

Till now:

→ Pmf:

→ Discrete Uniform distⁿ —
All values have
equally occurrence.

→ Bernoulli distribution

↳ 2 possible outcome

→ Binomial → n- bernoulli trial.

→ Geometric — first success after n trials

→ Poisson → Avg no of occurrence in a time interval

*

Probability density fn

↳ for continuous random
variable

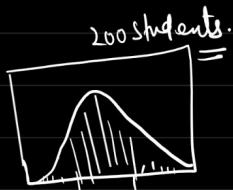
* Normal distribution | Gaussian distribution | Bell-shape distribution

In statistics, a **normal distribution** or **Gaussian distribution** is a type of continuous probability distribution for a real-valued random variable. The general form of its probability density function is

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

→ Most of the naturally occurring / human generated follow a Normal distribution.

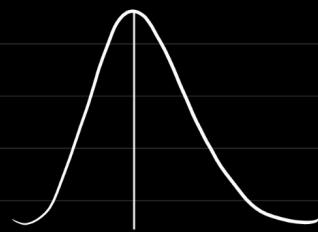
Why ??



→

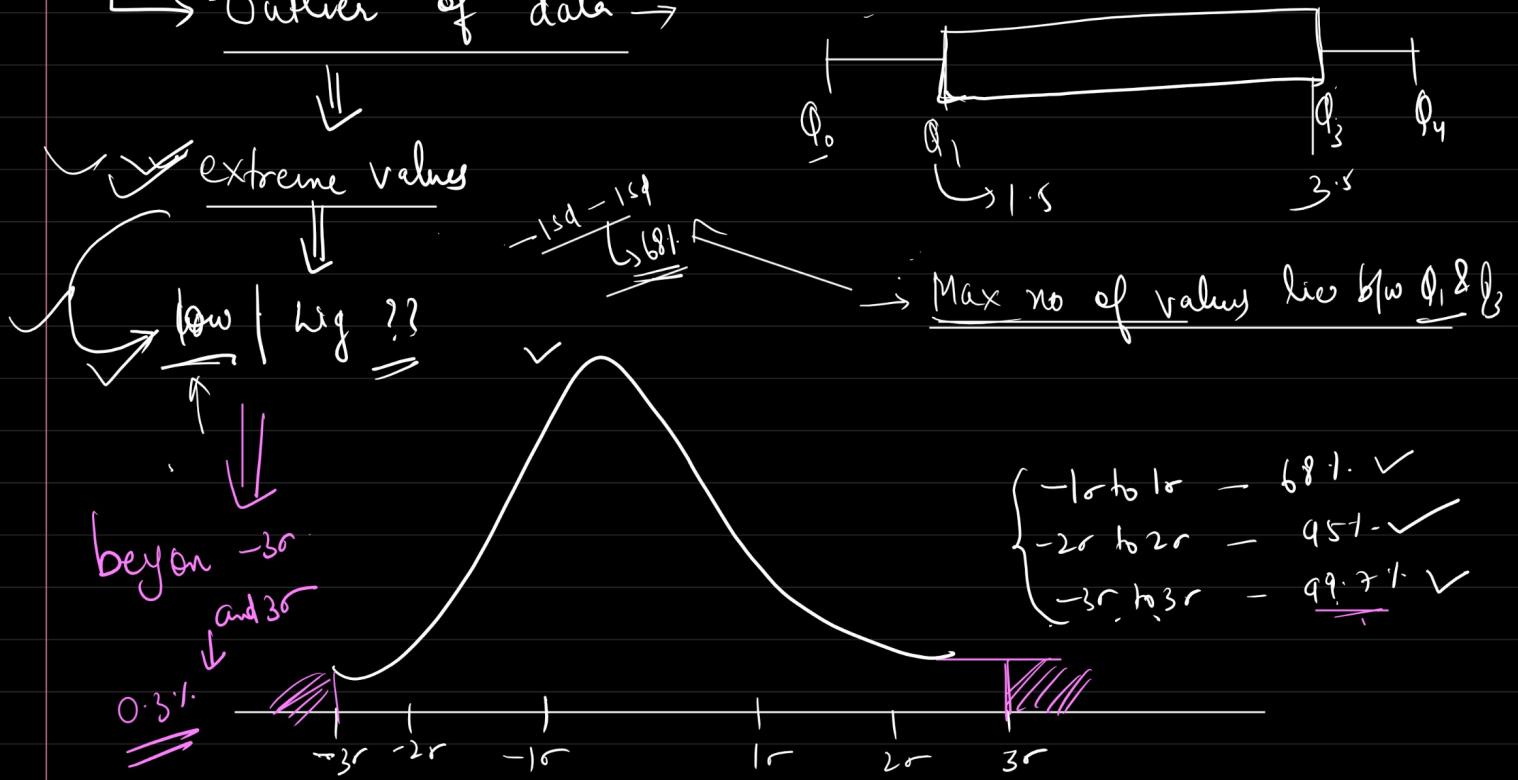
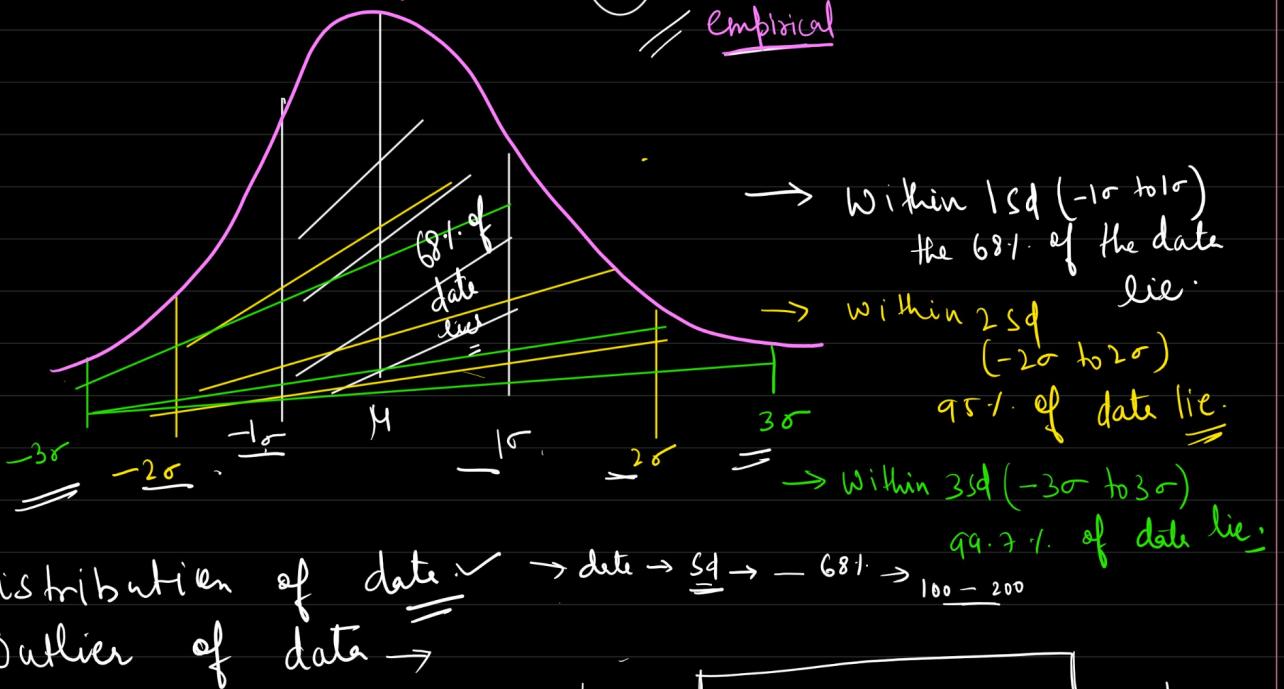
Measure the height of
Students

Characteristics of N.D



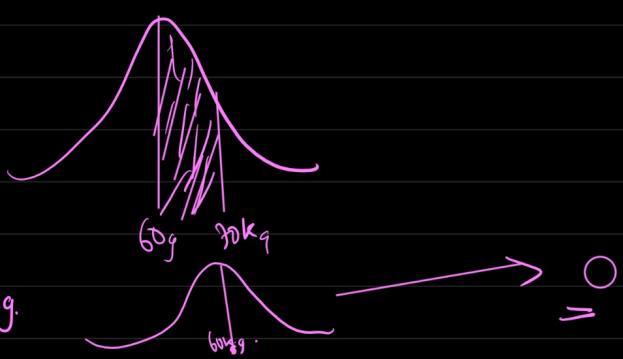
- Symmetric
- Skewness = 0
- Kurtosis = 3
- Mean = median = mode. μ

7.5th wonder of world: figure \rightsquigarrow 68-95-99.7%



* prob densit

↳ Calculate the prob \rightarrow person will fall between 60 kg & 70 kg.



- Height of people
- Weight
- Measurement of error
- Score
- Blood Pressure.

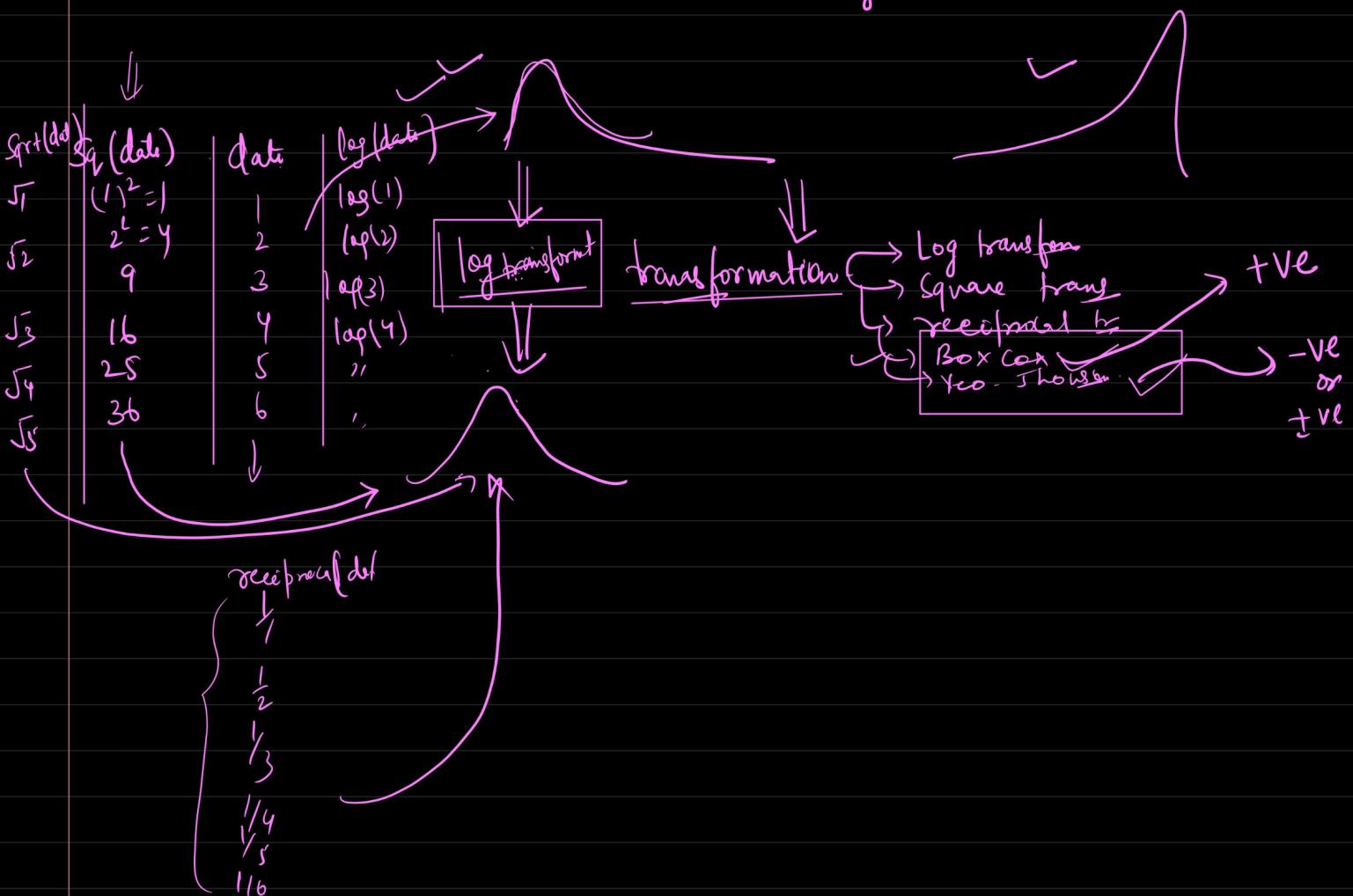
~~Mean is 0, Std = 1 for Normal distn.~~

* industry use case of normal distn

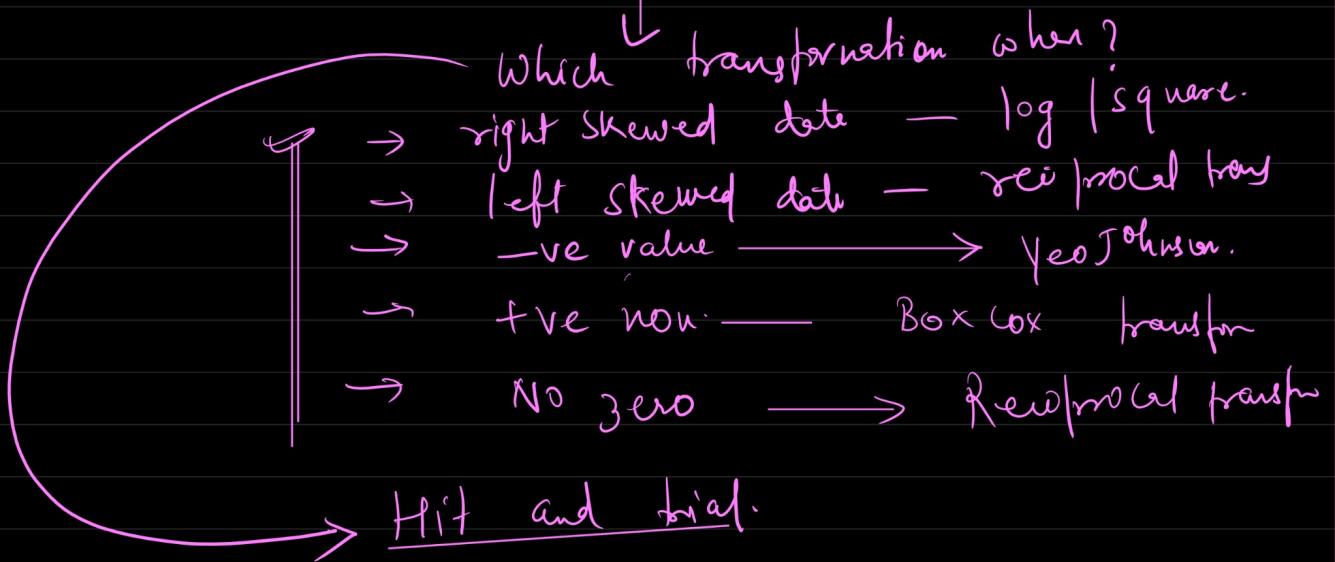
↳ Beyond $3\sigma \rightarrow$ outliers.

↳ Some of ML statistical models, requires the data to be Normally distributed.

- ① Most of data follows N.D
- ② Outlier identification.
- ③ Prediction will be wrong



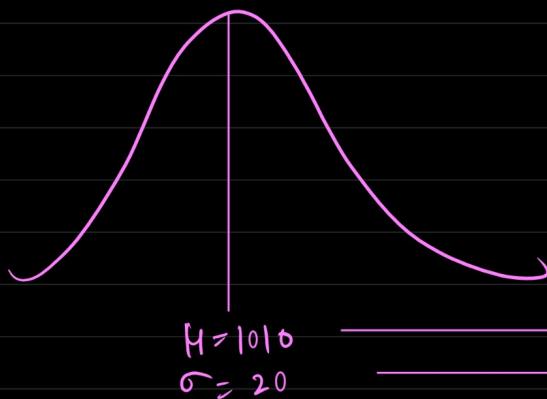
→ if you are modelling statistical ML models, then data should follow normal distⁿ and if not do transformation.



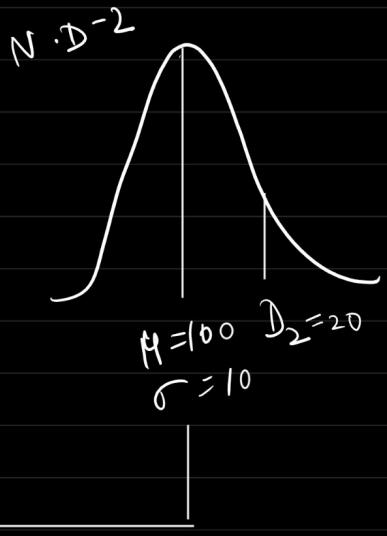
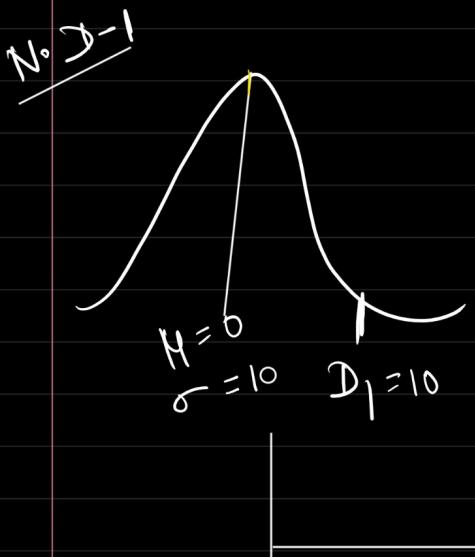
* for ND mean = 0, $\sigma = 1$

↓
true or not.

✓ Mean & std devⁿ for Normal distribution
Can be anything



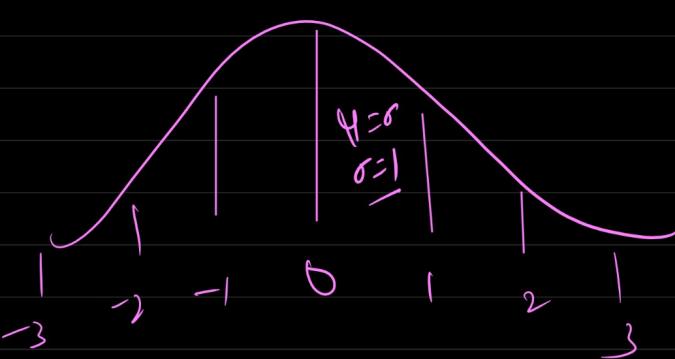
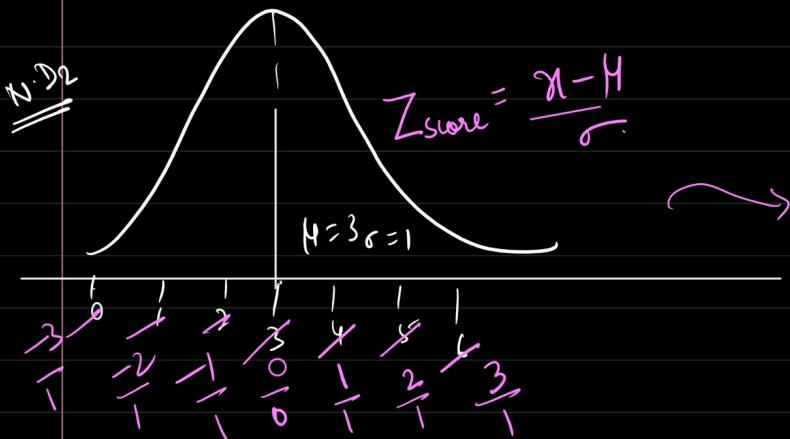
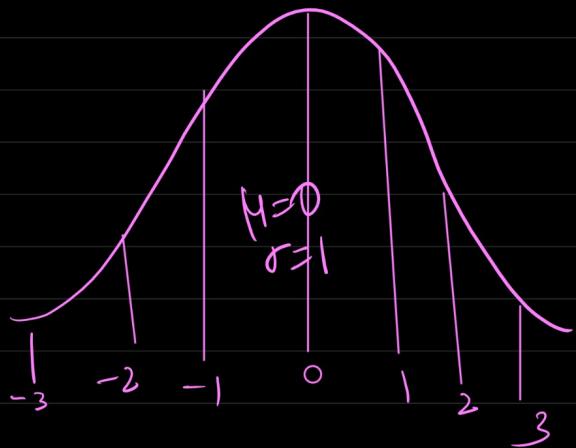
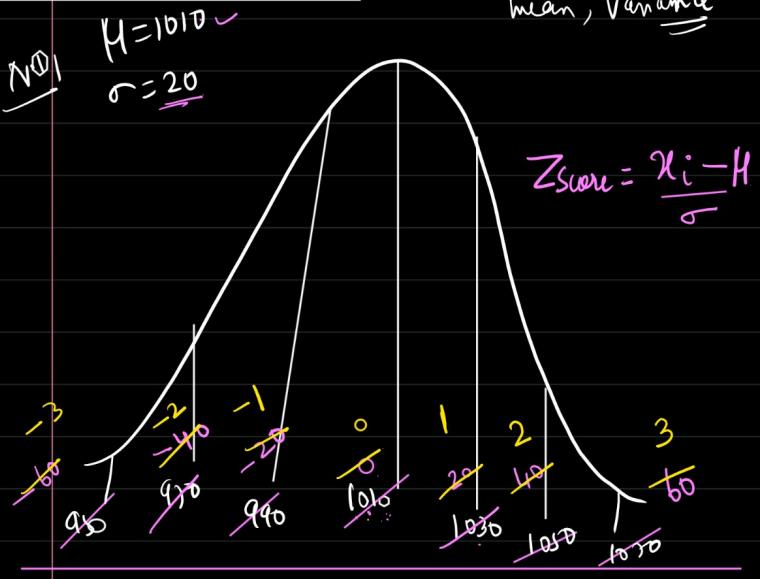
→ mean of the data is 1010
 → On an avg the data point is 20 units away from mean



→ which dp among D_1, D_2 is closer to μ ?
 ||
We can't say

Standard Normal distn Scales of both the distn is not same:
 ↓
 mean, Variance

two classes
 SC \hookrightarrow Arts.

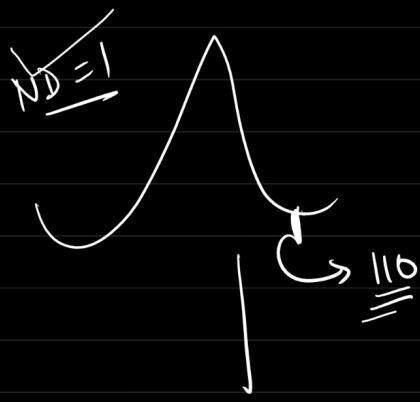
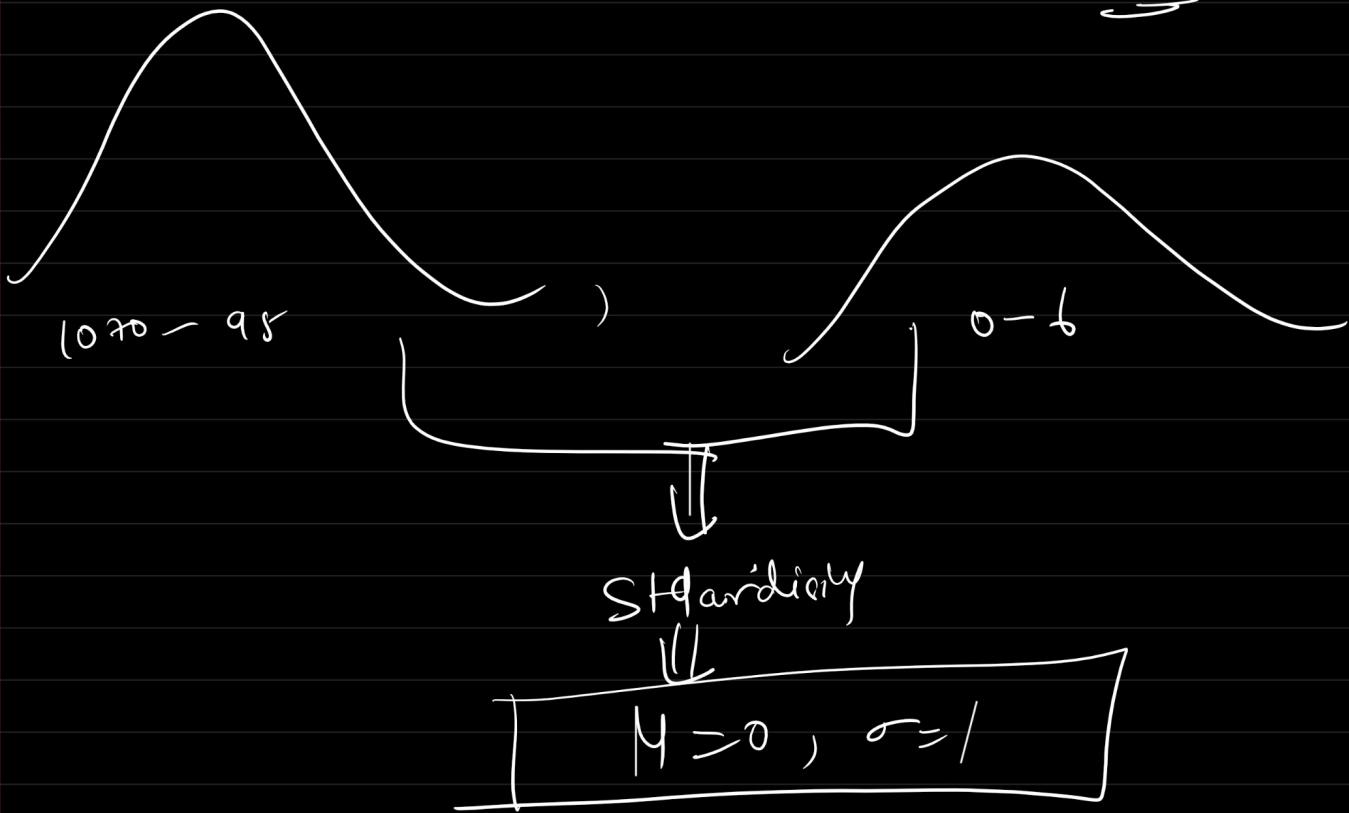


Std N.D

↳ Standardisation of data

$$\mu = 0, \sigma = 1$$

* The normal distⁿ which has $\mu = 0, \sigma = 1$ is
Std normal distⁿ



Std

$$\mu = 0,$$

$$\sigma = 1$$

$$Z_{95\%} = \frac{X - \mu}{\sigma} = \frac{ND - 100}{20} = \frac{90 - 100}{20} = -5$$



$$Z_{95\%} = \frac{200 - 0}{1} = 200$$

Whichever Z score is highest, the dp will be far from mean.

Q $N(\mu = 50, \sigma = 20, D_1 = 110)$

Can you tell how many standard deviation away this data point is lying away from mean value?

$$\rightarrow Z\text{score} = \frac{x-\mu}{\sigma} = \frac{110-50}{20} = \frac{60}{20} = 3 \text{ s.d away from mean}$$

Q $N(\mu = 100, \sigma = 10, D_2 = 200)$

Can you tell how many std dev away D_2 is lying away from mean?

$$\rightarrow Z\text{score} = \frac{200-100}{10} = 10 \text{ s.d away from mean}$$

Q Which dp from D_1 , $2D_2$ is more away from mean value? $\rightarrow D_2$.

Scenario -1

$$Z\text{score} = 3 \text{ cms}$$

It means that the data point is 3 std.

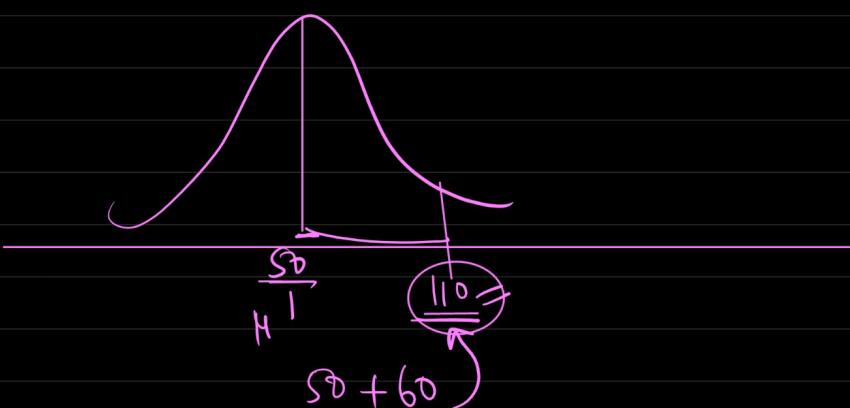
dev away from mean = $\sigma = 20$ $\rightarrow 3 \cdot \underline{\text{std}}$

Scenario -2

$$\sigma = 3 \text{ cm}$$

each and every data point on an avg is 3cm away from mean

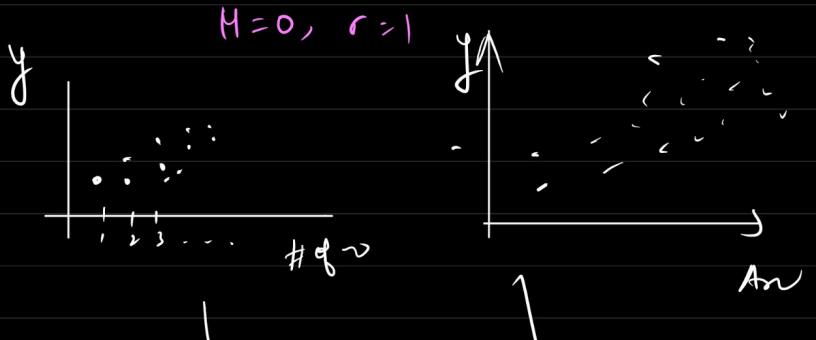
$$D_1 \rightarrow 20 \times 3 = 60 \text{ units away from mean}$$



* industry use case of std N.D

$$\hookrightarrow \text{Standardization} = \frac{x - \mu}{\sigma}$$

$\frac{x-\mu}{\sigma}$	# of room	Area of house (sq ft)	Price of house (y) (in Cr)
1.5	1	1100 sq ft	5
1	2	900 sq ft	3
2	1	2000 sq ft	8
1	2	900 sq ft	-
5	1	950 sq ft	-
8	1	1100 sq ft	-
3	1	11000	-
9	1	-	-
10	1	-	-



Both data have different scales.

→ learning parameter will also vary much in terms of scale.

Scaling (optional)



~~Robust scalar~~ Scaling

→ Standardization ($\mu=0, \sigma=1$)
(outliers)

→ Normalization (0 and 1)
(min-max scalar)

$$x_{\text{norm}} = \frac{x_i - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

1	$\frac{1-1}{5-1} = 0$
2	$\frac{2-1}{5-1} = \frac{1}{4} = 0.25$
3	$\frac{3-1}{5-1} = \frac{2}{4} = 0.5$
4	$\frac{4-1}{5-1} = \frac{3}{4} = 0.75$
5	$\frac{5-1}{5-1} = 1$

$$\frac{x_{\text{min}} - 1}{x_{\text{max}} - 1} = 0$$

$$\frac{2-1}{5-1} = \frac{1}{4} = 0.25$$

$$\frac{3-1}{5-1} = \frac{2}{4} = 0.5$$

$$\frac{4-1}{5-1} = \frac{3}{4} = 0.75$$

$$\frac{5-1}{5-1} = 1$$

Why Scaling?

→ Better interpretation

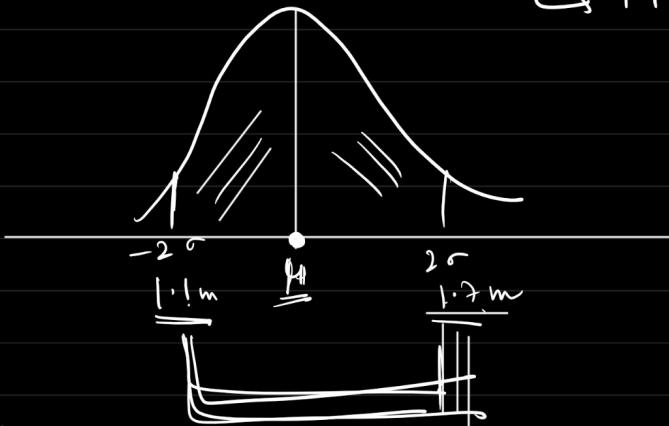
→ faster optimisation (gradient descent)

Q. 95% of students at school are between 1.1 m and 1.7 m tall. Can you calculate mean and std dev.??

→ ~~-2σ to 2σ~~ 95% → 1.1 m and 1.7 m tall

→ ~~-2σ to 2σ~~
↳ 1.1 to 1.7

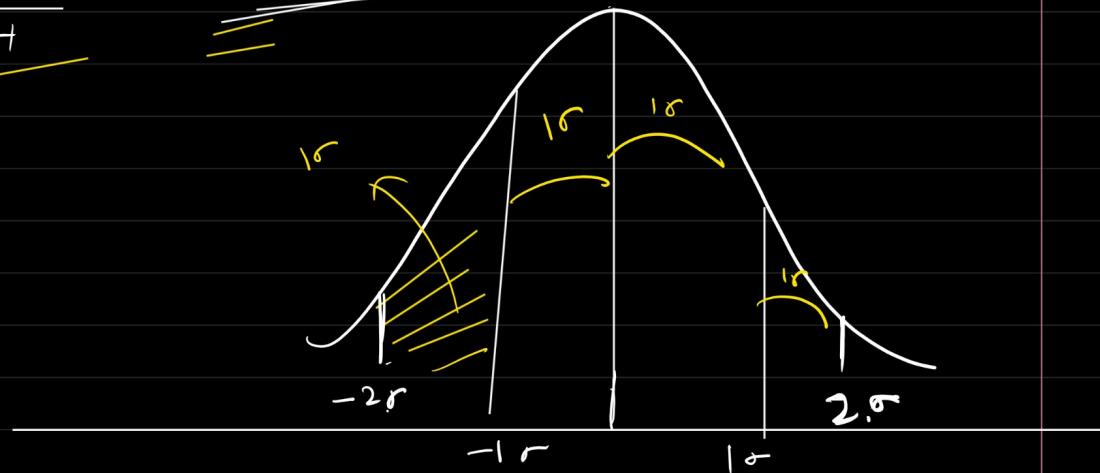
$$\frac{1.1 + 1.7}{2} = 1.4 \text{ m}$$



95% is ~~2 s.d.~~ 2 s.d.

$$4 \text{ s.d.} = 1.7 - 1.1$$

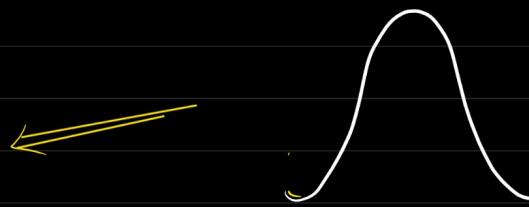
$$1 \text{ s.d.} = \frac{1.7 - 1.1}{4} = 0.15 \text{ m}$$



* Central limit theorem (8th wonder of world)

* for a N.D

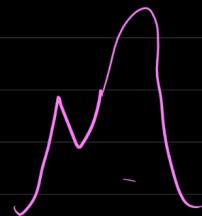
- Symmetric
- Skewness = 0
- Kurtosis = 3
- mean = median = mode.
- 68-95-99.7 rule.



$$\text{std}(N.D) - \mu = 0, \sigma = 1$$



Haphazard / irregular

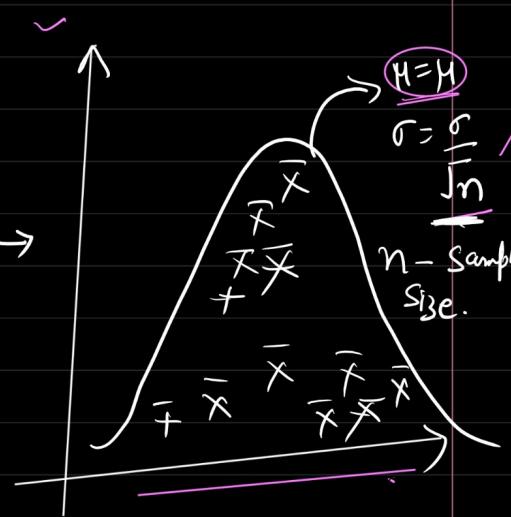


Central limit theorem → A way to convert any form of distribution to Normal distribution.

Sample with replacement.



$$\begin{aligned} S_1(50) &\rightarrow \bar{X}_{S_1} \\ S_2(50) &\rightarrow \bar{X}_{S_2} \\ S_3(50) &\rightarrow \bar{X}_{S_3} \\ &\vdots \\ S_{20k}(50) &\rightarrow \bar{X}_{S_{20k}} \end{aligned}$$



→ The CLT states that if you have any pop with μ, σ , and you take sufficiently large random sample from population with replacement, then the distribution of Sampling mean will be approximately

normally distributed.

* Sampling mean of population (μ, σ) will approximate to a normal distribution.

Population (μ, σ) → large no of Sample → Sampling → ND

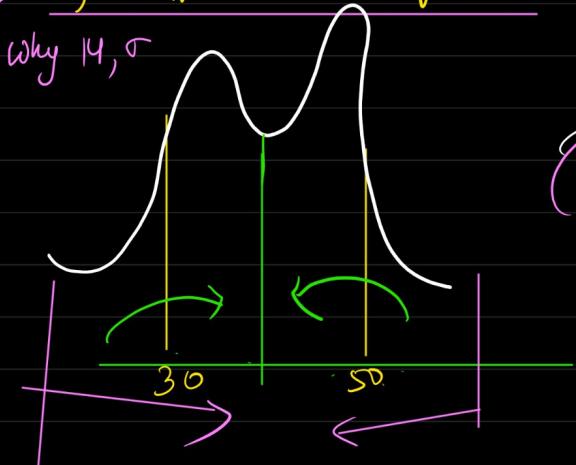
$\mu = \bar{\mu}$
 $\sigma = \sigma / \sqrt{n}$

Two conditions of CLT

- ✓ → The no of samples should be large enough.
- The Sample size should be equal to and greater than 30 (Except pop in ND)

Why $n \geq 30$?? → Sampling mean follow Z_distrib $\rightarrow z_{\text{sur}} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$

CLT ??
 Why sufficient large no of sample?

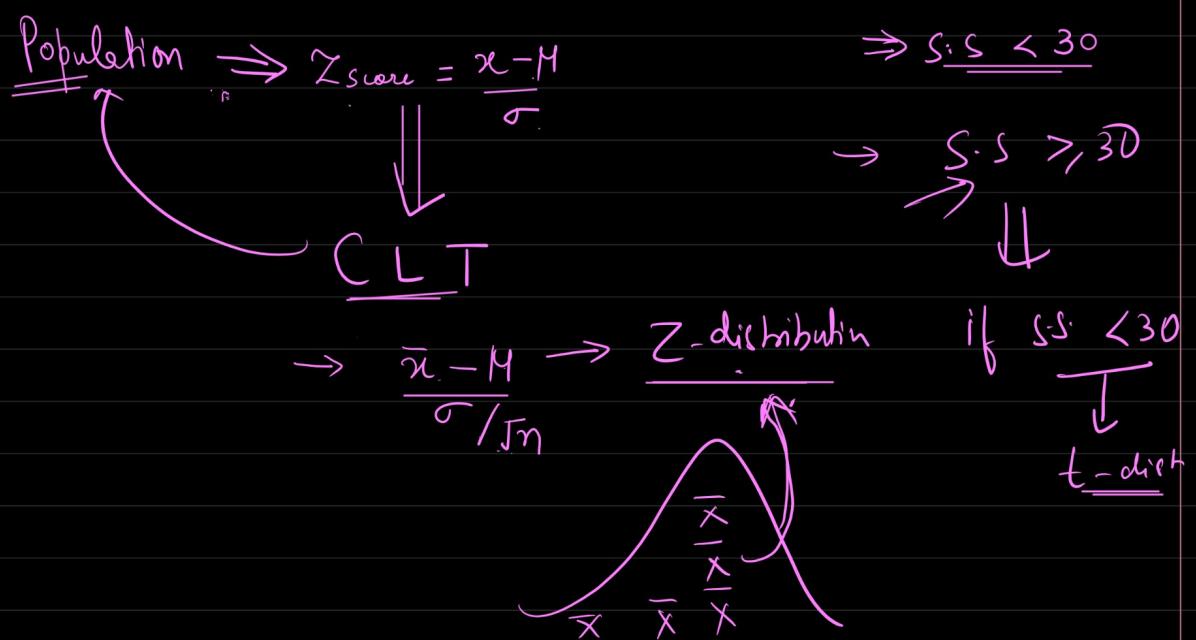
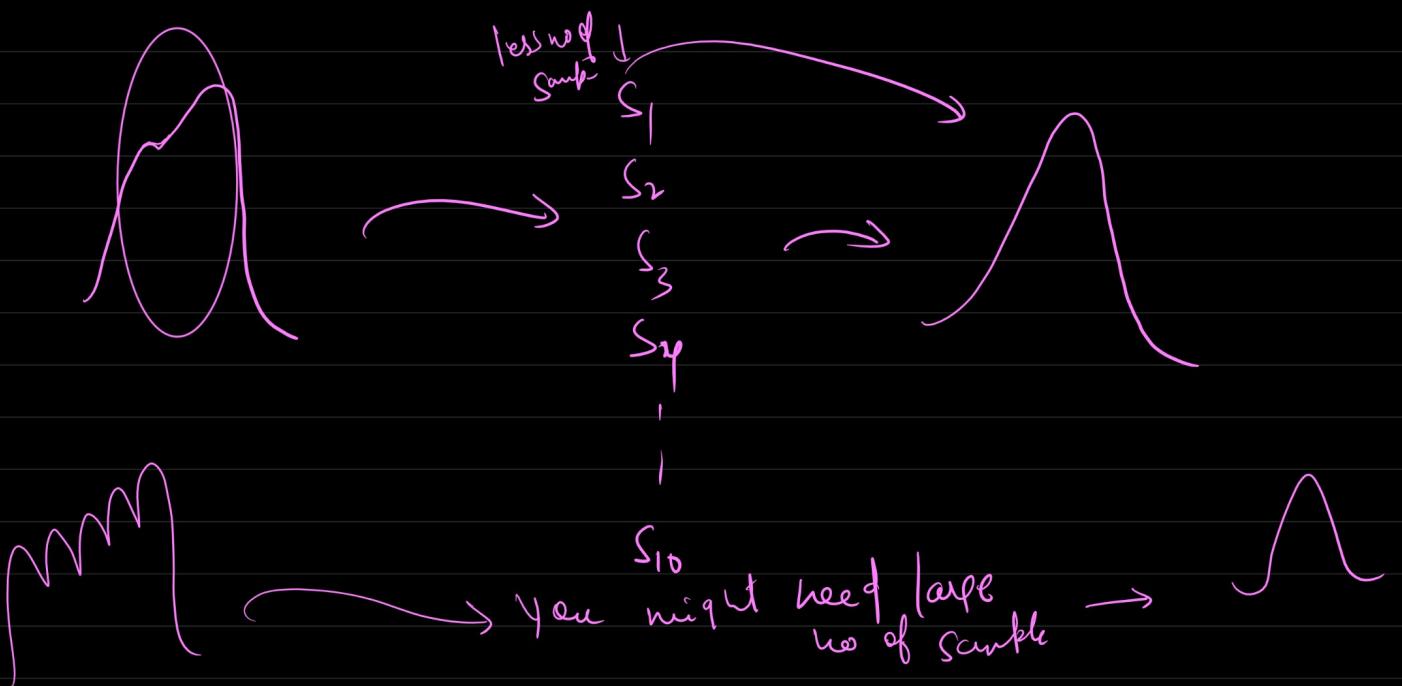


$$\begin{array}{c} \bar{X}_1 \\ \bar{X}_2 \\ \bar{X}_3 \end{array} \rightarrow \frac{30 \text{ & } 50}{40}$$



HYP / FIT

$\xrightarrow{\text{N.D}}$ Most (68%) → will lie in the catn to 15%



Why μ is μ after CLT?

Why σ is σ / \sqrt{n} after CLT?

(without)

$$2 \rightarrow 2$$

$$3 \rightarrow 3$$

$$\frac{5}{5} \rightarrow 4$$

$$S_1(1)$$

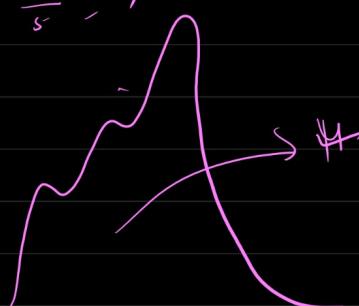
$$S_2(1)$$

$$S_3(1)$$

$$S_4(1)$$

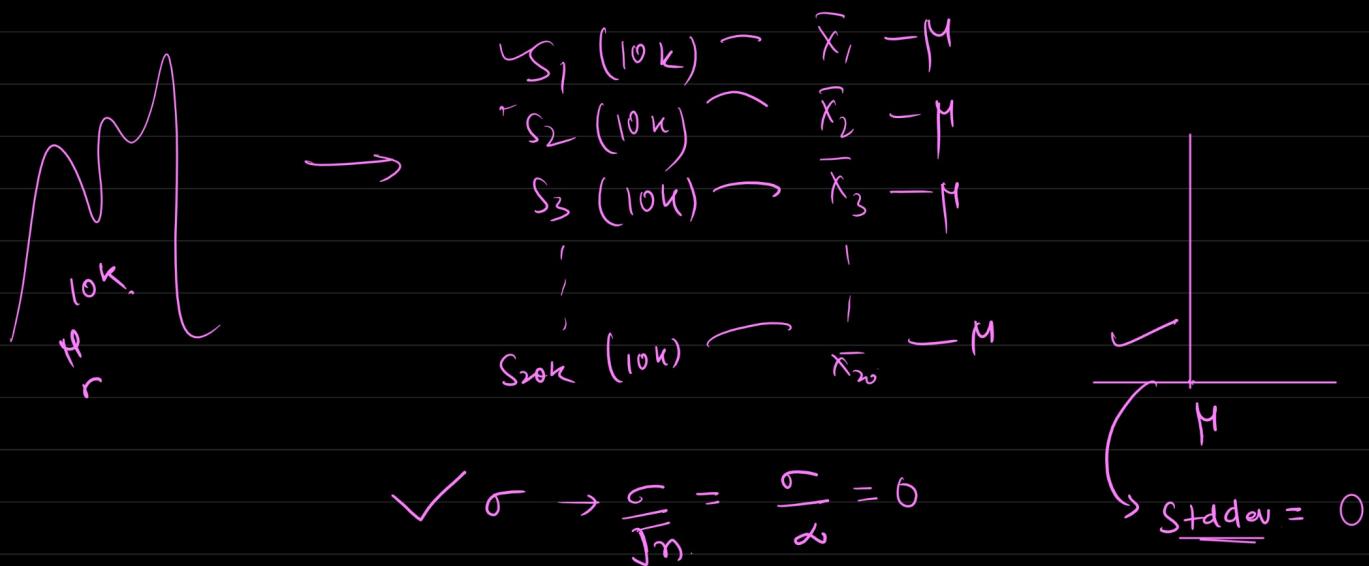
!

\rightarrow





$n = \infty$
 $S(\sigma) \Rightarrow$ Whatever no of elements are there in population, you take all element in sample.



Q You have a population with $\mu = 100$ and std dev $\sigma = 20$. If you have sample size of 100 from this population. What is the prob that sample mean will be less than 105.
 $\Rightarrow \mu = 100, \sigma = 20, n = 100, \bar{x} = 105$

$$Z_{\text{corr}} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{105 - 100}{20 / \sqrt{100}} = \frac{5}{2} = 2.5$$



Using Z table $\rightarrow \text{Area} = 0.89$

area upto Z

