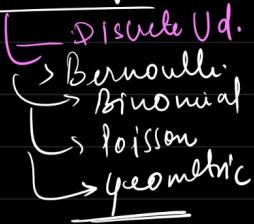
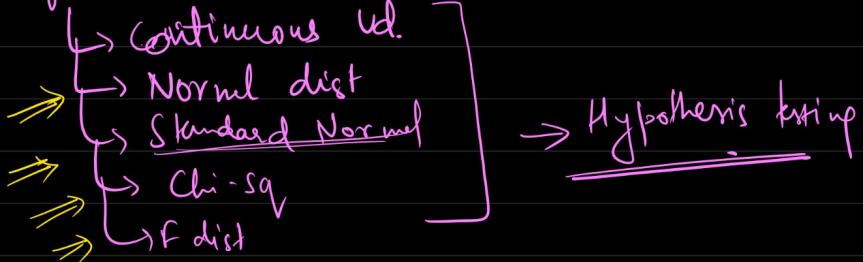


* Probability mass fn



* Probability fn →



→ Hypothesis testing

* continuous Uniform distn

$$\text{mean} = \frac{1}{2}(a+b)$$

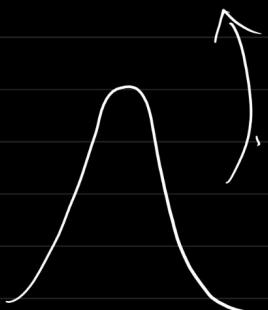
$$\text{varian} = \frac{1}{12}(b-a)^2$$



* Normal distn

↳ Gaussian distn

↳ Bell shaped distn.



→ Most of the naturally occurring human generated data follow a Normal distn.

Characteristics of NP

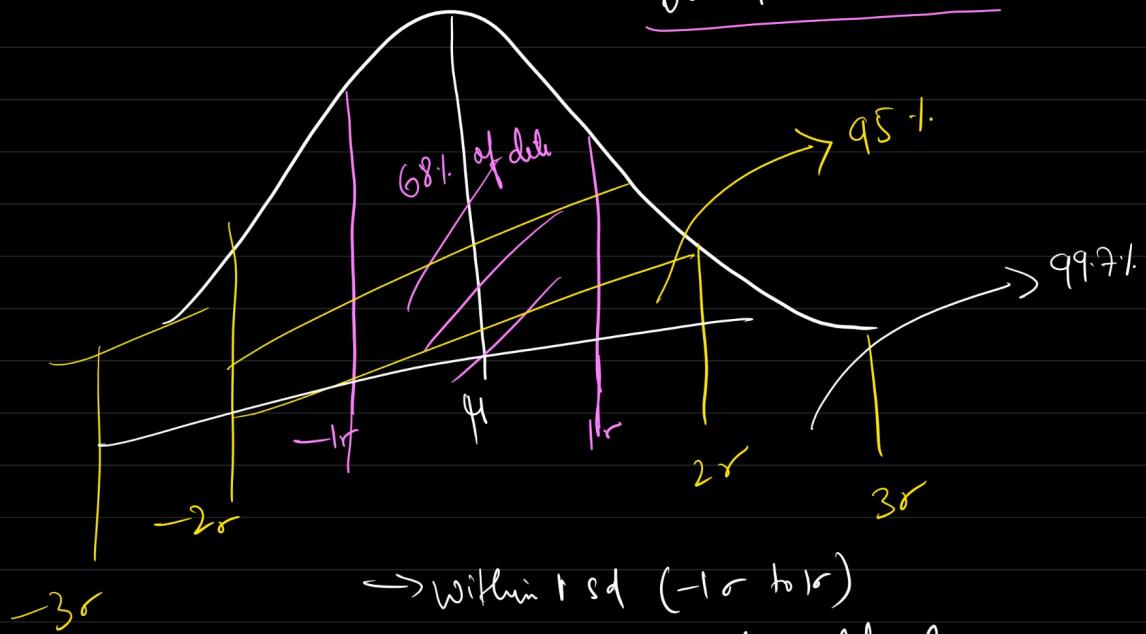
→ Symmetric ✓

→ Skewness = 0

⇒ Kurtosis = 3

→ Mean = Median = mode

rule:- $68 - 95 - 99.7\%$.

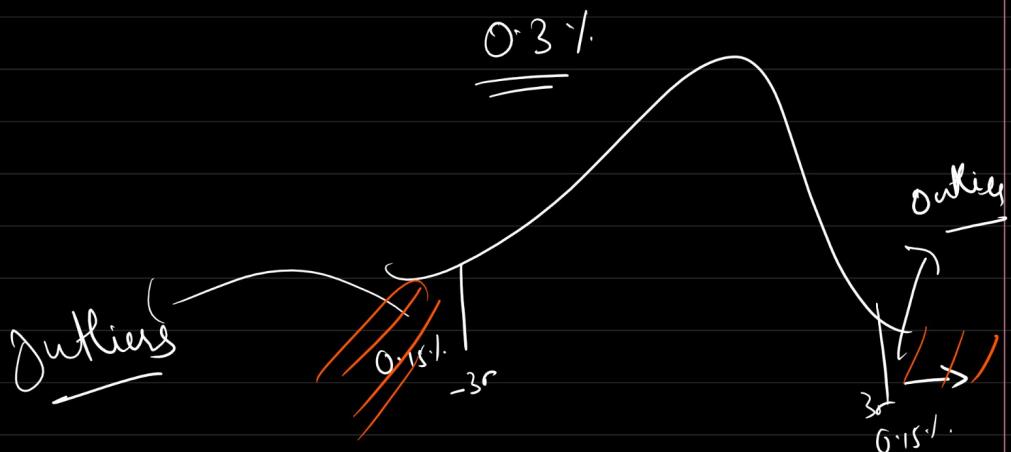


→ Within 1 sd (-1σ to +1σ)

68% of the data lie.

→ 2sd = 95.4%

→ 3sd = 99.7%



2.6% dP 3σ → beyond }
 -3σ → beyond } 3 outliers
3 outliers

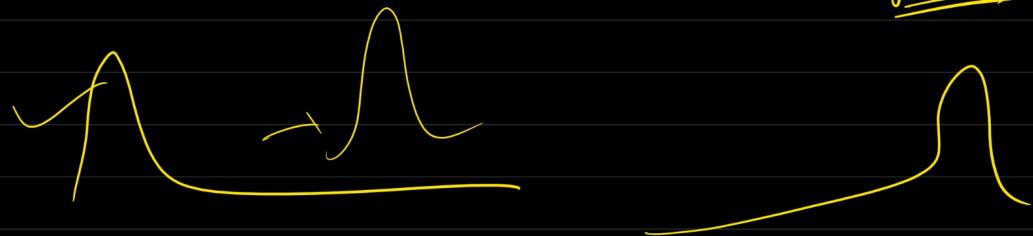
Example

- eg. Ht of people
- eg. Wt .. "
- eg. Score
- eg. BP
- eg. Errors
- eg. IQ

* Use Case

Beyond $3\sigma \rightarrow$ Outlier.

Required the data to be Normally distⁿ



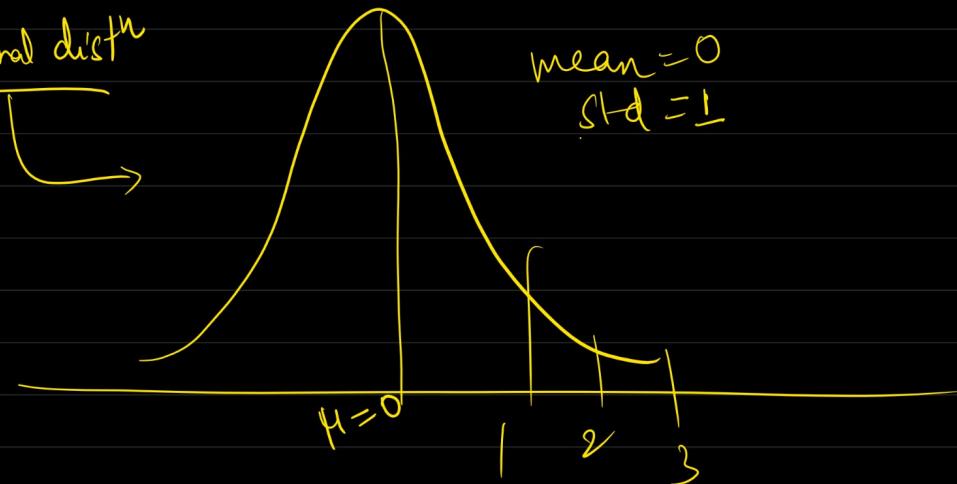
transf ⁿ	log base ⁿ →
square m →	
recip ⁿ →	

{	right skewed data	- log / square
	left " "	- reciprocal
	-ve value	- Ave Thomson
	+ve	- Box Cox transform
	No zero	- Reciprocal loss

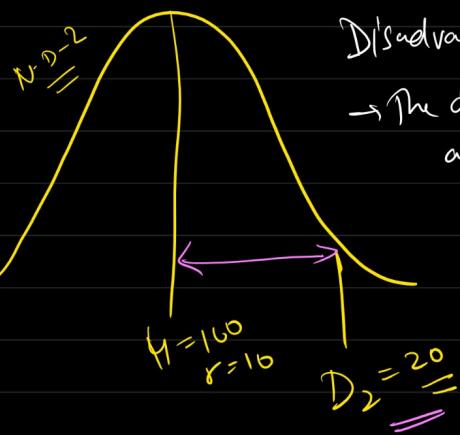
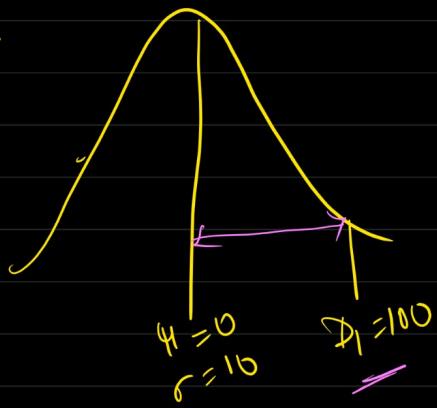
Nd - mean & sd

* Standard Normal distⁿ

mean = 0
std = 1



Why?



Disadvantage of SD
→ The data is on a different scale

Use Case

- Scaling
- To compare two disⁿ

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

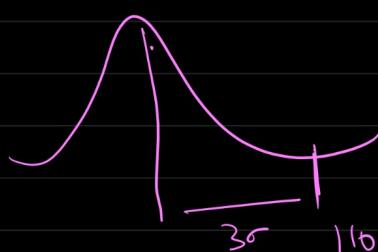
How many s.d away the dp is from mean?

→ Zscore = $\frac{x-\mu}{\sigma} = \frac{110-50}{20} = \frac{60}{20} = 3$ s.d away from mean.

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

$\frac{200-100}{10} = 10$ s.d away.

far from fint.



Scen-1

Zscore = 3 cm. ($Z_{SD} = \frac{x-\mu}{\sigma}$)
It means that
the data point
is 3 std away from mean

Scen-2

(Standar dev.) = 3 cm.
Each and every dp
is on an avg
is 3 cm away from
mean.

Use SND → Scaling (optional)

Computer

$$\begin{array}{l}
 2 \times 3 = 6 \\
 5 \times 2 = 10 \\
 2 \times 4 = 8 \\
 3 \times 2 = 6 \\
 \hline
 88 + 92 =
 \end{array}$$

ML		
# of rooms	Area of house	Price of house
1	1000	5
2	900	3
1	2000	8
2	1	
2		
5		
8		
2		
1		

ML → To learn pattern from data

Mathematical relation

(mathematical calc)

Scaling

$$\begin{aligned}
 & \begin{array}{l} 1100 \rightarrow \\ 900 \\ 200 \end{array} & \begin{array}{l} 1100 - 800 \\ \hline 100 \end{array} \\
 & \xrightarrow{\quad} & \\
 & \begin{array}{l} 900 \\ 200 \end{array} & \begin{array}{l} 900 - 800 \\ \hline 100 \end{array} \\
 & \xrightarrow{\quad} & \\
 & M = 800 & \\
 & \sigma = 100 &
 \end{aligned}$$



Scaling

Standardisation (SND)

Normalisation

Q Q5. If students at school are b/w 1.1m and 1.7m tall. mean & std dev ??

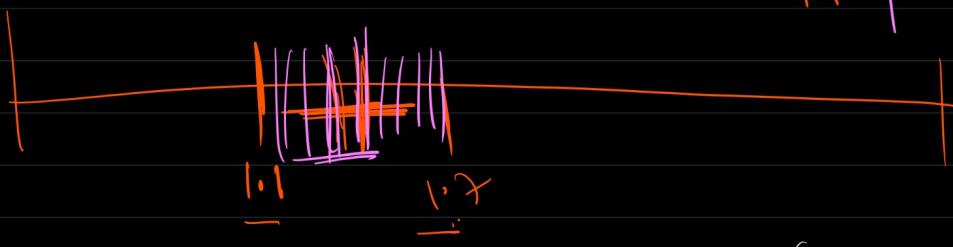
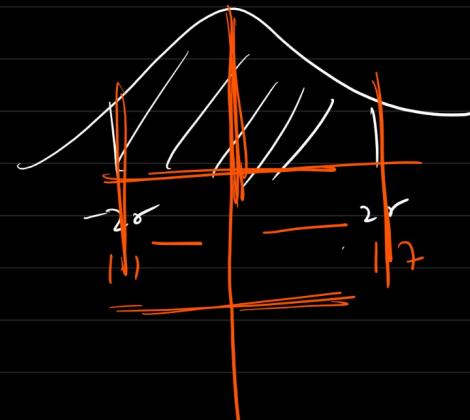
$$\underline{95.1} = 1.1m - 1.7m$$

$$\begin{array}{r} 68 - 95 - 99 - 7 \\ \hline \end{array}$$

$$\text{Mean} = \frac{11+17}{2} = 14$$

$$140 = 1 \cdot 7 - 11$$

$$\sigma = \frac{17-11}{4} = 0.15$$



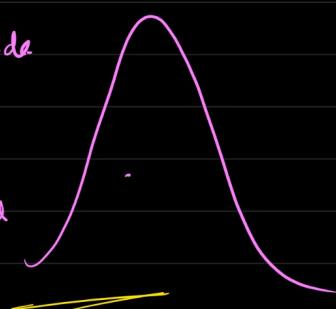
8th wonder of world → CLT (Central Limit Theorem)

(1)

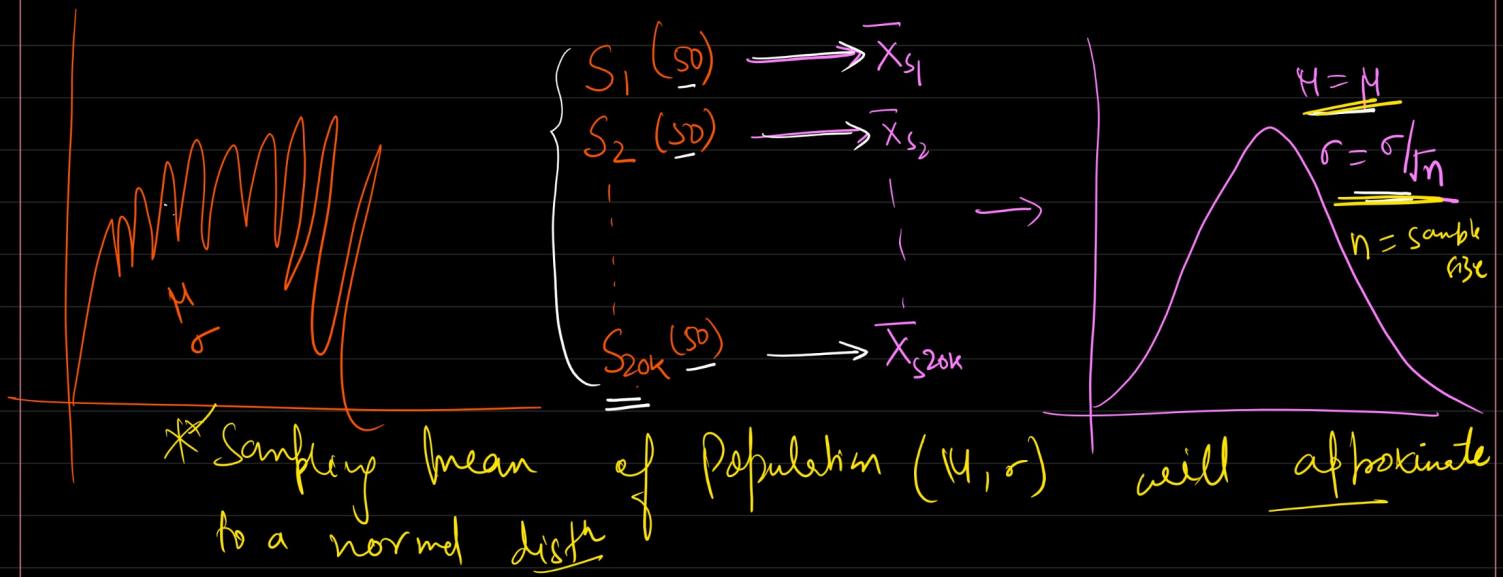
{

- Mean = median = mode
- Skewness = 0
- Symmetric
- $68 - 95 - 99.7.1.$ rule

—



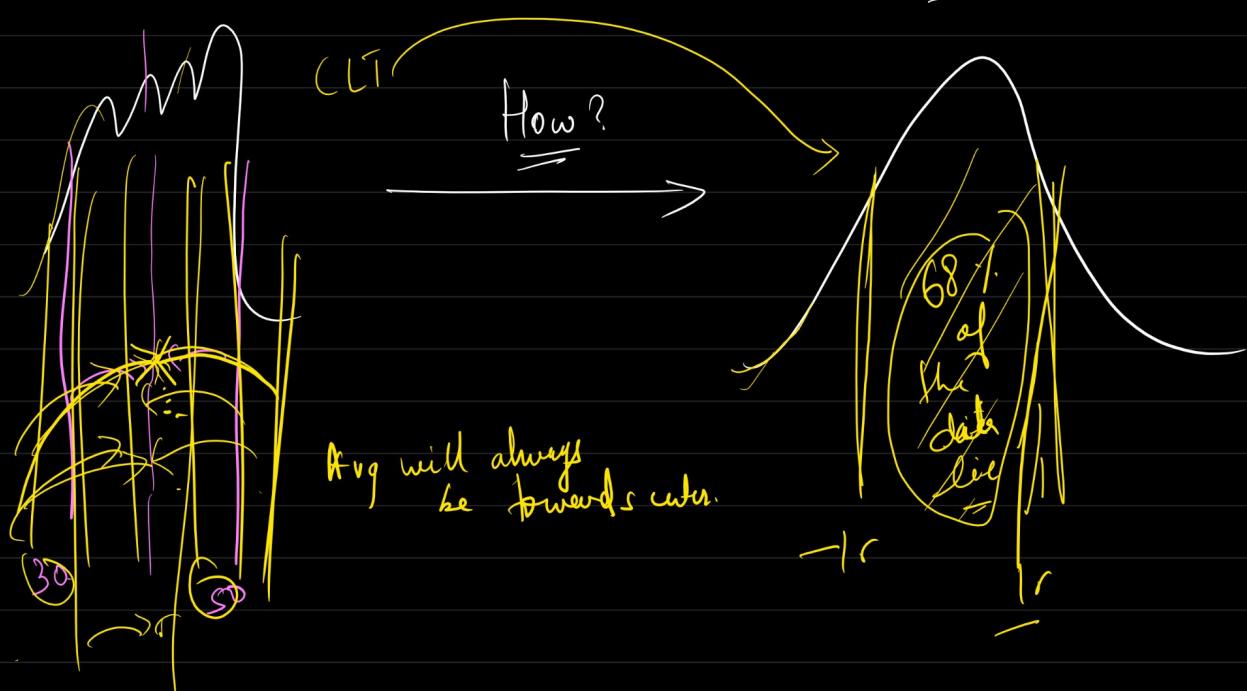
Central limit theorem \rightarrow A way to convert any form of distⁿ to Normal $\underline{d^n}$



* Two condition of CLT

\rightarrow The no of Samples Should be large enough

\rightarrow The Sample size > 30



STD

$$Z_{\text{sur}} = \frac{\bar{x} - \mu}{\sigma}$$

CLT

$$Z_{\text{sur}} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$\underline{S.S} < 30 \rightarrow \underline{t \text{ distn}}$$

Q You have a population
with $\mu = 100$

& $\sigma = 20$, if you have
 $S.S$ of 50 from this pop. what is

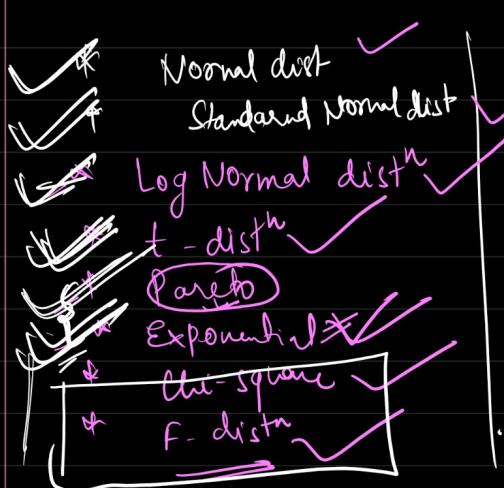
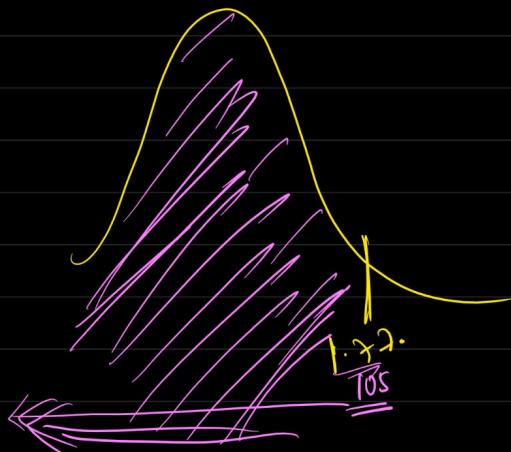
finds that sample mean will be less than
105?



$$\underline{\mu = 100} \quad \underline{\sigma = 20} \quad \underline{n = 50} \quad \underline{\bar{x} = 105}$$

$$Z_{\text{std}} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{105 - 100}{20 / \sqrt{50}} = \frac{5\sqrt{2}}{4} = 1.765$$

$$\begin{aligned} Z_{\text{dist}} &\neq Z_{\text{table}} \\ &= 1.77 \end{aligned}$$



① t distn

→ Shape is determined by dof

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

σ of popn
is given

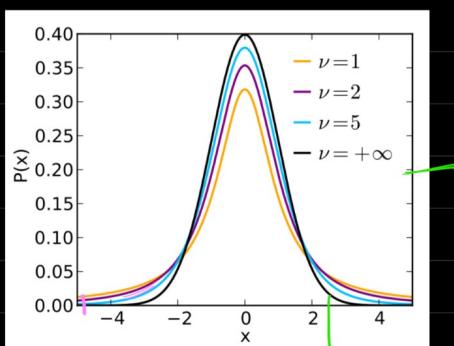
What if σ_{pop} not
given or $S.S < 30$

$$\underline{\sigma^2} \geq 30$$

↓
t-distribution

$$t \text{ score} = \frac{\bar{x} - \mu}{S / \sqrt{n}}$$

$S \rightarrow$ Sample stdn



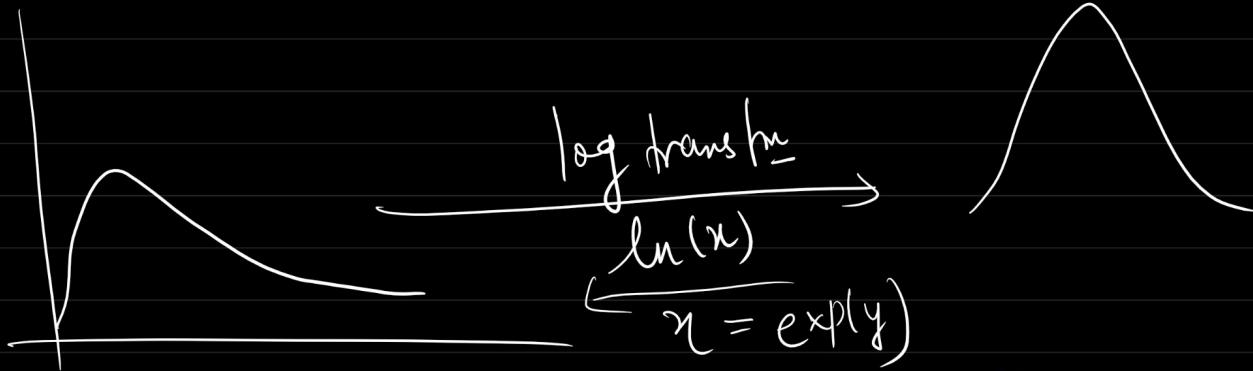
$$dof = 5 \cdot 5 - 1$$

+ dist - tails becomes thicker

shape - becomes thinner

② log Normal distn

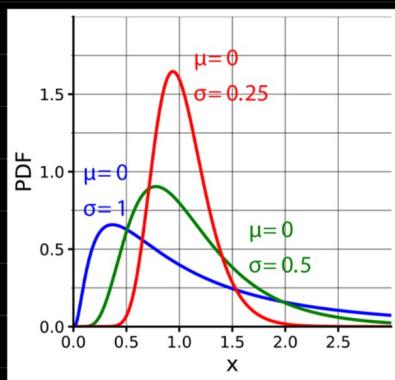
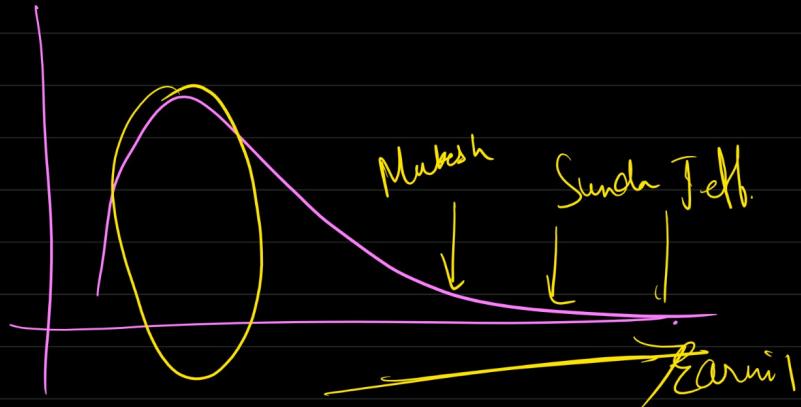
↳ A contin' prob distn whose logarithm is normally distributed.



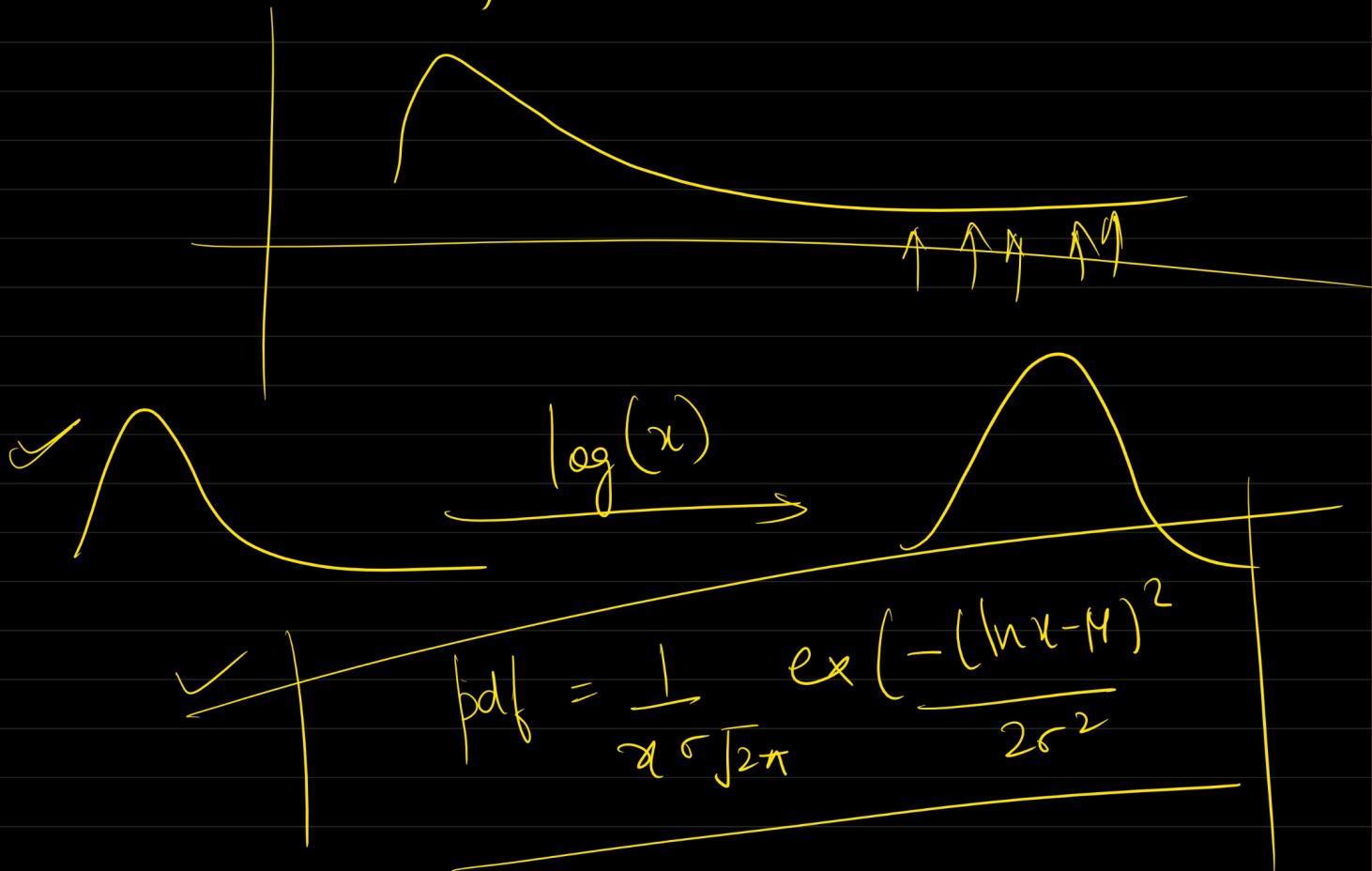
logarithm distn to normal distn

Example

Wealth distn of world



* People writing comments in Post



* Power law distn (Pareto principle)

RCB

20%



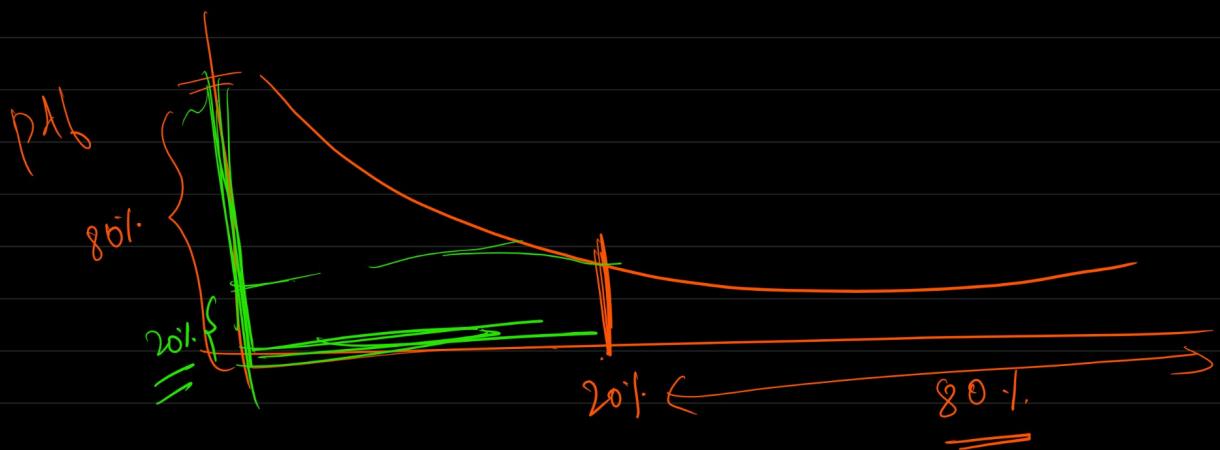
20% of the
people do
the 80% of the

→ 20% of employees
in a team
does 80% of work

→ 20%

80 | 20
form

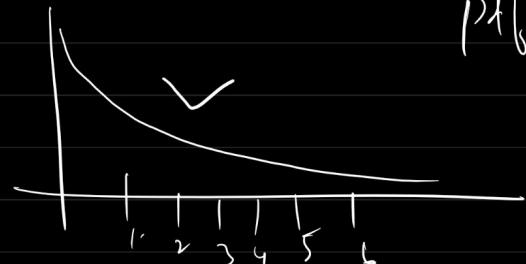
K G F
Kohli Gly
Fabs



$$Pf = \frac{\lambda x^{\lambda}}{x^{\lambda+1}}$$

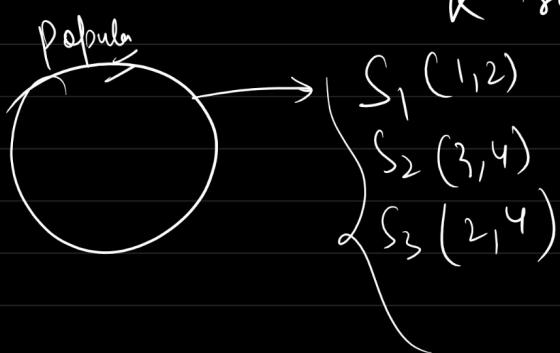
* Exponential distⁿ

- Radio active decay
- life span of battery

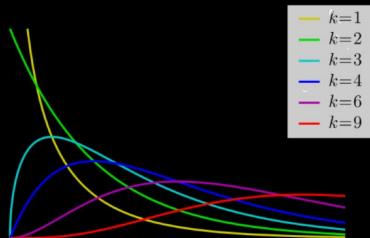
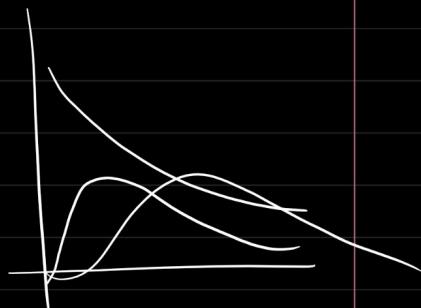


* Chi-square distⁿ

→ describe the distⁿ of Sum of Square of k random variable.



Square the sum of
$1^2 + 2^2$ add
$3^2 + 4^2$
$2^2 + 4^2$

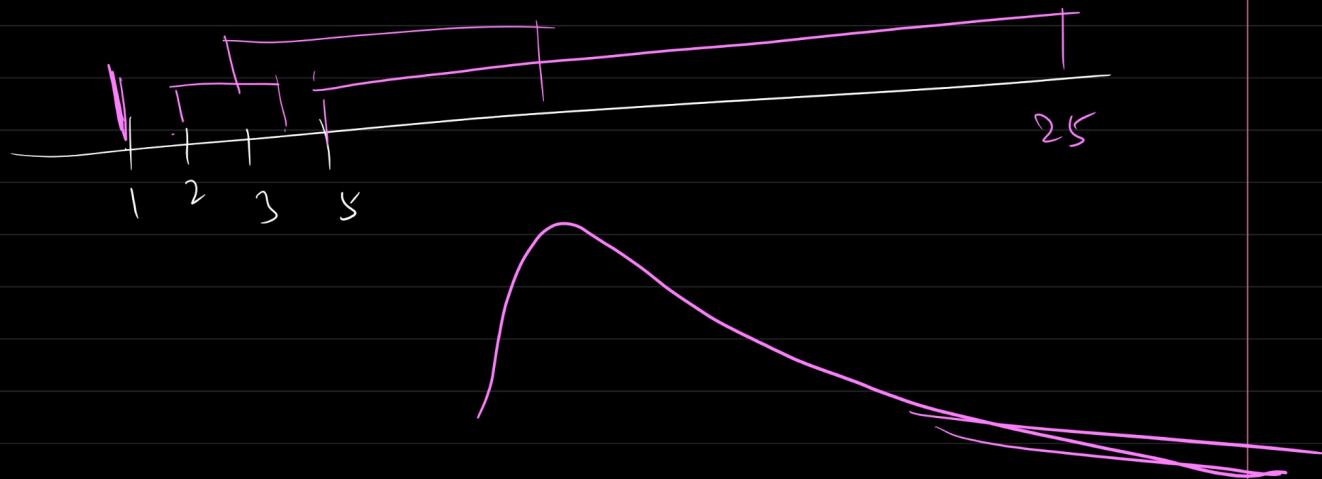


→ determined by k
 → Non-neg disⁿ
 → right skewed dist
 for small k - Pareto pr

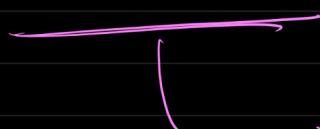
$$(0.2)^2 = \underline{\underline{0.04}}$$

$$2^2 - 4$$

$$3^2 - 9$$



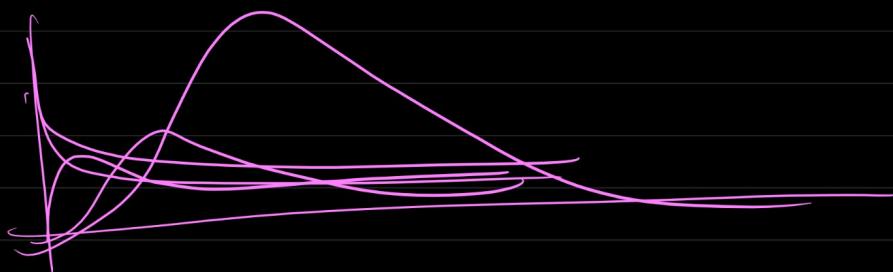
* F distn

 Compares Variances

$d_1 d_2, s_1 s_2$

Fisher Sneddr

$$X = \frac{s_1 | d_1}{\underline{\underline{s_2 | d_2}}}$$



$$dof = \text{Sample size} - 1$$

f dist → with change in No of Sample,
the shape changes

Bernoulli $\begin{cases} \text{Pass} \\ \text{fail.} \end{cases}$ $P^k (1-P)^{1-k}$

n no. of Bernoulli $\xrightarrow{\quad}$ Binomial distn

$$\xrightarrow{\quad} {}^n C_k P^k (1-P)^{n-k}$$