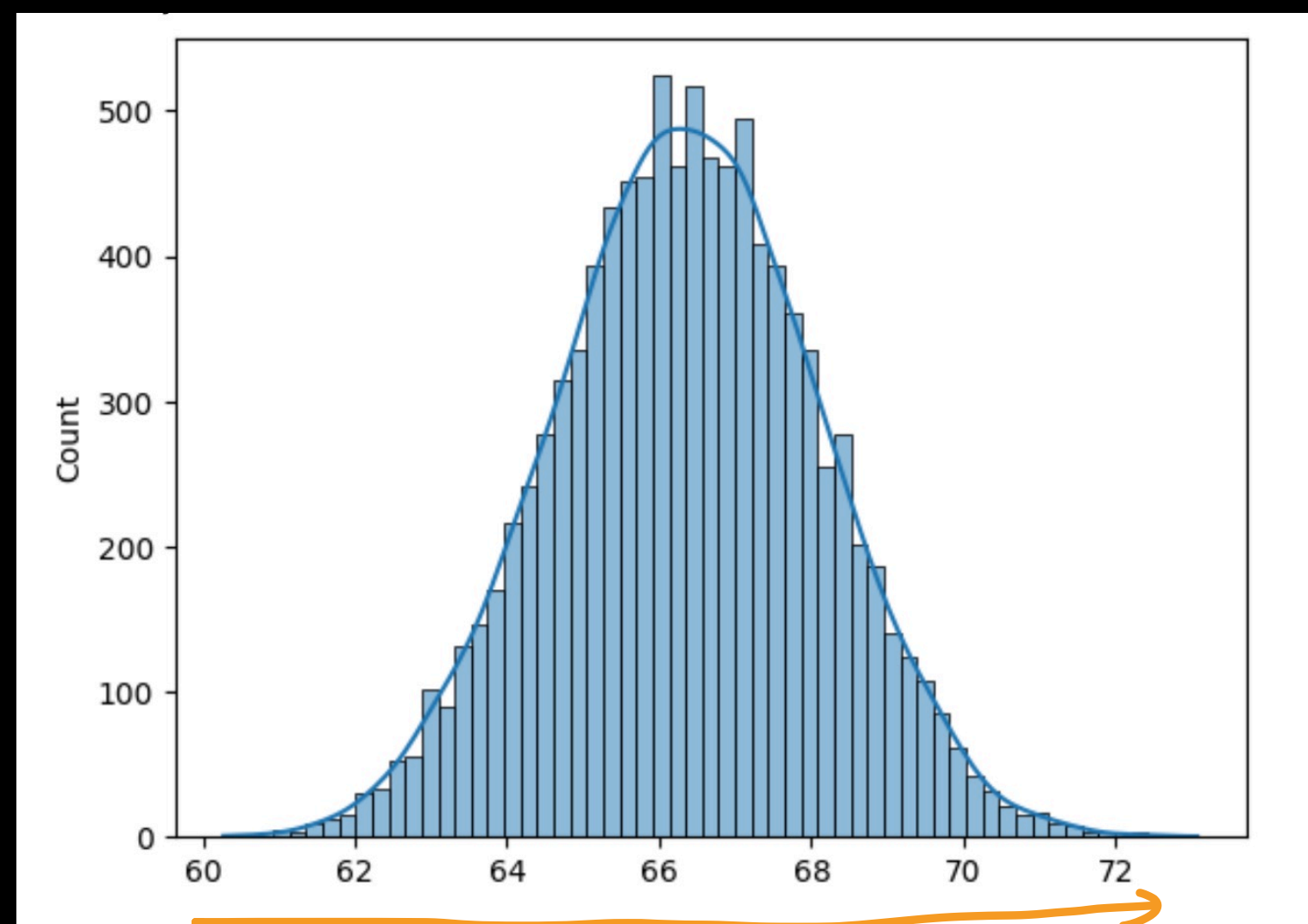


$$\begin{cases} n = 5 \\ n = 20 \end{cases}$$



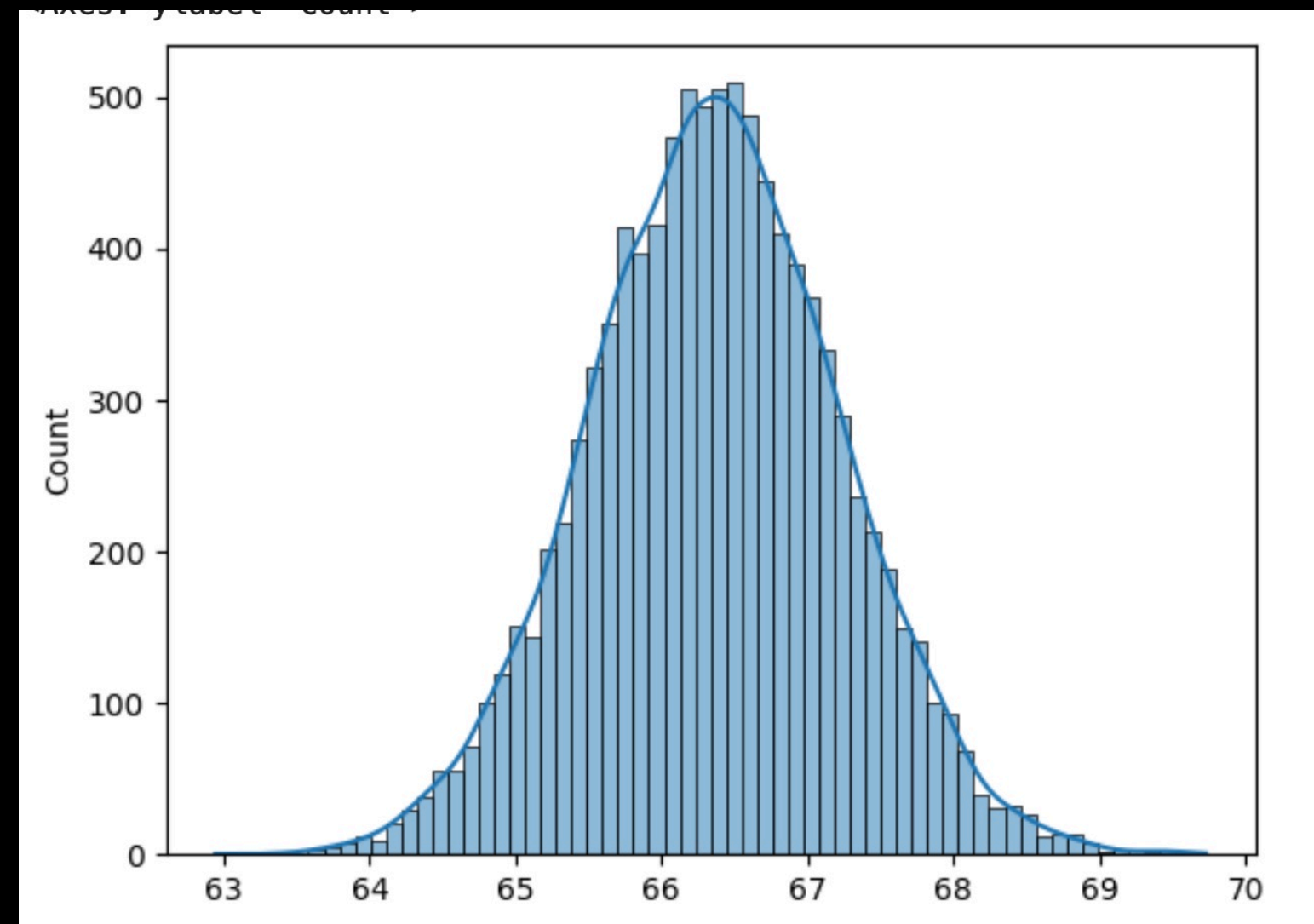
$$n \geq 30$$



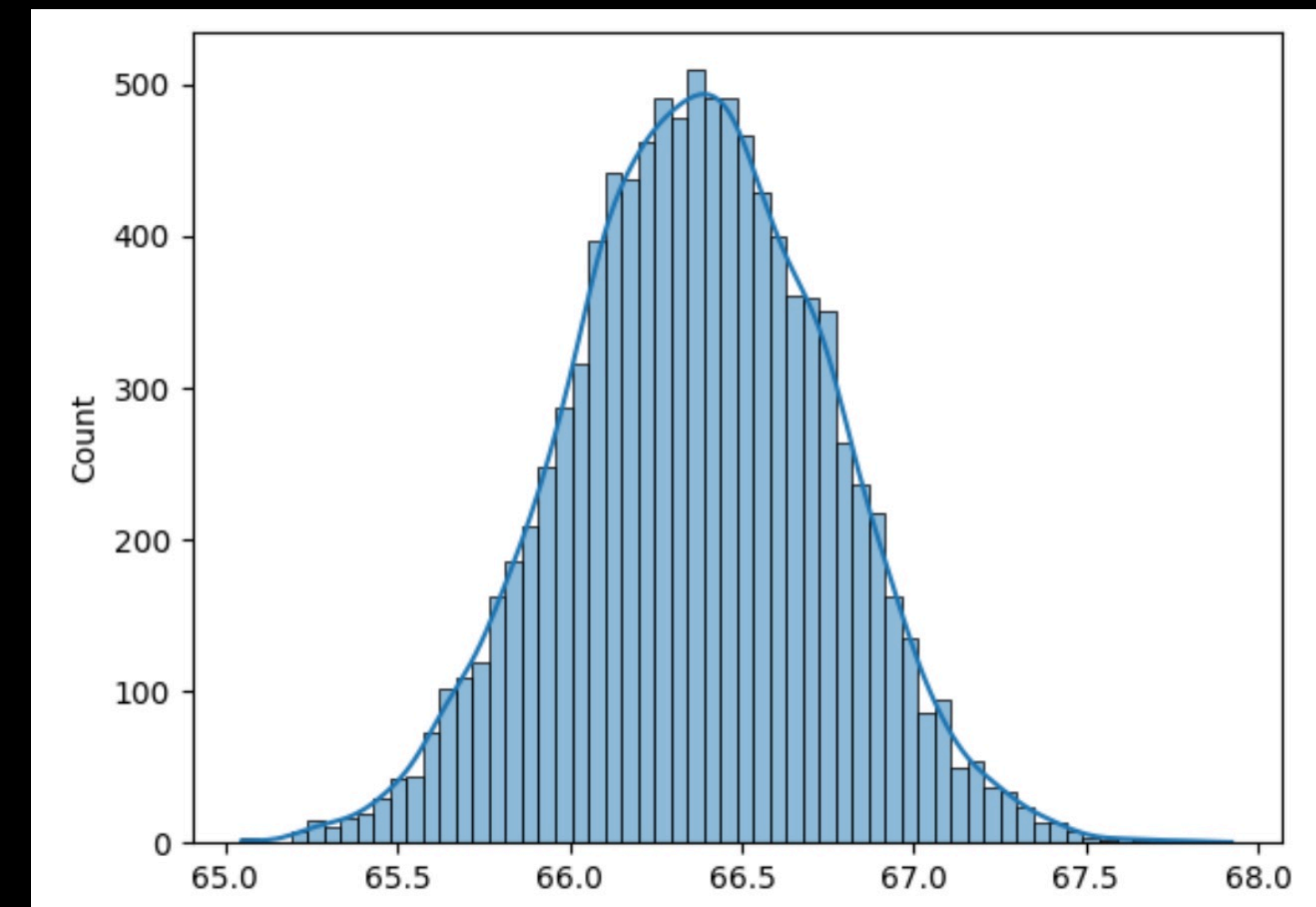


Sample-5

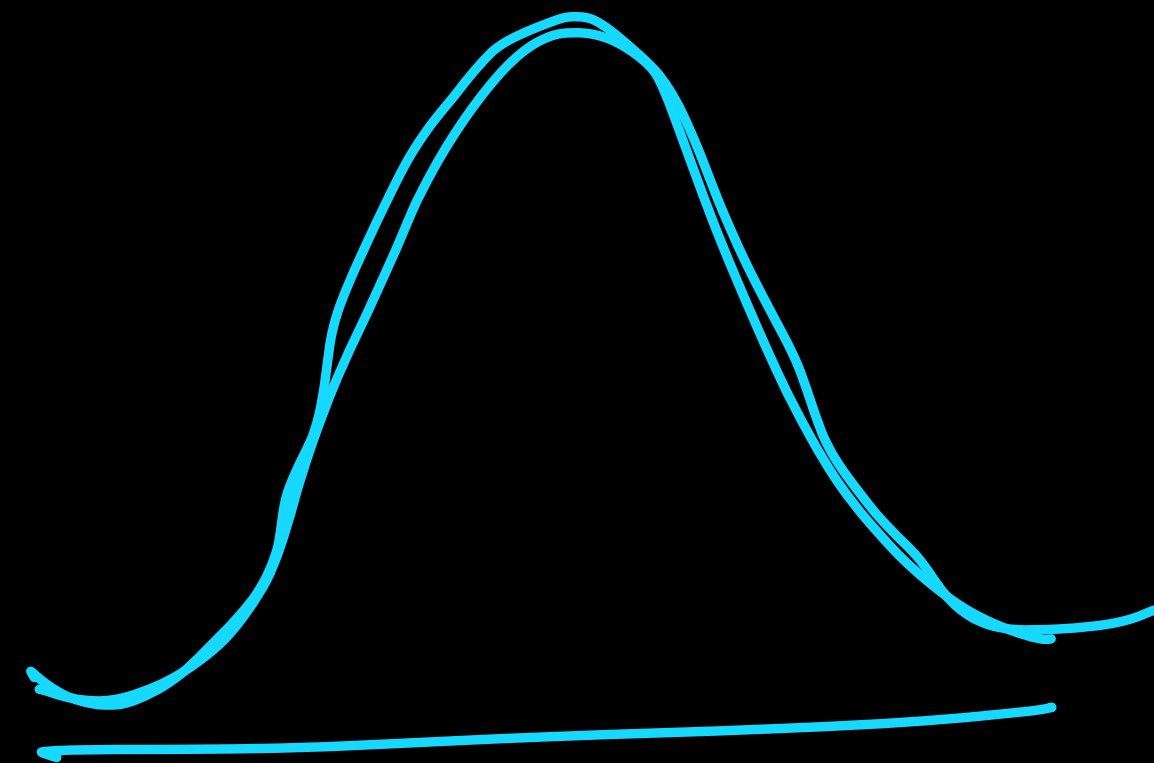
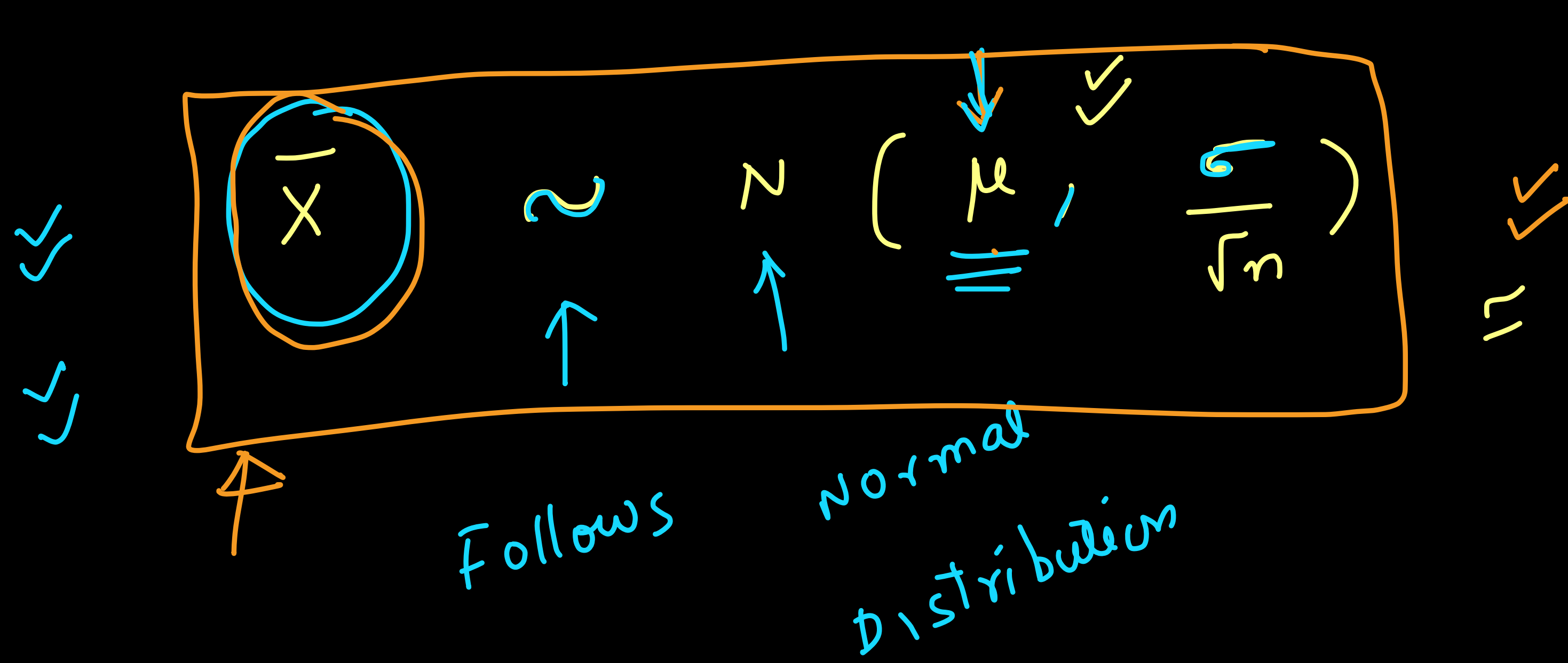
5/5m



20



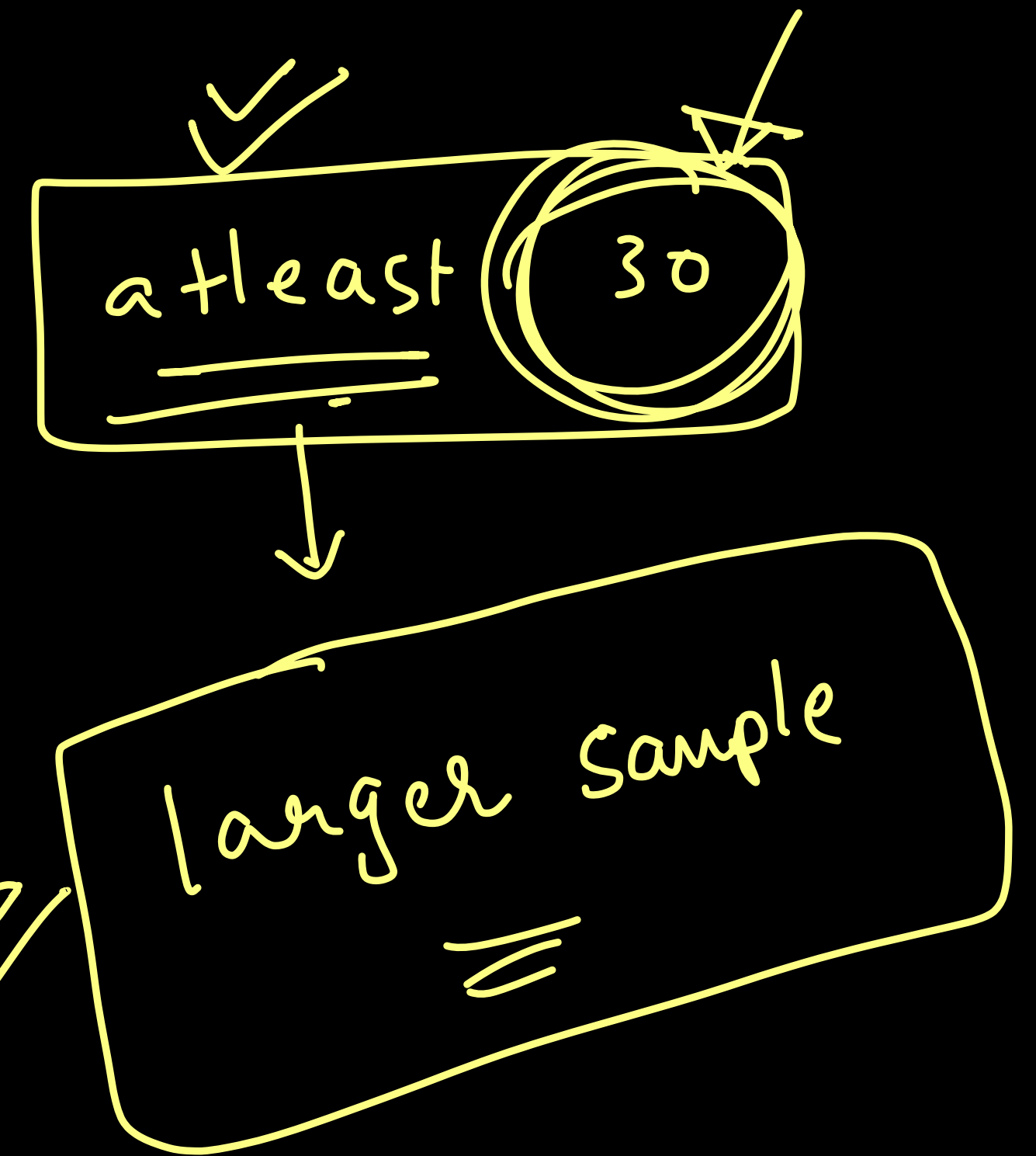
100



$n < 30$

\Rightarrow

more accurate



$n \geq 30$

conditions of the CLT ✓✓✓

To apply the central limit theorem, the following conditions must be met:

1. **Randomization**: Data should be randomly sampled, ensuring every population member has an equal chance of being included.

2. **Independence**: Each sample value should be independent, with one event's occurrence not affecting another. Commonly met in probability sampling methods, which independently select observations.

3. **Large Sample Condition**: A sample size of 30 or more is generally considered "sufficiently large." This threshold can vary slightly based on the population distribution's shape.

10000
2

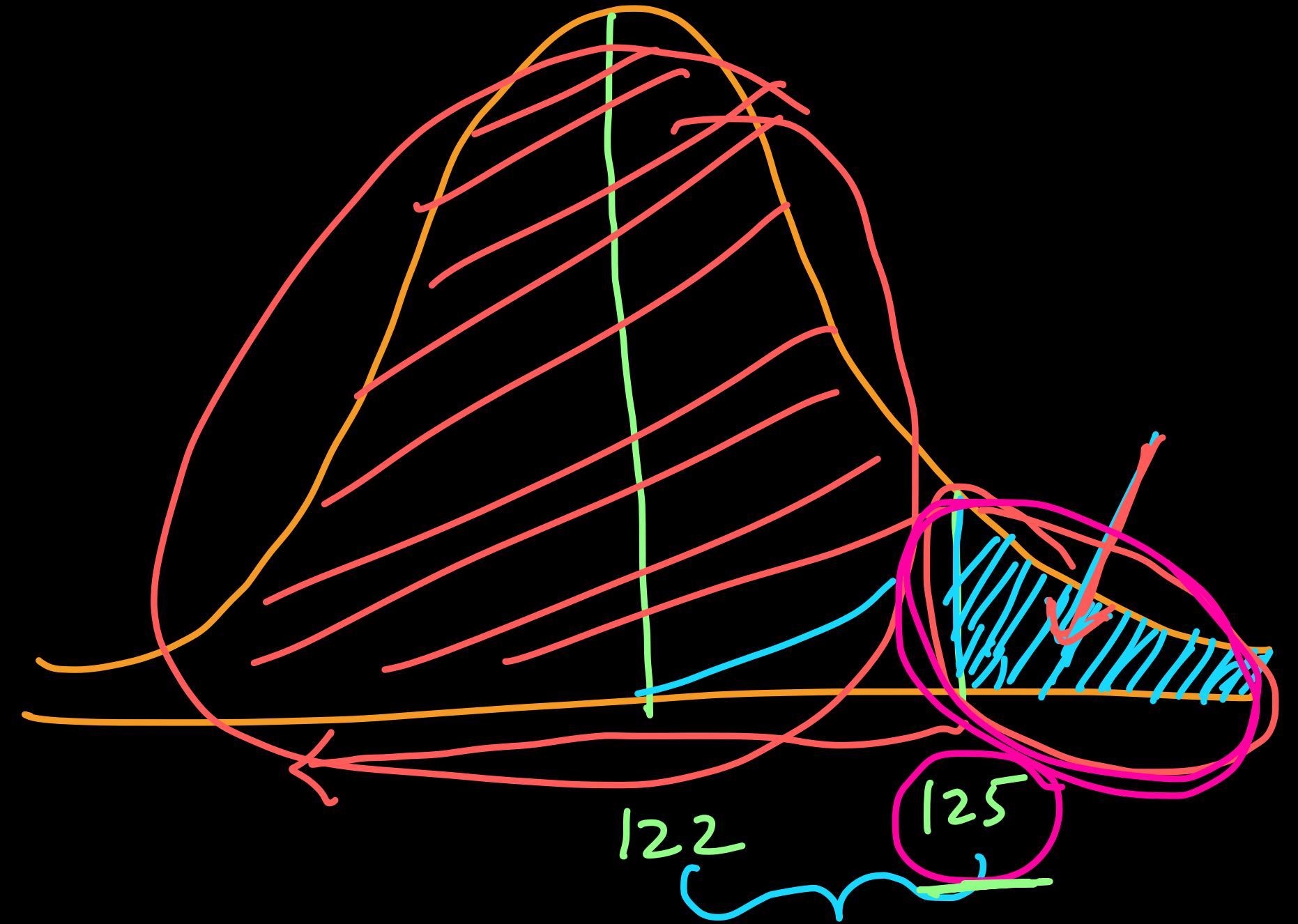
30

problem

Systolic blood Pressure of a group of people is known to have an average of 122 mmHg and standard deviation of 10 mmHg. Calculate the probability that the average blood pressure of 16 people will be greater than 125 mm Hg.

$$\begin{cases} \mu = 122 \text{ mmHg} \\ \sigma = 10 \text{ mmHg} \\ n = 16 \end{cases}$$

$$\begin{aligned} \text{S.D of Sample} &= \sigma / \sqrt{n} \\ \text{S. error} &= 10 / 4 = 2.5 \end{aligned}$$



$$Z = \frac{125 - 122}{2.5} = \frac{3}{2.5} = 1.2$$

In an ecommerce website, the average purchase amount per customer is \$80 with a standard deviation of \$15.

If we randomly select a sample of 50 customers, what is the probability that the average purchase amount in the sample will be less than \$75?

$$\mu = \$80$$

$$\sigma = \$15$$

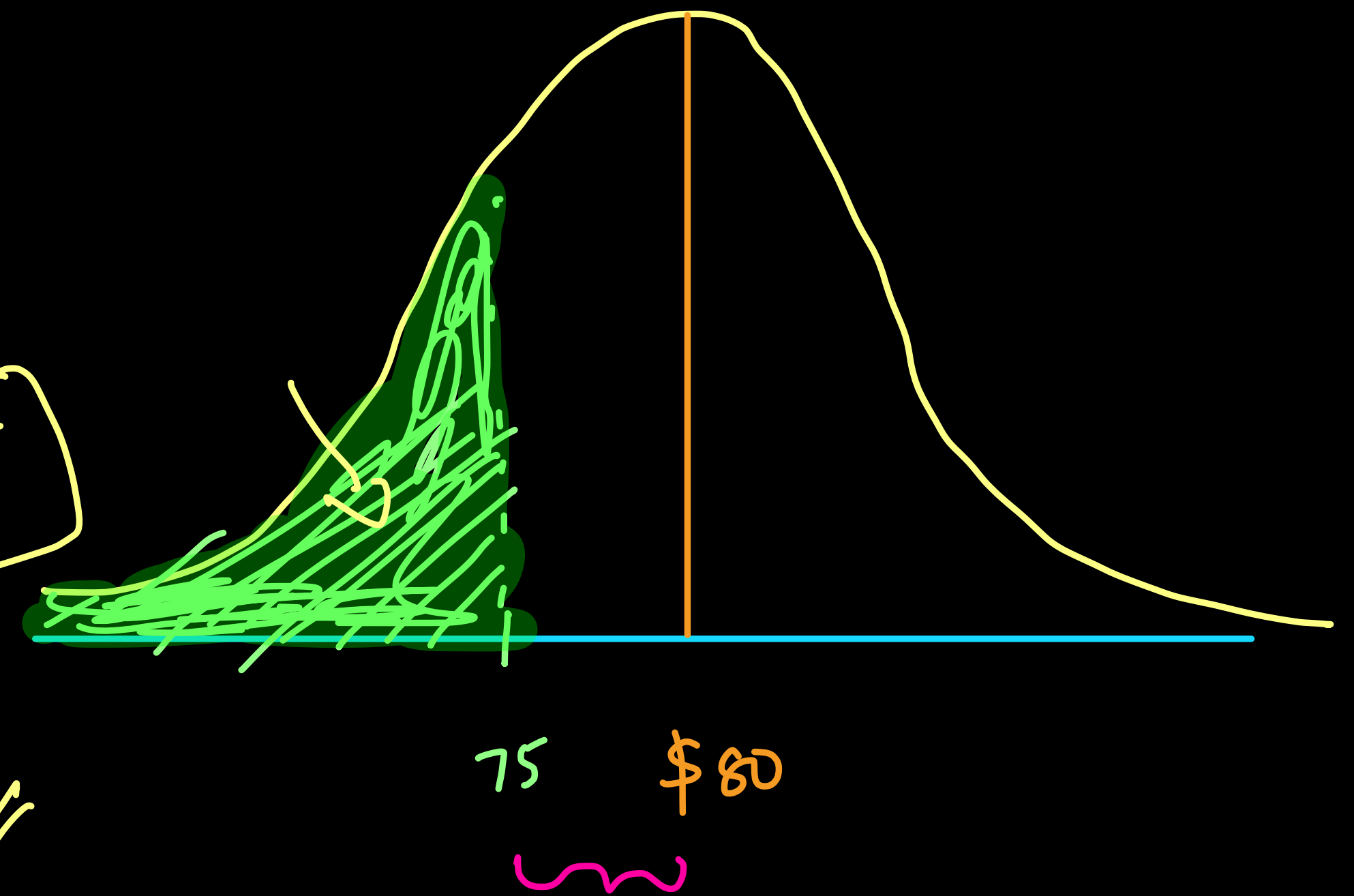
$$n = 50$$

$$z = \frac{75 - 80}{\left(\frac{15}{\sqrt{50}}\right)}$$

$$z = -2.35$$

$$S.D = \frac{\sigma}{\sqrt{n}} = \frac{15}{\sqrt{50}}$$

$$P\{x < 75\} = 0.009 \checkmark \checkmark \leftarrow \text{norm.cdf}(-2.35) \checkmark \checkmark$$



Weekly toothpaste sales have a mean of 1000 and std dev of 200. What is the probability that the avg weekly sales next month is more than 1110?

population
std dev $\leftarrow \sigma = 200$
✓ $n = 4$

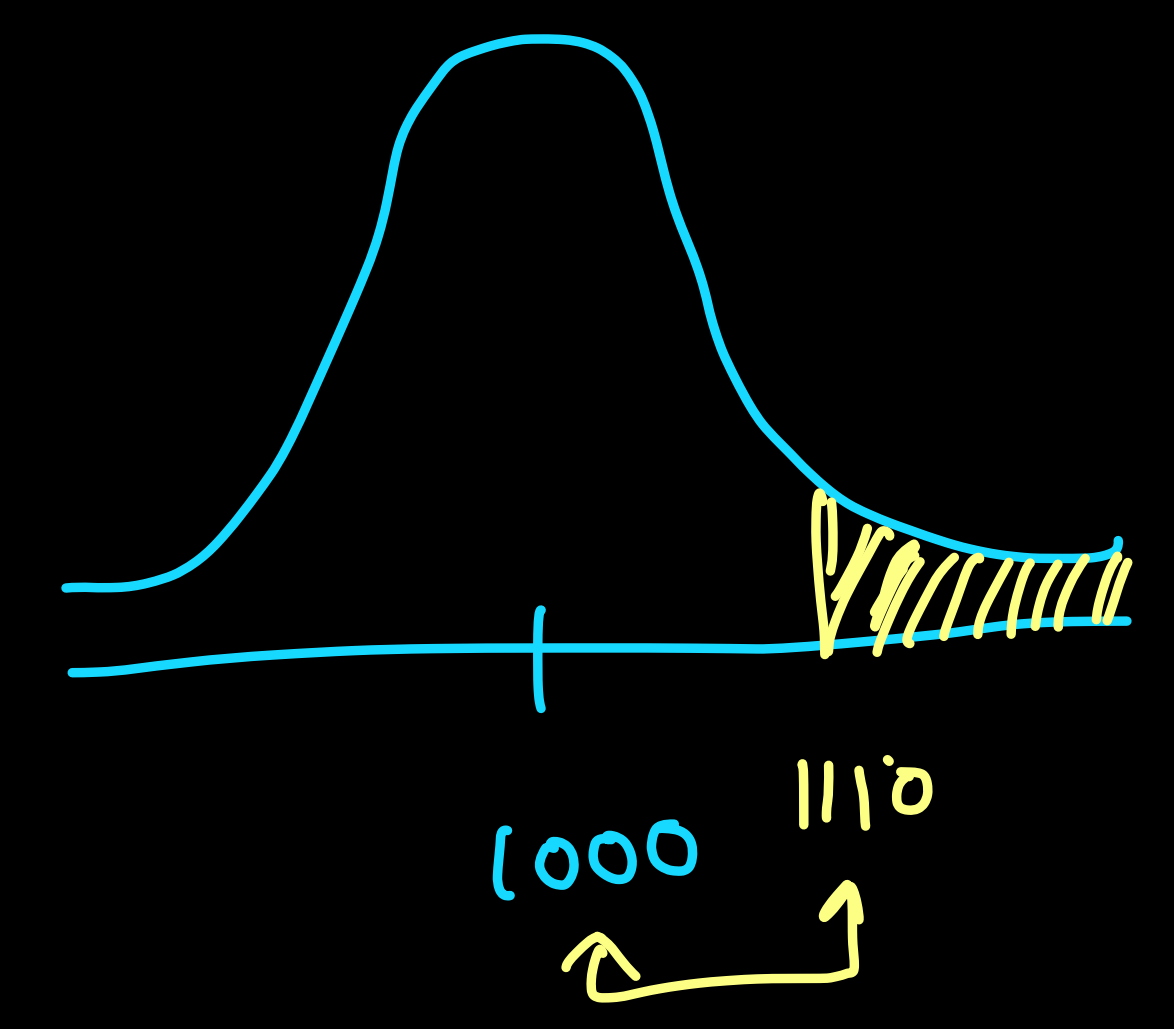
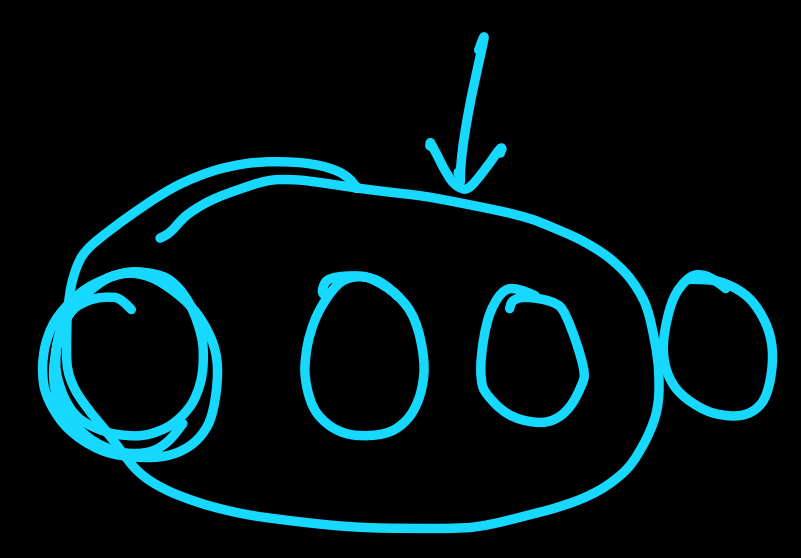
$$\sigma / \sqrt{n}$$

S.D of
sample

$$= \frac{\sigma}{\sqrt{n}} = \frac{200}{2} = 100$$

$$z = \frac{1110 - 1000}{100} = \frac{110}{100} = 1.1$$

✓
✓
1 - norm cdf(1.1)



Confidence Interval

Exit polls

Range

p_1
 p_2
 p_3

30

0-100

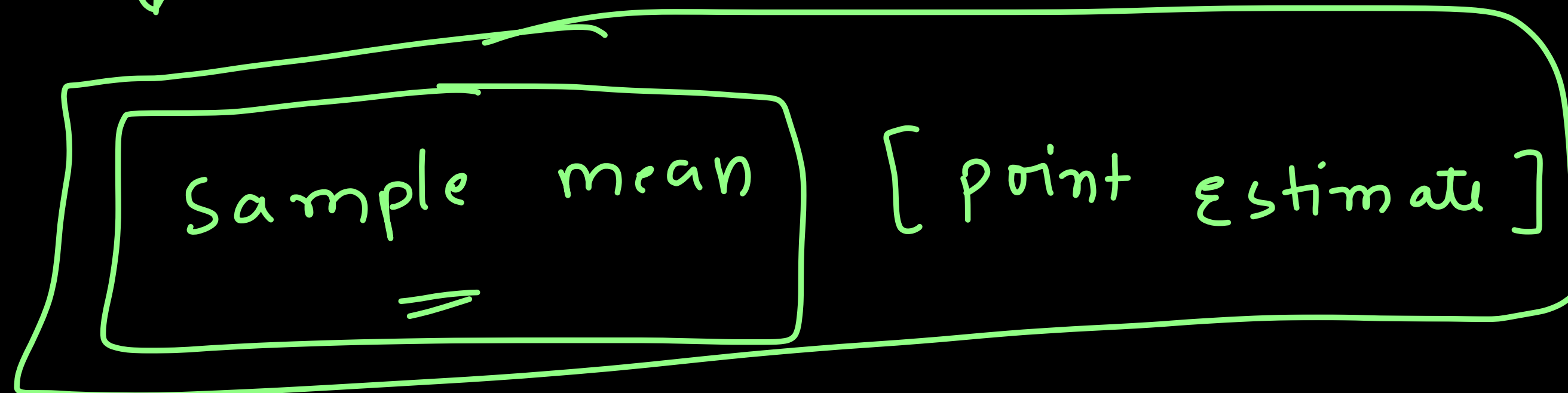
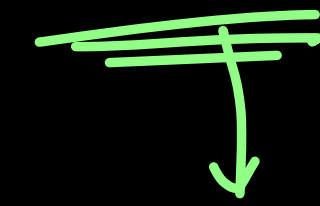
0-100

0-100

25-35

How do they decide the range

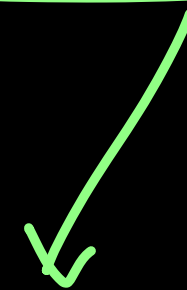
→ point estimate



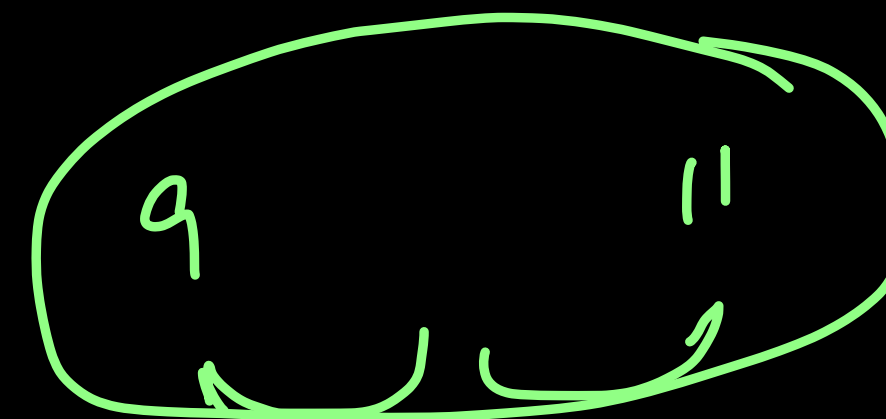
\pm

~~S.D~~

Some error



C.I



95% C.I →

99% C I

→

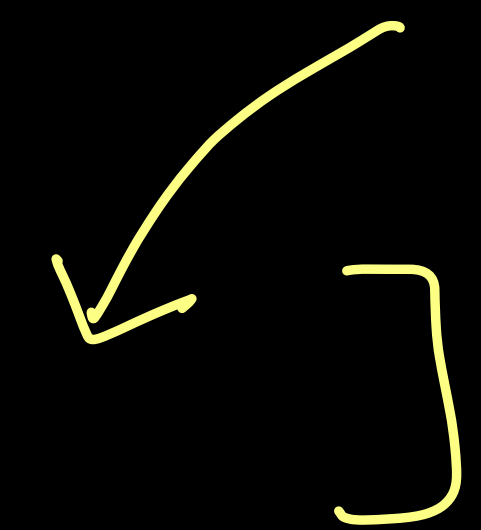
[,]



A

→ 95%

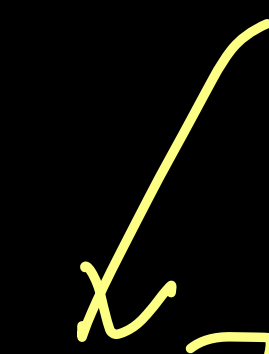
[



→ 90%

→

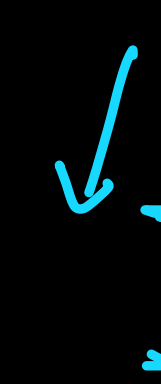
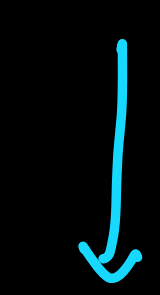
[.]



95%



[]



[64 66]

Q How to calculate C.I ?

95

✓ (1) Using CLT [central limit theorem]

(2) Using Bootstrapping

Example
Height

Mean height of sample of 100 adults = 65 inches

Sample mean

→ [population]

S.D

= 2.5

S.D of sample
 $\frac{6}{5}n$
=

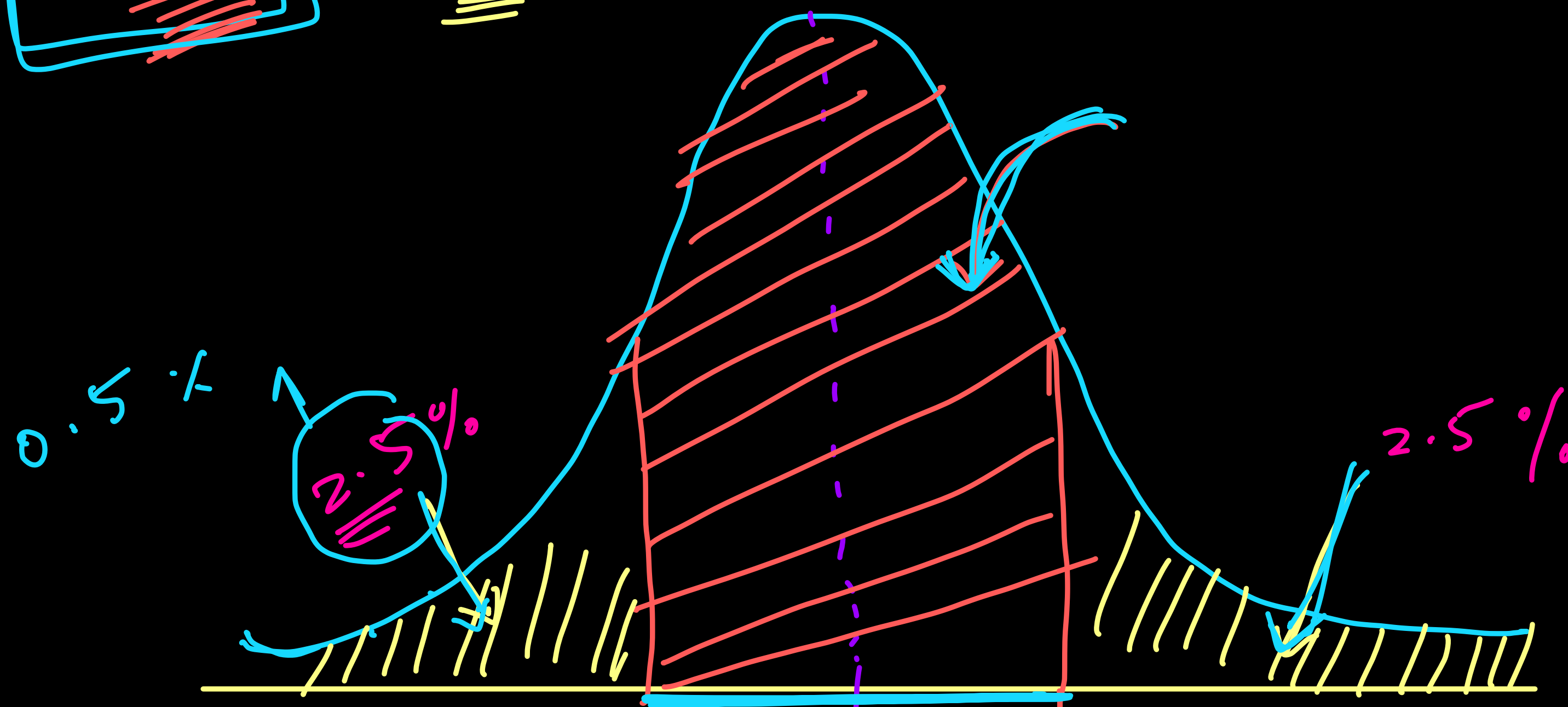
$n = 100$
 $\sigma = 2.5$
 $\bar{x} = 65$

S.D of sample = $\frac{\sigma}{\sqrt{n}}$
 $= \frac{2.5}{10} = 0.25$

95%

C I

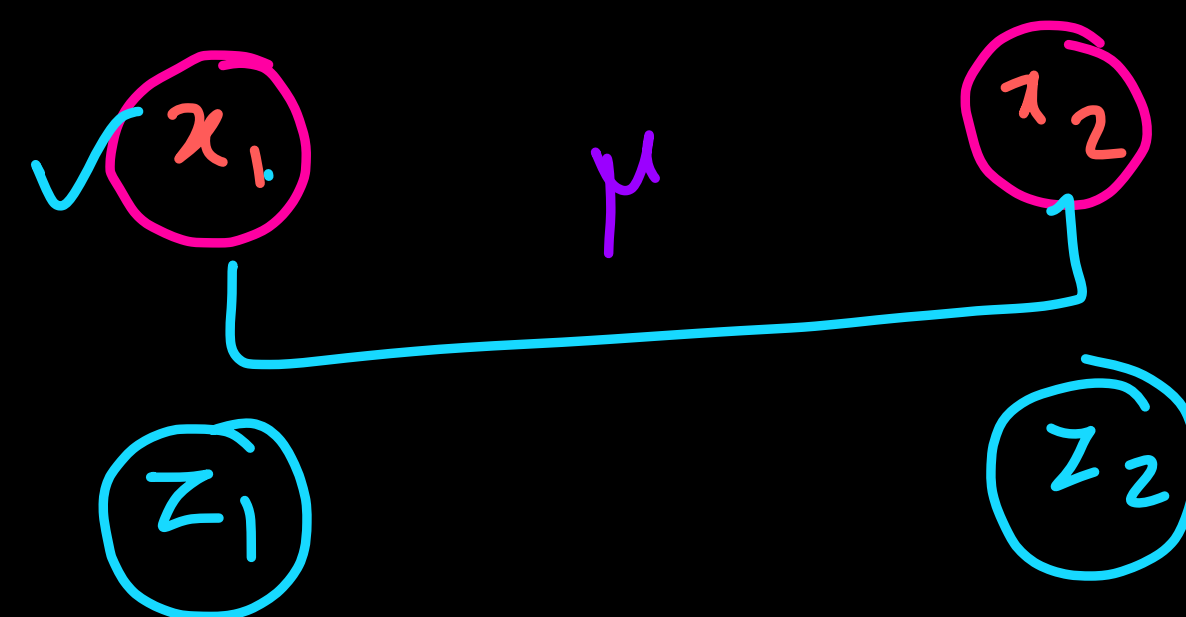
→ 99%



$$z_1 = \text{norm.ppf}(0.025)$$

0.5%

2.5%



~~norm.ppf~~

$$\text{norm.ppf}(1 - 0.025)$$


$$x_1 = 65 + (-1.95) * 0.25$$

=

$$\Rightarrow [64.51 \quad 65.48]$$

$$S.D = \frac{25}{\sqrt{100}}$$

$$= \boxed{0.25}$$

$$C.I = \bar{x} \pm z * \left(\frac{\sigma}{\sqrt{n}} \right)$$


$$\left[\bar{x} - \boxed{z} * \frac{\sigma}{\sqrt{n}} , \bar{x} + z * \frac{\sigma}{\sqrt{n}} \right]$$

$\bar{x} \rightarrow$ sample mean

$\boxed{z} \rightarrow$

$\frac{\sigma}{\sqrt{n}} \rightarrow$ standard error

The sample mean recovery time of 100 patients after taking a drug was seen to be 10.5 days with a standard deviation of 2 days. Find the 95% confidence interval of the true mean.

↑
population mean

$$\begin{aligned}n &= 100 \\ \bar{x} &= 10.5 \\ \sigma &= 2\end{aligned}$$

$$S.E = \frac{\sigma}{\sqrt{n}} = \frac{2}{10} = \underline{\underline{0.2}}$$

From a sample of 80 endangered birds, the average wingspan was found to be 45 cm, with a population standard deviation of 10 cm. What is the correct confidence interval of the mean wingspan of the entire population with 90% confidence.

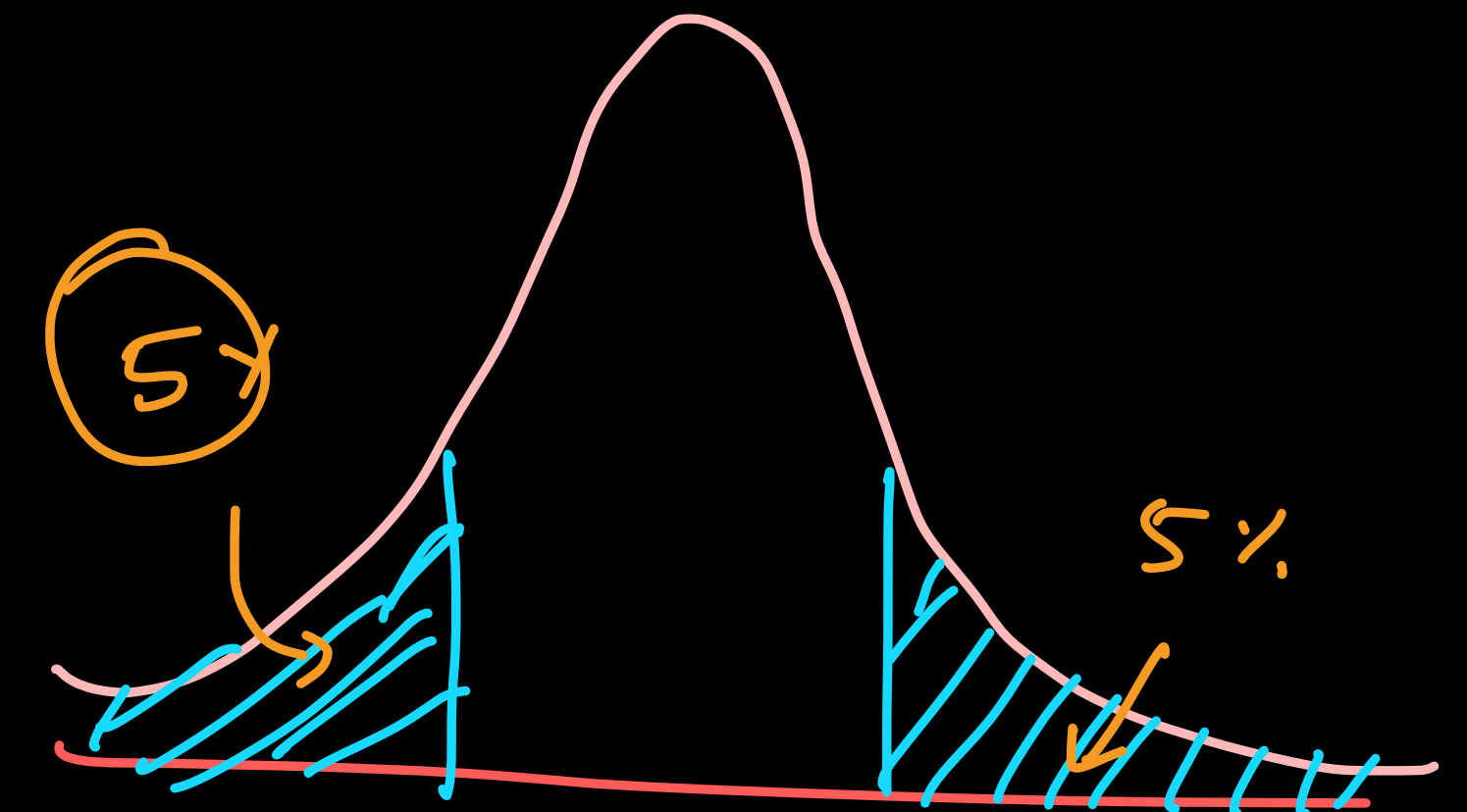
$$n = 80$$

$$\bar{x} = 45 \text{ cm}$$

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

$$S.E. = \frac{\sigma}{\sqrt{n}} = \frac{10}{\sqrt{80}}$$

90%



In a software project, the team estimates bug resolution time at an average of 6 hours with a standard deviation of 2 hours.

To estimate the mean resolution time with 99% confidence, the project manager samples 25 resolved bugs.

What is the correct confidence interval?

$$n = 25$$

99%.

$$\bar{x} = 6$$

$$\sigma = 2$$

$$S.E. = \frac{2}{\sqrt{25}} = \frac{2}{5} = \boxed{0.4}$$

$$z_1 = \text{norm.ppf}(0.005)$$

$$z_2 = \text{norm.ppf}(1 - 0.005)$$

$$x_1 = 6 + \boxed{z_1} \times 0.4$$

$$x_2 = 6 + \boxed{z_2} \times 0.4$$

Normal

2-3

CLT \rightarrow 2-3
 \downarrow

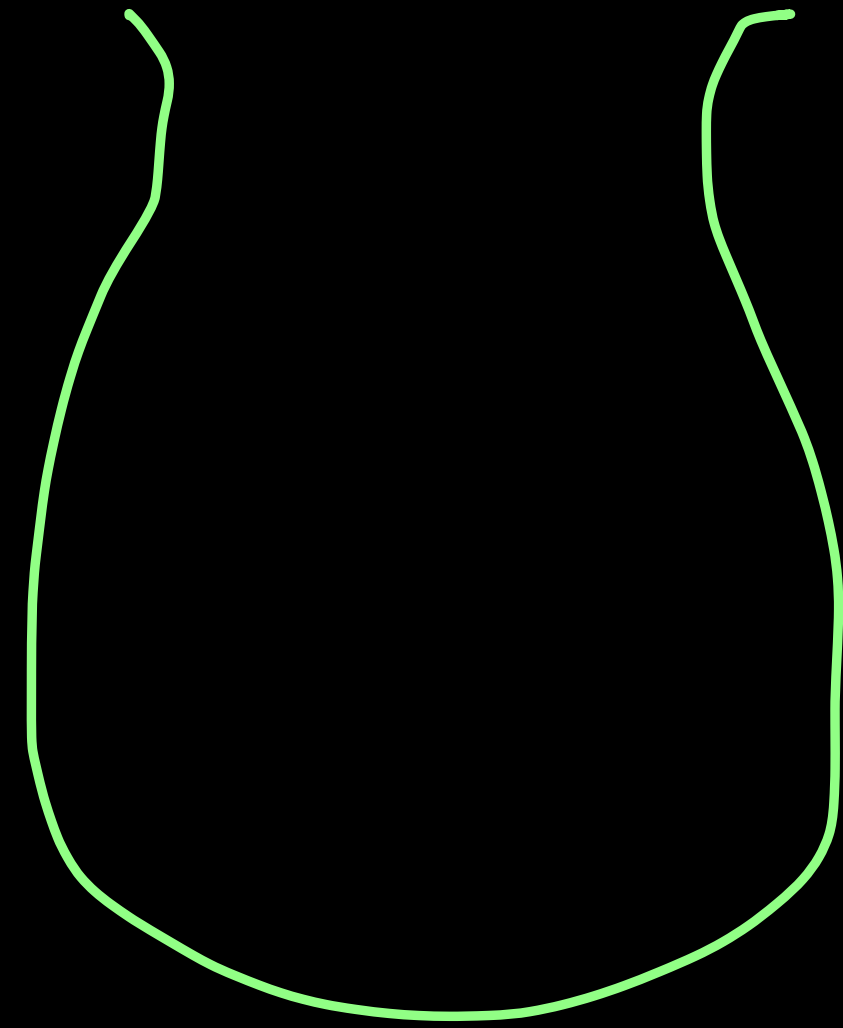
CI \rightarrow 2-3

CLT

$$\bar{X} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

✓✓

C. I \rightarrow using CLT



10 time



$$\left\{ \begin{array}{l} x_1 \\ x_2 \\ \vdots \\ x_{100} \end{array} \right. = \left[\begin{array}{l} \\ \\ \\ \end{array} \right]$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \Rightarrow \boxed{\bar{x}_1}$$

$$= \bar{x}_2$$

Sample means

\sim

Normal distribution

→ Sample

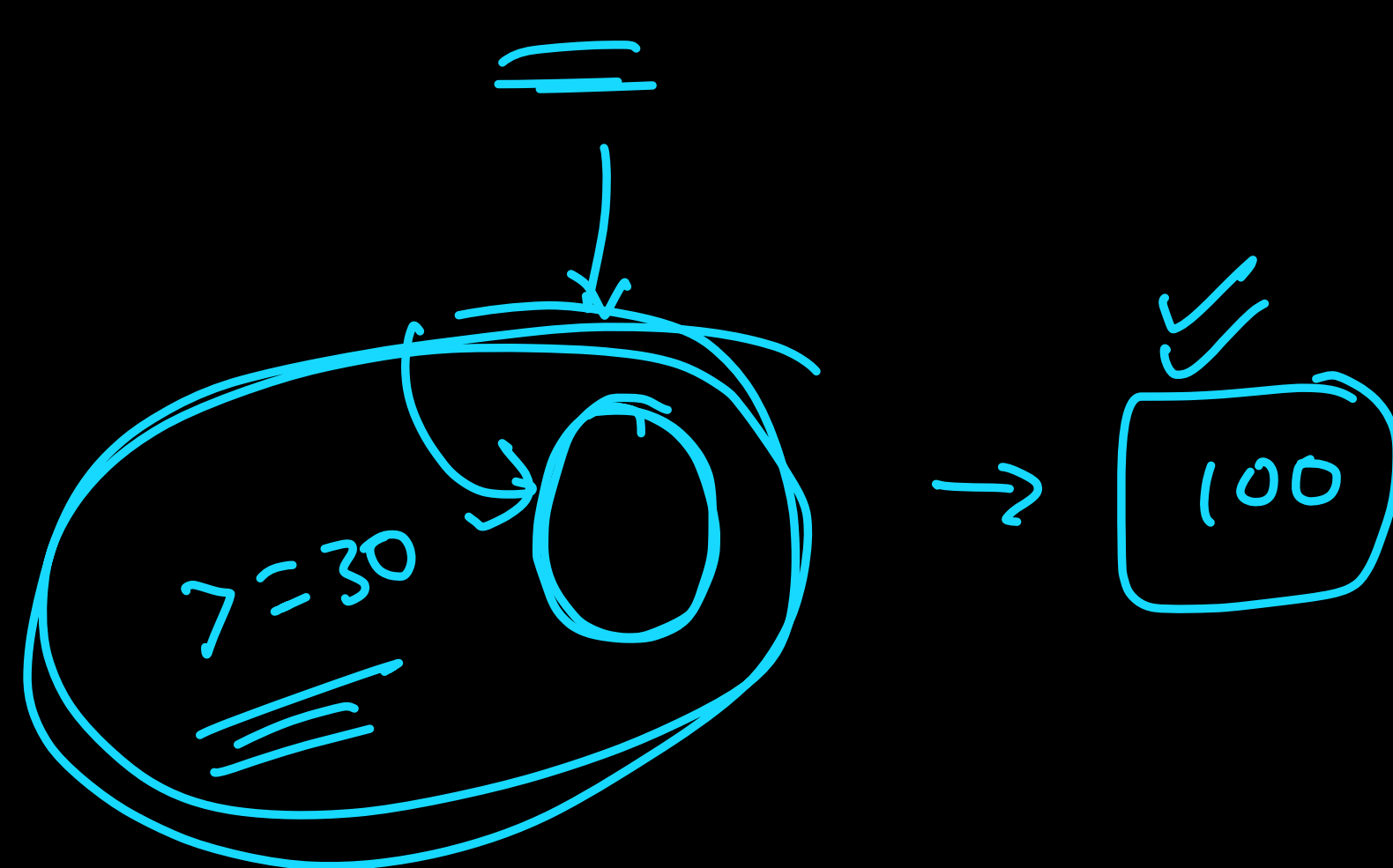


Diagram illustrating the standard deviation calculation:

$n =$ $\sigma =$

$\Rightarrow \sigma \Rightarrow S.D = \sigma / \sqrt{n}$

\Rightarrow S.D of Sample

Diagram illustrating the standard deviation calculation:

\Rightarrow $\sigma =$

σ / \sqrt{n}

Diagram illustrating the standard deviation calculation:

$\frac{3.84}{\sqrt{100}}$

$\Rightarrow 0.384$

