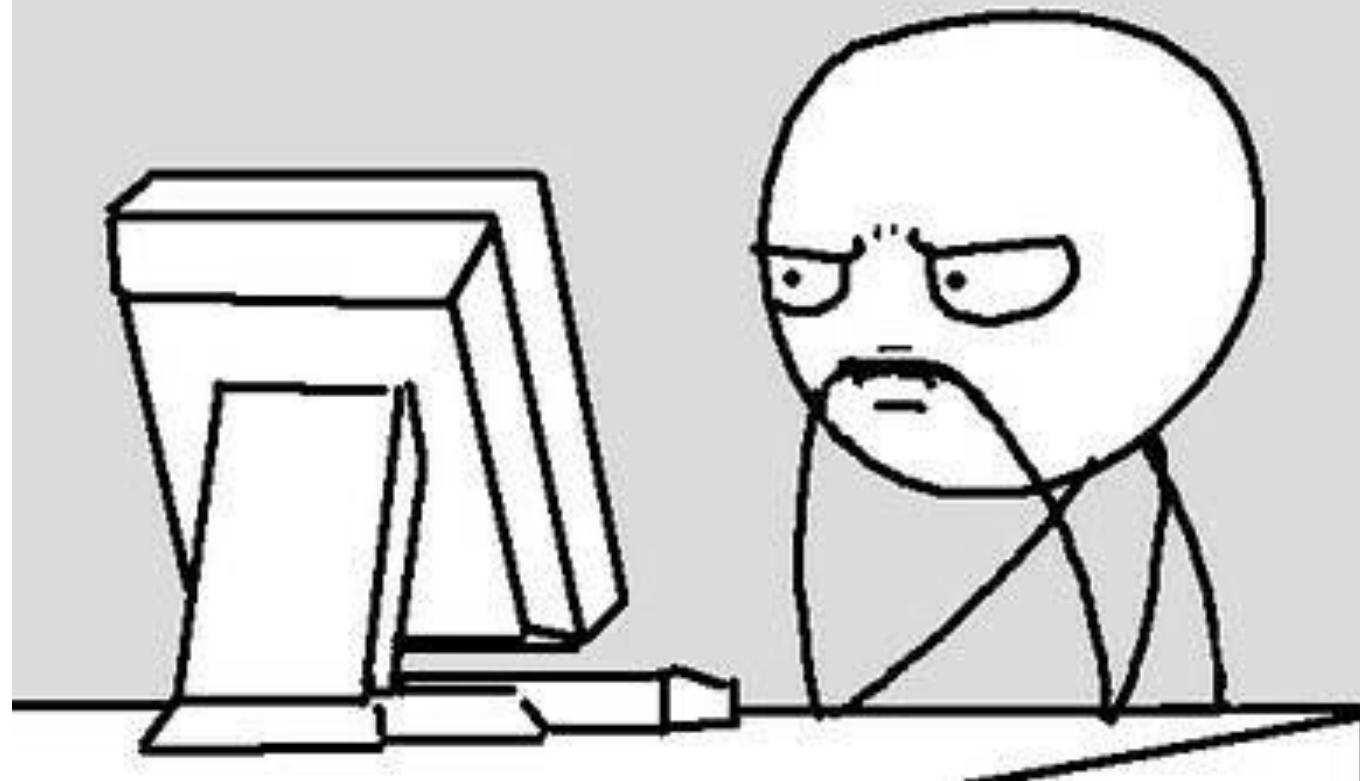


# LET'S WAIT

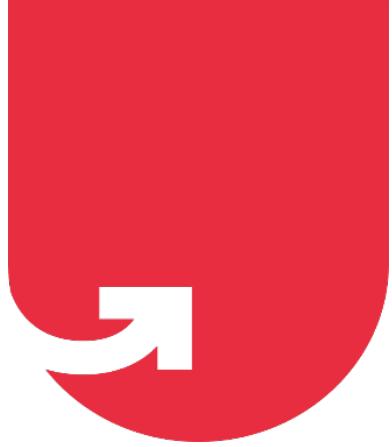


**Course : Machine Learning**

**Lecture On :Inferential  
Statistics**

**Instructor : Shivam Garg**





#LifeKoKaroLift

# Data Science Certification Program

# Today's Agenda

- ① Random Variable
- 2 Distributions – Normal & Binomial
- 3 { Inferential Statistics      *examples*
- 4 Central Limit Theorem
- 5 Doubt Session      {

*Expected Value*

# Random Variable: Numerical formulation

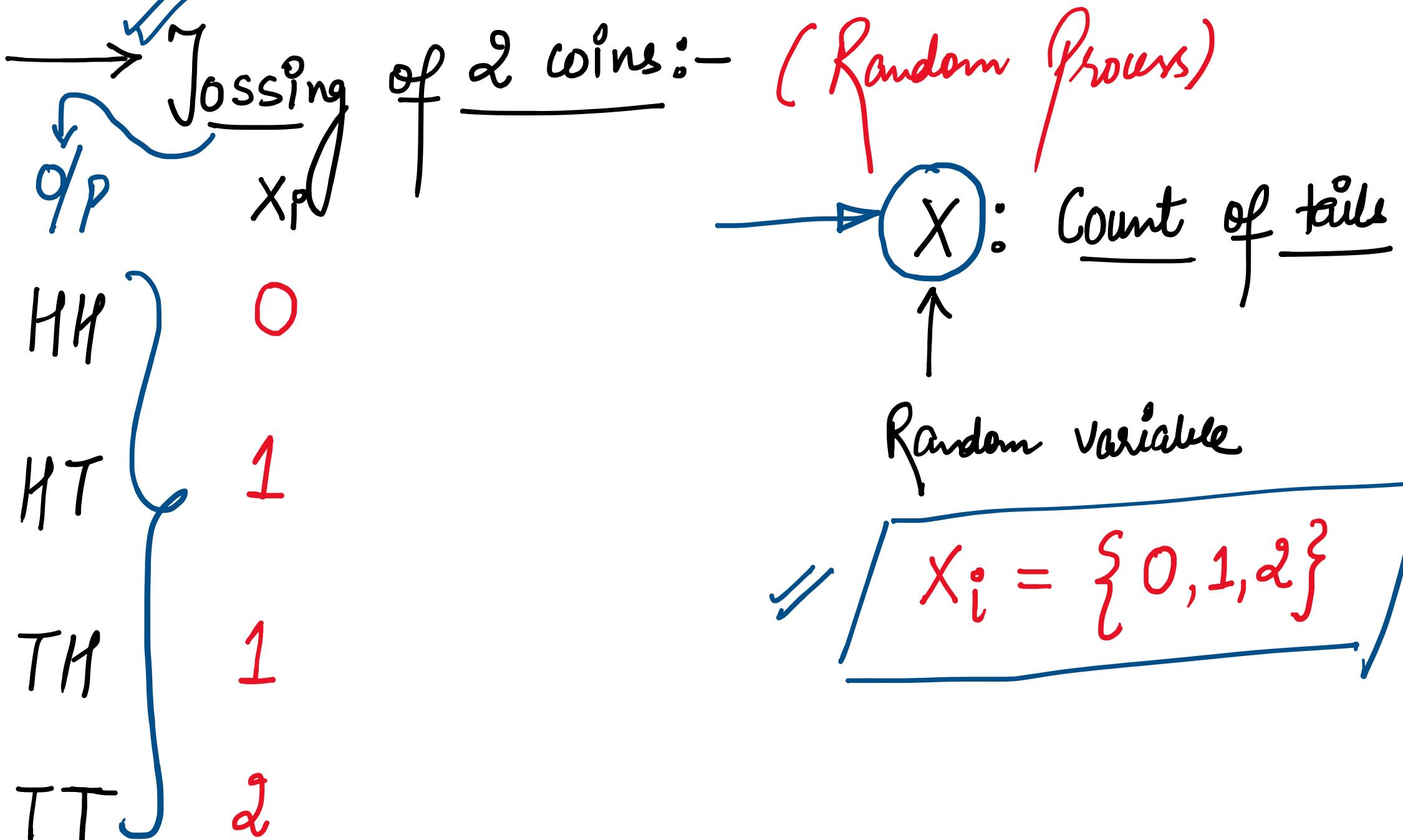
upGrad

\* It is a variable which maps the output of random process to some numerical values.

→ Random process: An event whose O/p is completely random.

Ex: ① Tossing a coin  $\{H, T\}$

② Rolling a die  $\{1, 2, 3, 4, 5, 6\}$



## Random Variable:

upGrad

{ Expected Value :- Average Value of random variable if the random process is repeated multiple times.

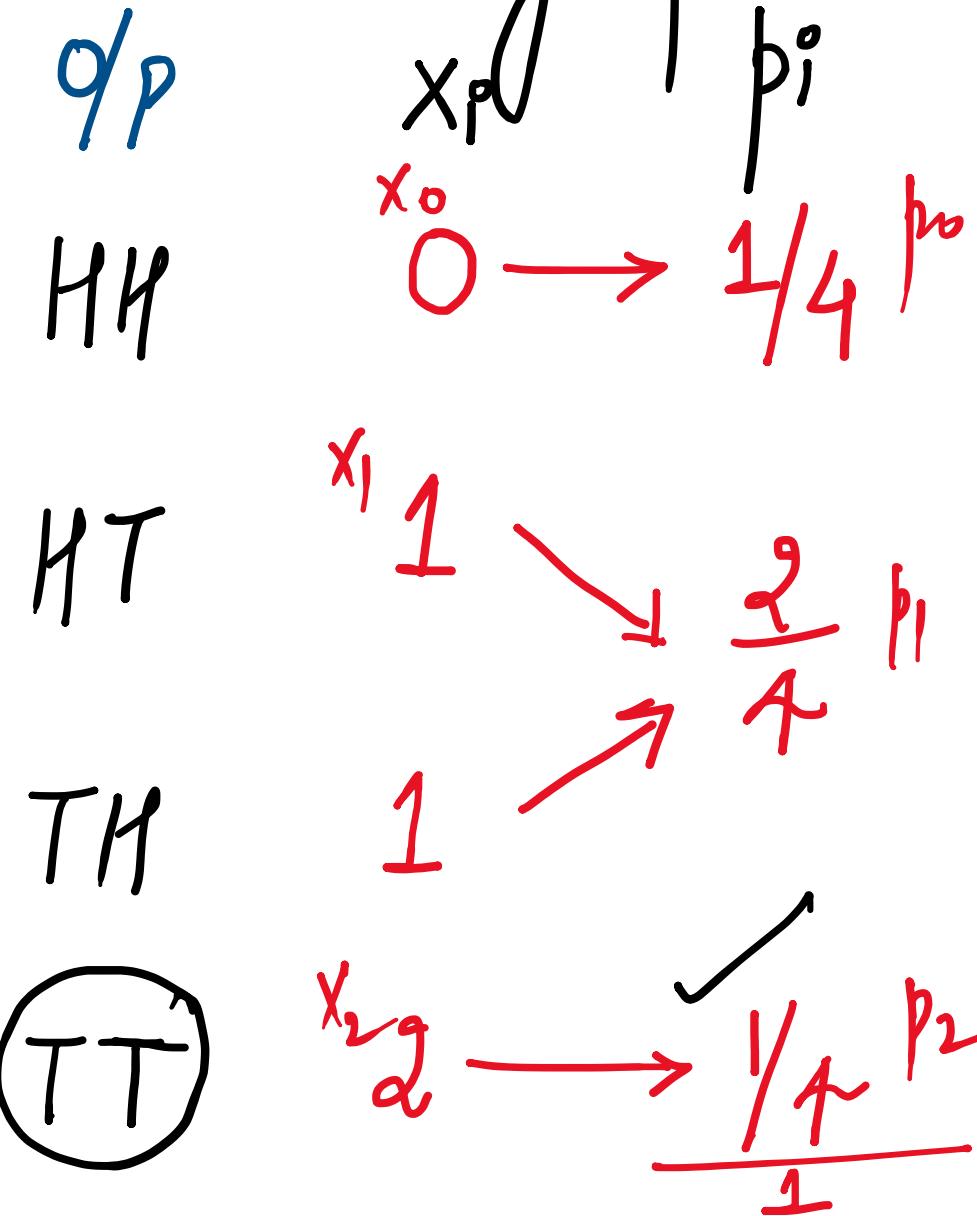
Mathematically,

$$E[X] = \sum x_i p_i$$

$x_i$  = Random Variable }  
 $p_i$  = Corresponding prob. }

$$E[X] = X_0 p_0 + X_1 p_1 + X_2 p_2 + \dots$$

→ Tossing of 2 coins :- (Random Process)  $X$ : Count of tails

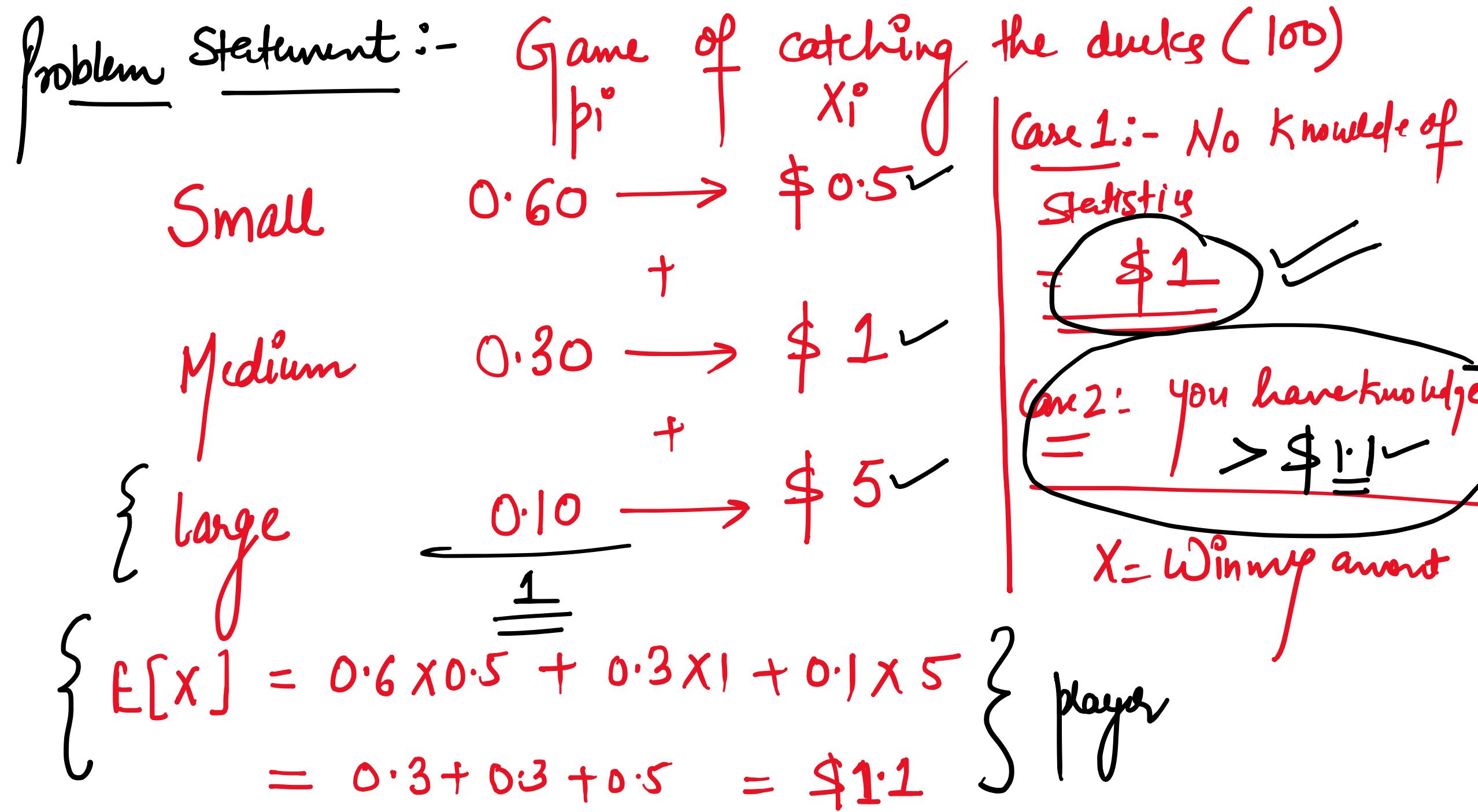


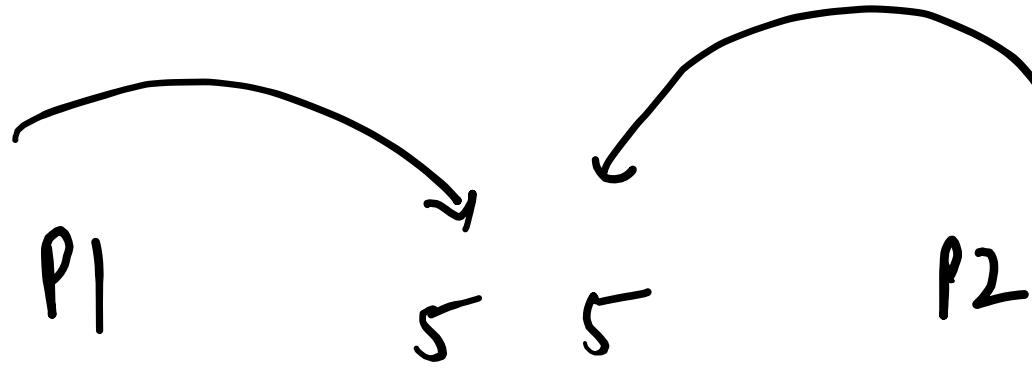
$$E[X] = \sum x_i p_i = x_0 p_0 + x_1 p_1 + x_2 p_2$$

$$= 0 \times \frac{1}{4} + 1 \times \frac{2}{4} + 2 \times \frac{1}{4}$$

$$E[X] = 1$$

If you toss 2 coins on an average you will be getting one tail.





50%

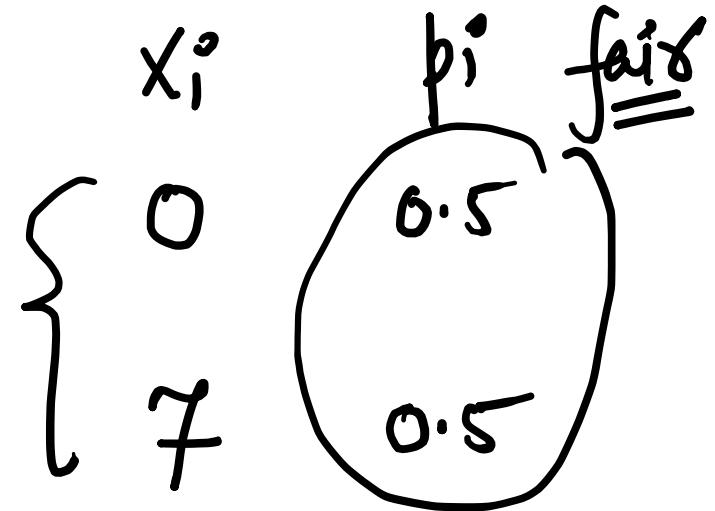


50%

$$E[X] = 0.5 \times 0 + 7 \times 0.5$$

$$= \underline{\text{₹}3.5} - \underline{\text{₹}5}$$

$$E[X] = -\underline{\text{₹}1.5}$$



Distributions

Continuous (Numerical)

Normal Dist

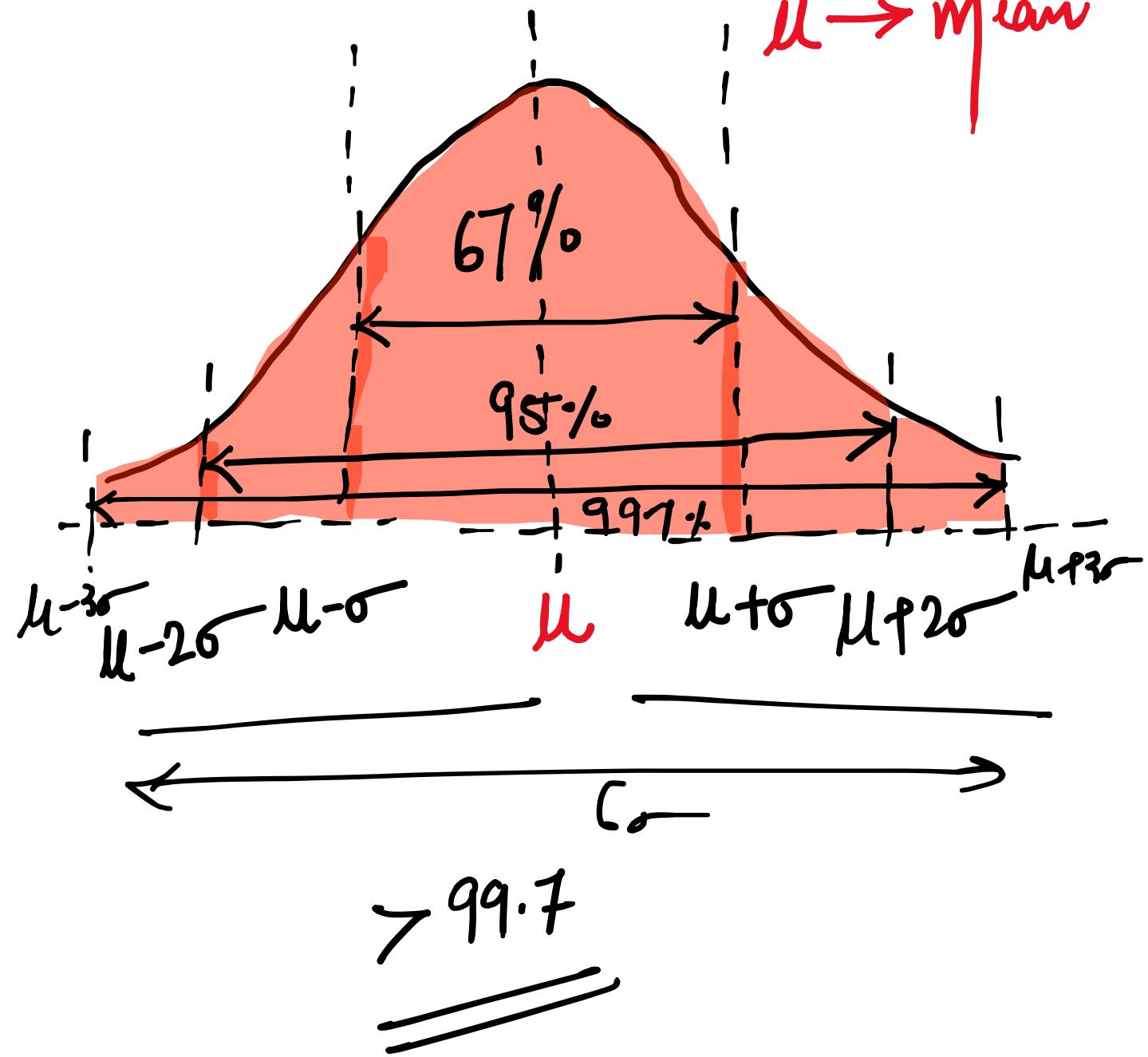
Categorical (Discrete)

Binomial Dist

→ Normal Dist :-

- \* Bell shape curve
- \* Symmetry at centre

$\sigma \rightarrow$  Std (Standard deviation)  
(sq root of variance)



## → Binomial Distribut<sup>n</sup> :- (Discrete distribut<sup>n</sup>)

- \* There should be only 2 categories (binary)  $nC_r = \frac{n!}{r!(n-r)!}$
- \* No. of trials should be known.
- \* All trials should be identical to each other.

$$\rightarrow P(X=r) = {}^n C_r p^r (1-p)^{n-r}$$

$n \rightarrow$  no. of trials  
 $r \rightarrow$  value of random variable  
 $p \rightarrow$  probability

→ Prob. of getting 2 tails while tossing 2 coins.

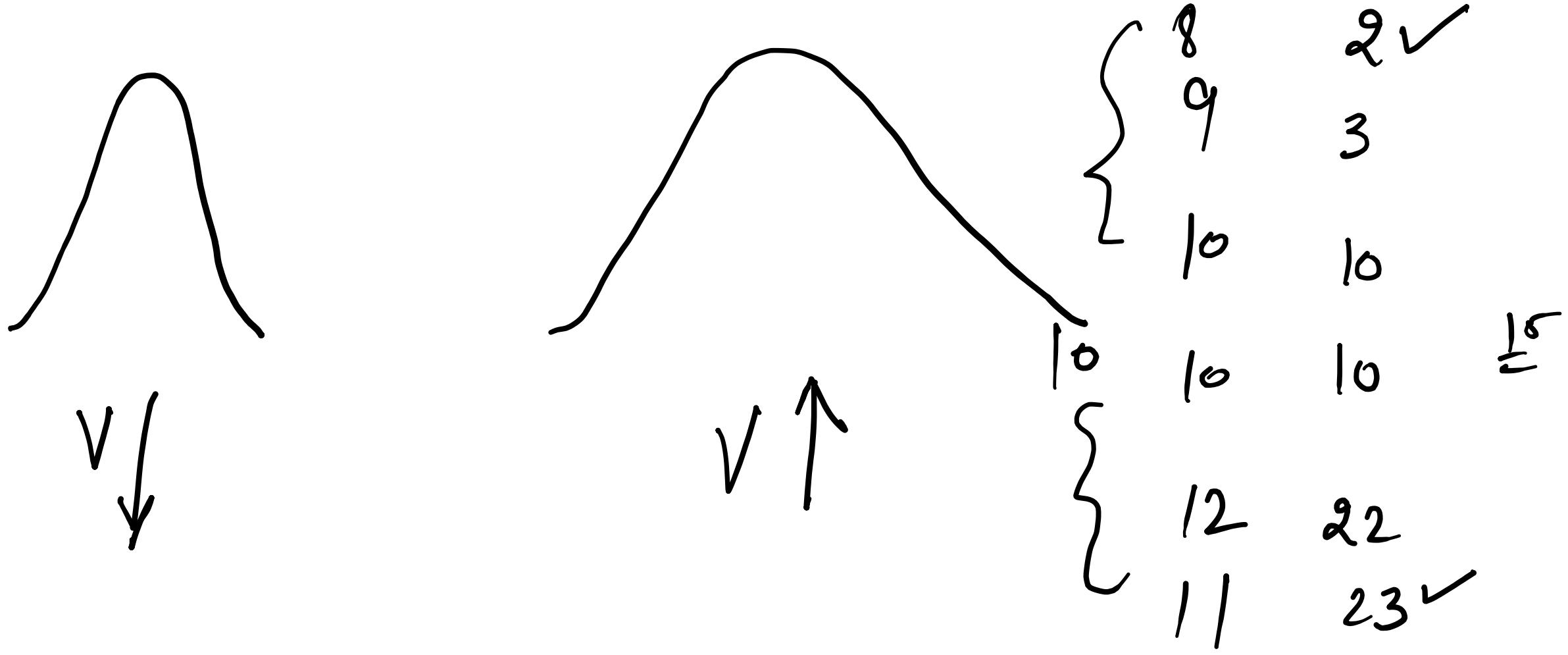
$$n = 2$$

$$r = 2$$

$$p = \frac{1}{2}$$

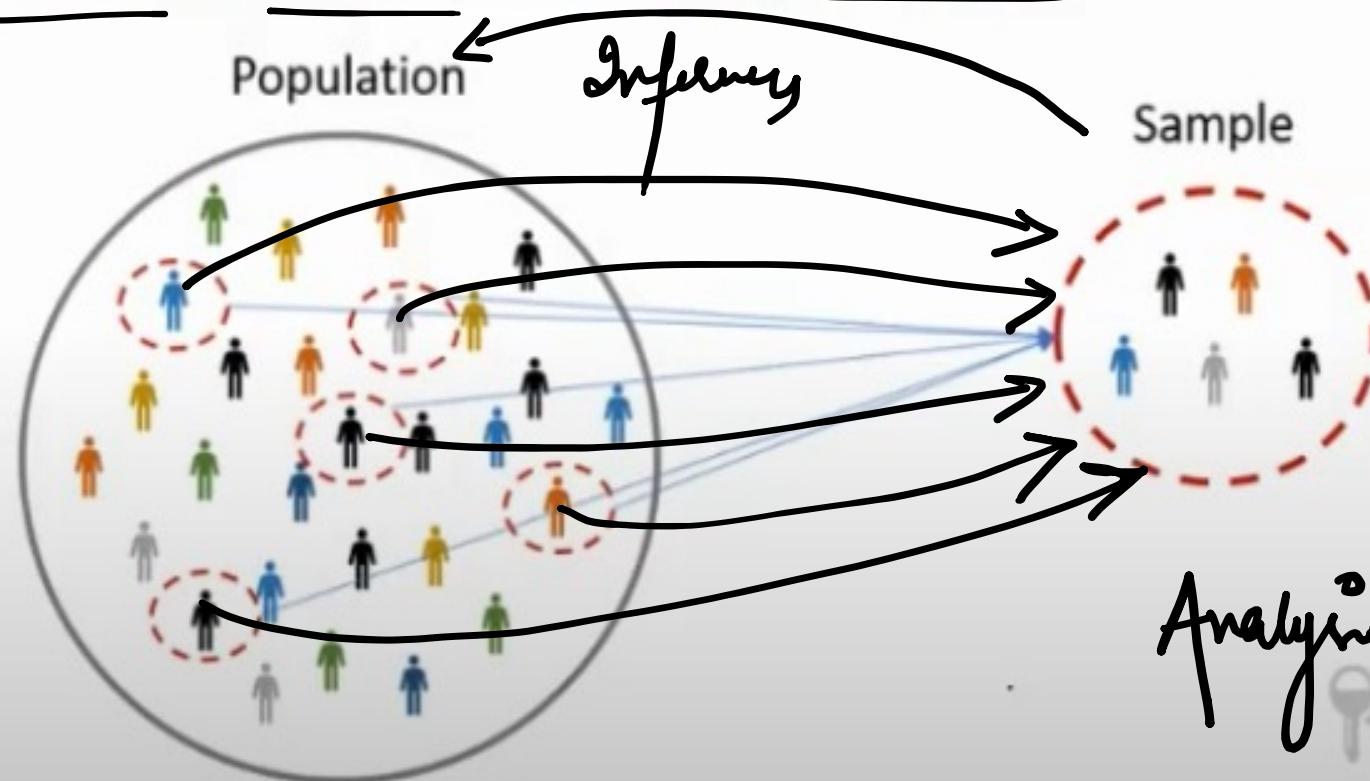
$$P(X=2) = {}^2C_2 \left(\frac{1}{2}\right)^2 \left(1 - \frac{1}{2}\right)^{2-2}$$

$$= \cancel{\frac{2!}{2!}} \times \frac{1}{4} \times 1 = \underline{\underline{\frac{1}{4}}}$$



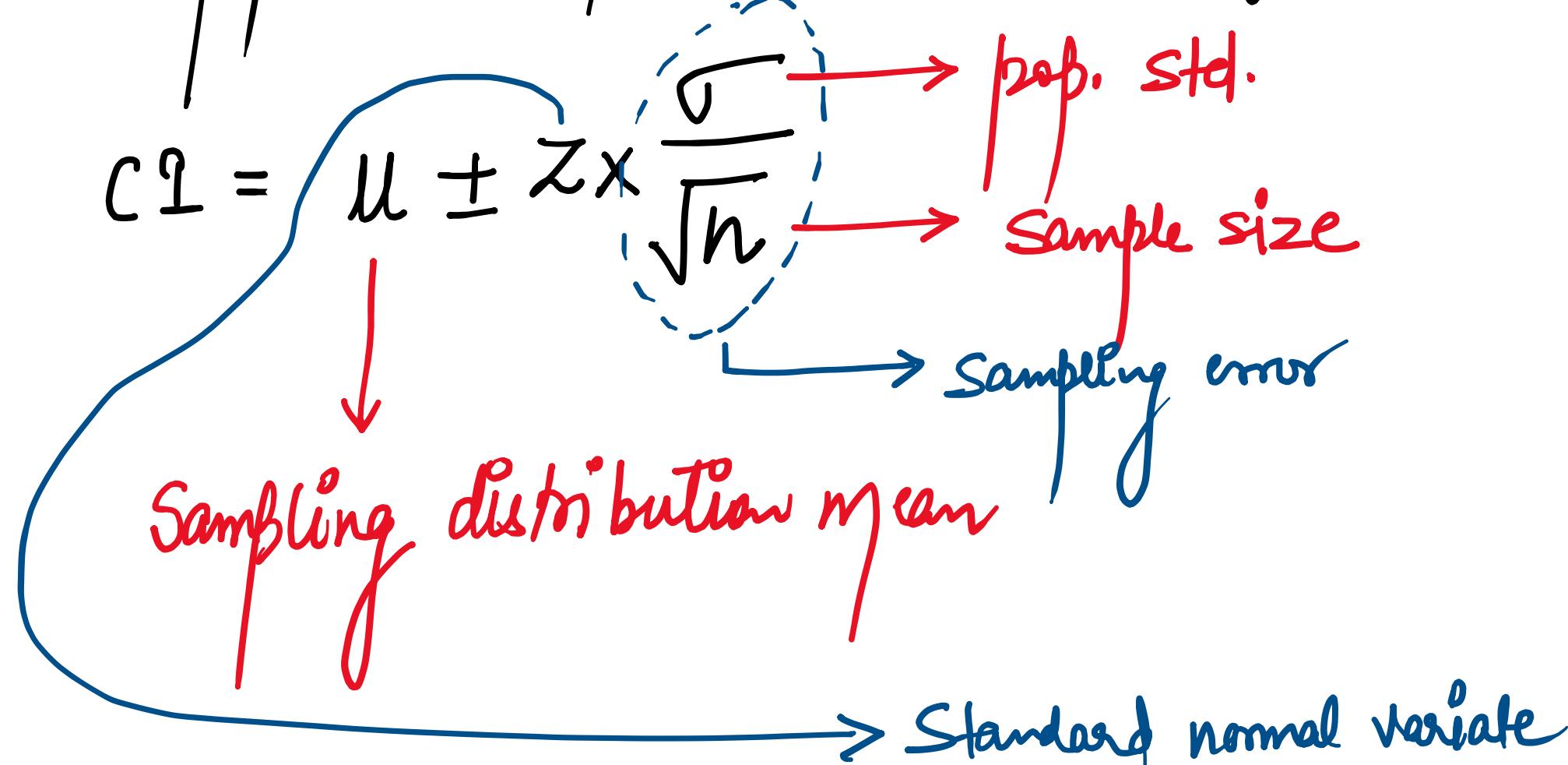


**The Problem:** Let's say that, for a business application, you want to find out the average number of times people in urban India visited malls last year. That's 400 million (40 crore) people! You can't possibly go and ask every single person how many times they visited the mall. That's a costly and time-consuming process. How can you reduce the time and money spent on finding this number?



Activate Windows  
Go to PC settings to activate Windows.

→ Population mean will lie in confidence interval of



$Z \rightarrow$  confidence level

$\text{CL}$

$\underline{\underline{90\%}}$

$Z$

$1.65$

$\underline{\underline{95\%}}$

$\underline{\underline{1.96}}$

$\underline{\underline{99\%}}$

$2.56$

{ by default }

~~hypo~~  $\rightarrow$  ~~Why?~~

Z table

$$\text{Value} = x + \frac{(100-x)}{2}$$

$$= 90 + \frac{10}{2} = \underline{\underline{95}}$$

$\underline{\underline{0.95}}$

# Central Limit Theorem: (*Assume → normal distribution*)

upGrad

A simple random sample of 50 adults women is obtained, and each person's red blood cell count (in cells per microliter) is measured. The sample mean is 4.63. The population standard deviation for red blood cell counts is 0.54. Construct the 95% confidence interval estimate for the mean red blood cell counts of adults.

$$z = 1.96$$

$$\mu = 4.63$$

$$\sigma = 0.54$$

$$n = 50$$

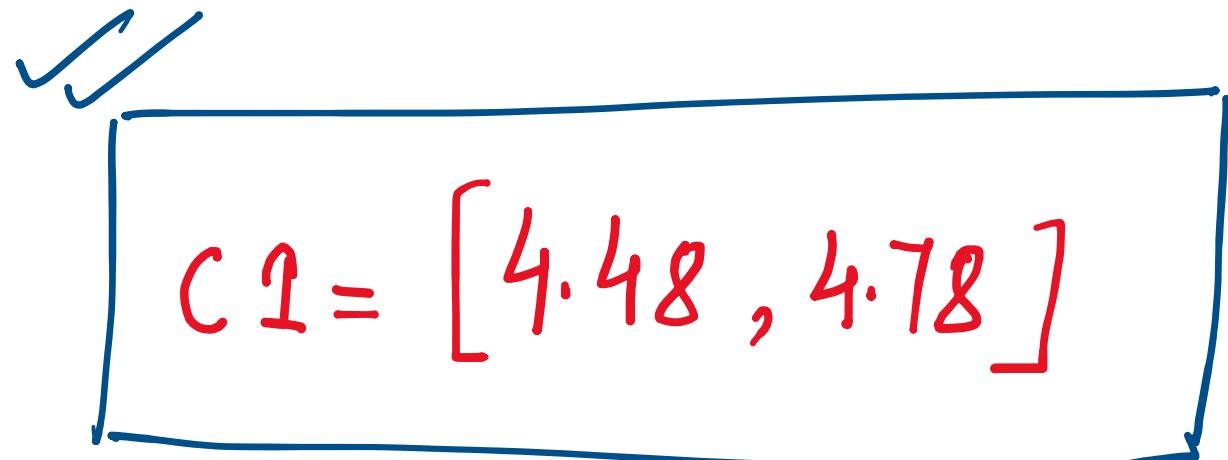
$$CI = \left[ \mu \pm z \times \frac{\sigma}{\sqrt{n}} \right]$$

$$= \left[ 4.63 \pm 1.96 \times \frac{0.54}{\sqrt{50}} \right]$$



Activate Windows  
Go to PC settings > Activation

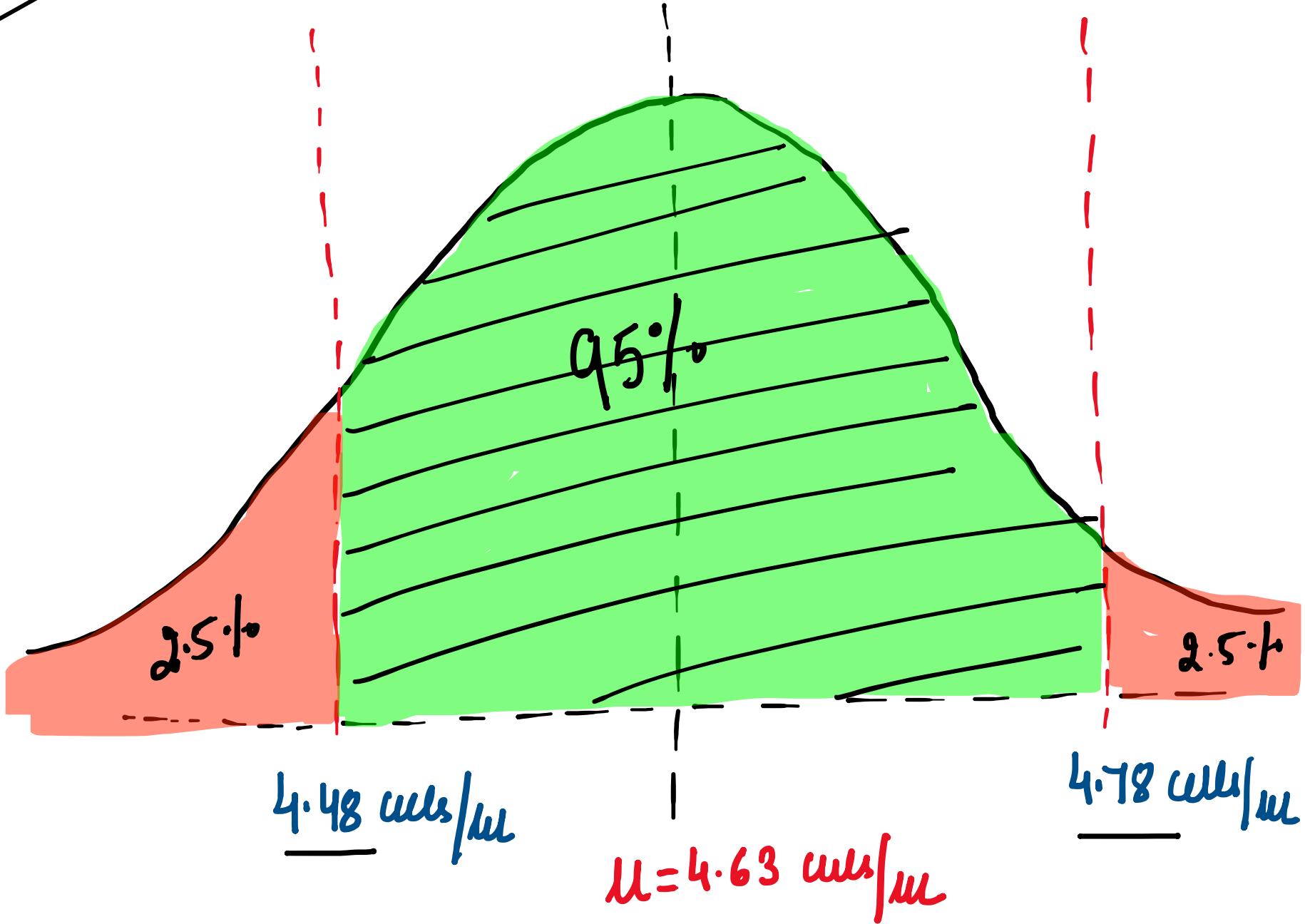
$$CI = 4.63 \pm .15$$



$$\{ C.I = [4.48, 4.78] \}$$

→ Population mean red blood cell count per micro liter  
will lie in the range of 4.48 cells/ $\mu$ l to 4.78 cells/ $\mu$ l

& we are 95% confident for the same.



## → Drawback of Inferential Statistics:-

\* Whatever inference is given as output of inferential statistics, is not validated anyway of the concept itself.

