

## Basic Statistics

### 1. Population and Samples

Research Question: Stress levels in teenagers

Population: All teenagers in the world

Sample: 200 teenagers

Quality of biscuits made in a factory

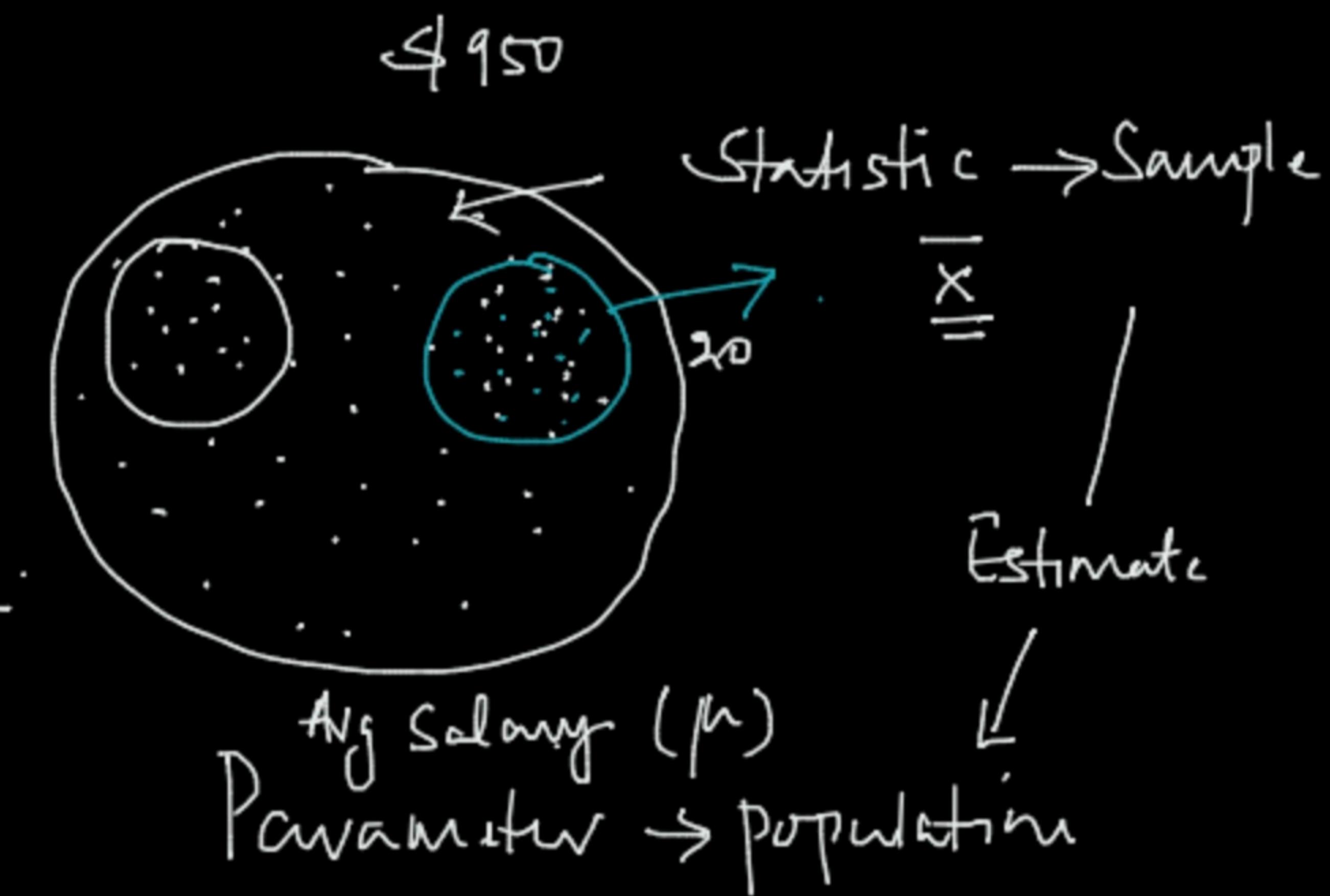
Population: All biscuits made in the factory

Sample: 100 biscuits

### Good Sample:

1. Representative → Represent the entire population

2. Random → Every datapoint should have equal chance of being selected



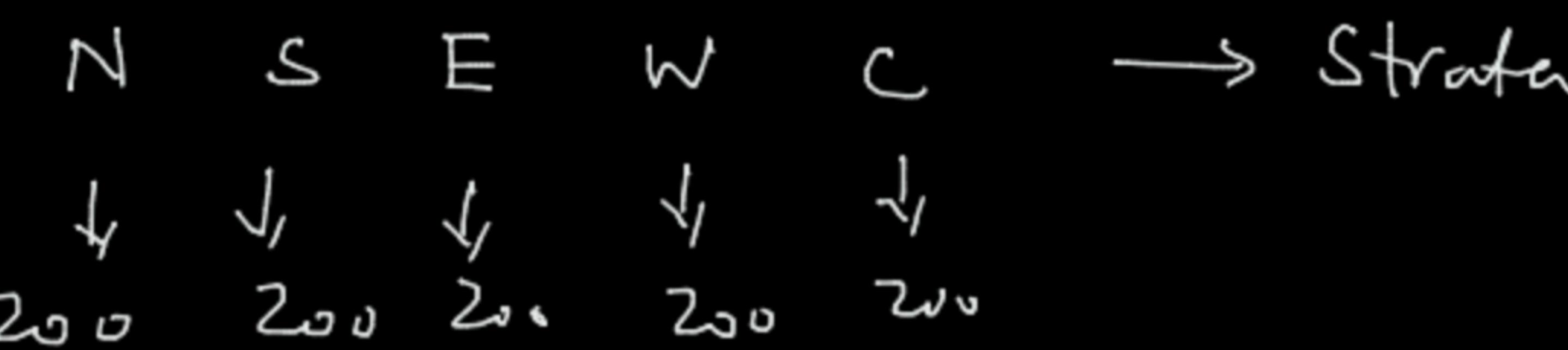
"Statistic" serves as an estimate for the

"Parameter"

# Food Habits of Indians

Population: All Indians

Sample: 100 → Chennai



## 1. Simple Random Sampling (SRS)

- Blind folded Selection

- When no info abt the population

Estonia → SRS

## 2. Stratified Sampling

- Some info avail abt. the population



## 3. Snowball Sampling

- sensitive info

- 

large Samples

# Types of Statistics

## Descriptive Statistics

### Measures of

- ✓ 1. Frequency ✓  
— Counts, Frequency / proportion
- ✓ 2. central tendency — Most Likely Value  
— Mean, Median, Mode
- ✓ 3. Dispersion / spread.  
— Range, Standard deviation / Variance
- ✓ 4. position  
— Percentiles, Deciles, Quartiles
- 5. Shape  
— Skewness, Kurtosis

## Inferential Statistics

— Draw some conclusions

— Hypothesis testing

Age	Course Under	Cntry
1	M	-
2	M	-
3	F	-
.	.	.
.	.	.
$n = 100$		Avg.

[10:45 am]

# 1. Arithmetic Mean

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$$

 $n = 5$ 

Median (E)

$$\bar{X} = \frac{1}{n} (x_1 + x_2 + x_3 + x_4 + x_5)$$

Marks (X)

$$x_1 = 45$$

$$x_2 = 70$$

$$x_3 = 54$$

$$x_4 = 66$$

$$x_5 = 75$$

$$\underline{\underline{310}}$$

$$\bar{X} = \boxed{62}$$

outliers affect the mean  
↓

Extreme Values

$$8 - DW = .8L$$

$$1 - CE = \cancel{45} \text{ (no)}$$

$$1 - CT = 80$$

$$\begin{array}{|c|} \hline 93 \\ \hline 90 \\ \hline 97 \\ \hline 95 \\ \hline \end{array}$$

$$\rightarrow \frac{5}{\underline{\underline{380}}} = \boxed{76.1} \rightarrow$$

Avg  $\Rightarrow 93.75$  ✓

Median  $\rightarrow 93.5$  ✓

$$64 + 165 + 80 \rightarrow \cancel{150} \text{ th } \frac{244}{10} \rightarrow 24.4$$

→ QSL

$$\leftarrow \quad \downarrow \quad \rightarrow$$

$$5 \quad 90 \quad \boxed{93} \quad 95 \quad 97$$

$$5 \quad 75 \quad 90 \quad \downarrow \quad 93 \quad 95 \quad 97$$

$$\leftarrow \quad \boxed{91.5} \quad \rightarrow$$

$$50.1 \quad 50.1$$

2. Median — unaffected by outliers

— arrange no. in order

— take central value

3. Mode  $\rightarrow$  Most frequent Value

70 M ←

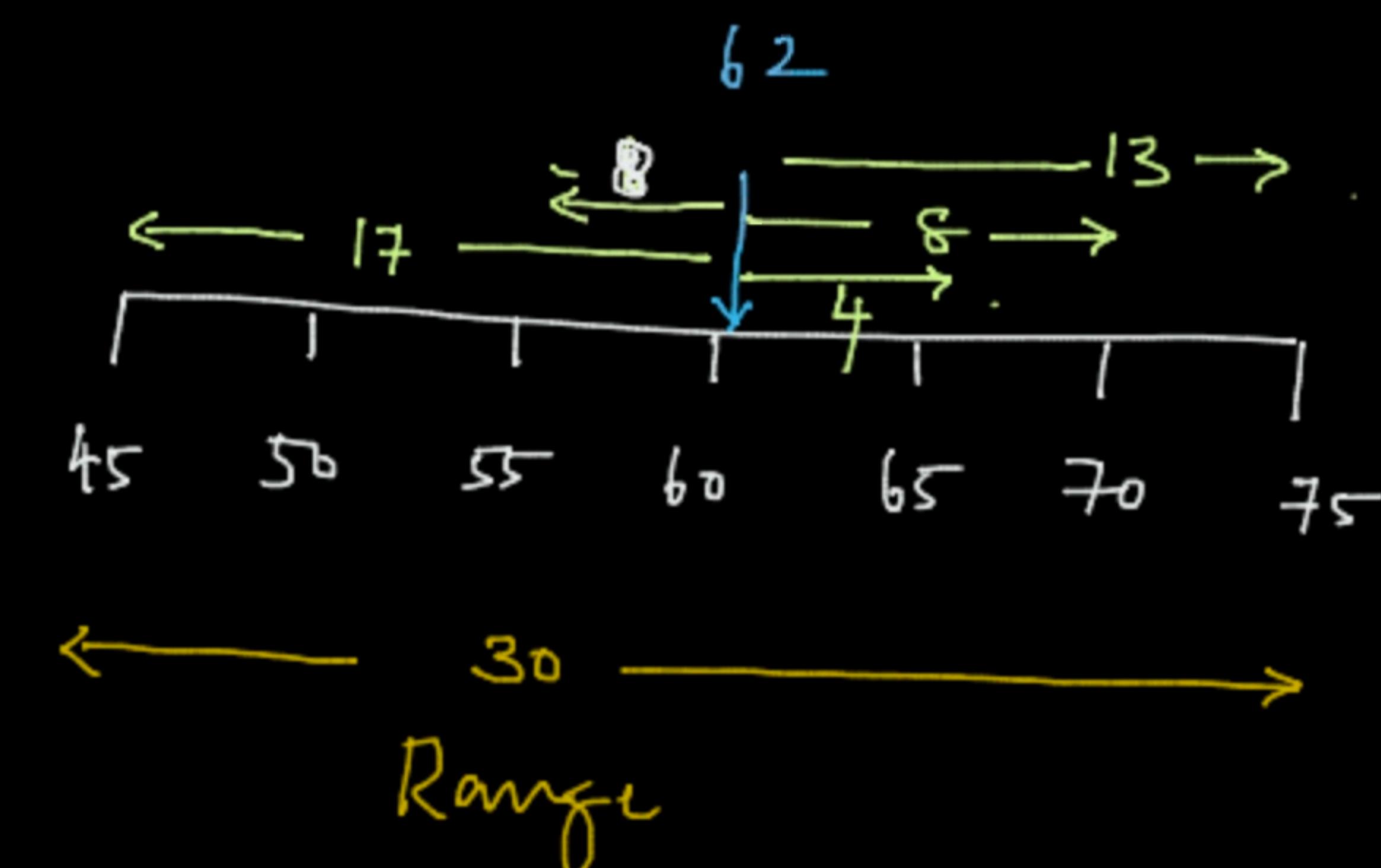
30 F

As	30
Am	30
Tn	30
La	15

1, 2, 2, 3, 4, 5  $\rightarrow$  2.

160  $\rightarrow$  37

$$\frac{2}{5} = 0.4 \times$$



Marks (x)	$(x - \bar{x})$	$(x - \bar{x})^2$
45	-17	289
70	+8	64
54	-8	64
66	+4	16
75	+13	169
$\bar{x} = 62$		$\frac{602}{5} = 120.4$

Measures of dispersion / spread.

$$\begin{array}{cc} \overline{\text{Class 1}} & \overline{\text{Class 2}} \\ \sqrt{50} & \checkmark \\ \sigma = 5 & \sigma = 30 \end{array}$$

a. Range:

$$= \text{Max} - \text{Min}$$

$$= 75 - 45$$

$$= 30$$

(b) Variance ( $\sigma^2$ )

$$= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

(c) Standard Deviation — on an avg how far are ✓

[the datapoints from the central Value -

$$\sigma = \sqrt{\text{Variance}}$$

$$\sigma = 11$$

Population

No. of Records

N

Avg.

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i$$

Variance

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

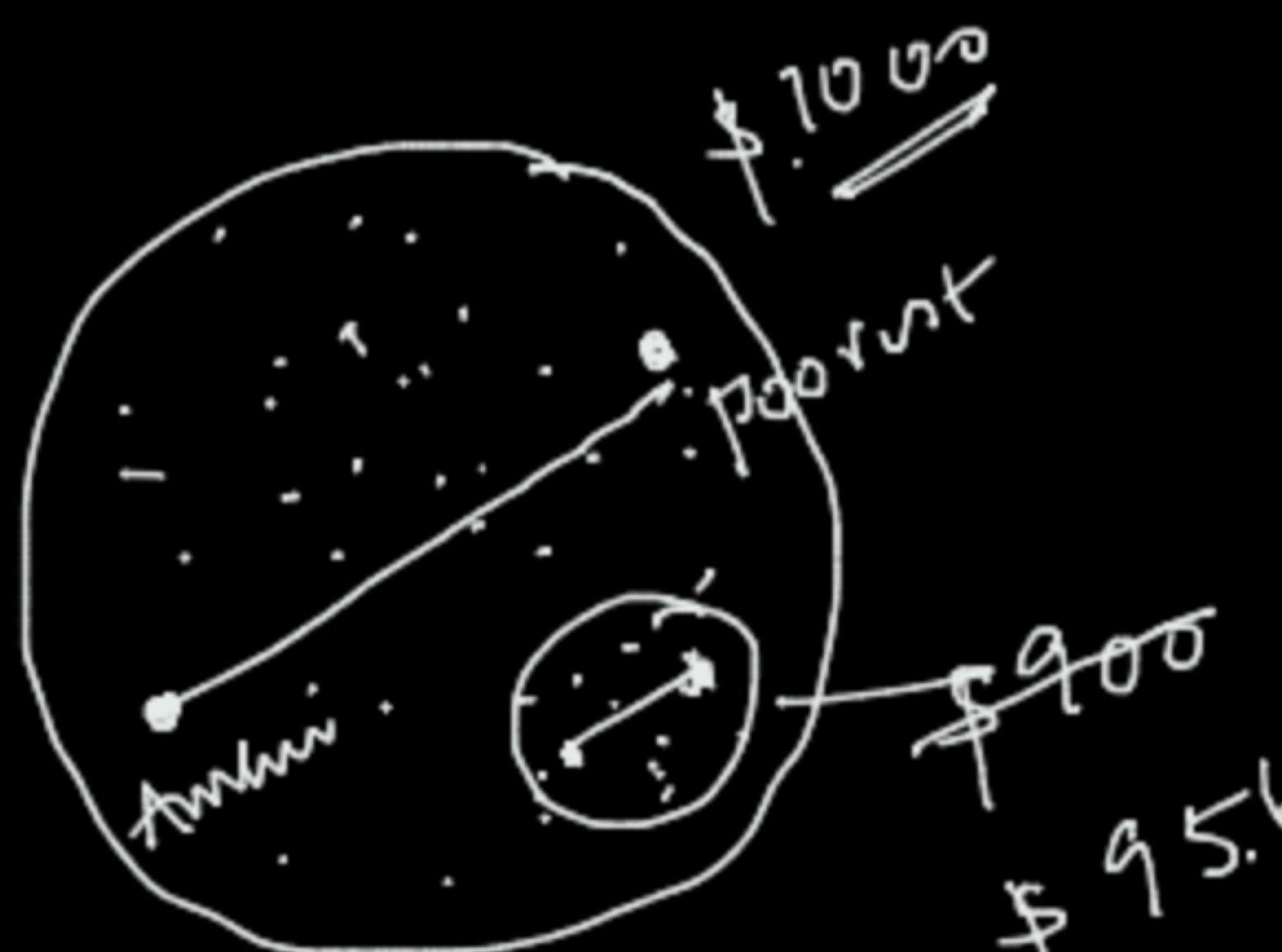
Sample

n

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

$$s^2 = \frac{1}{(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2$$

↳ Bessel's correction



Avg income  
( $\mu$ )

Sample Variance underestimates Population Variance

$n-1 \rightarrow$  Degrees of Freedom

$\rightarrow$  The no. of datapoints that are free to vary given a constraint

Constraint

$$\boxed{\text{Total} = 100} \quad \checkmark$$

$$n = 10$$

$$\boxed{\text{Total} = (9)}$$

$$\left\{ \begin{array}{l} p_1 = 0 \\ p_2 = 13 \\ p_3 = \\ \vdots \\ p_9 = 20 \rightarrow 85 \\ p_{10} = x \end{array} \right.$$

$$df = 9$$

$$\left[ \begin{array}{l} x_1 = 2 \checkmark \\ x_2 = 4 \checkmark \end{array} \right] \quad \frac{x_3}{2}$$

Variance  $\rightarrow$  Information

Area	Bed	Price
1500	2	50/-
1800	2	60/-
2200	2	75/-
2500	2	83/-

$$\hookrightarrow \text{Var} = 0$$

← Measures of position →

→  $\boxed{65}$

→  $90^{\text{th}}$  percentile

$85\% \rightarrow 40^{\text{th}}$  percentile



:

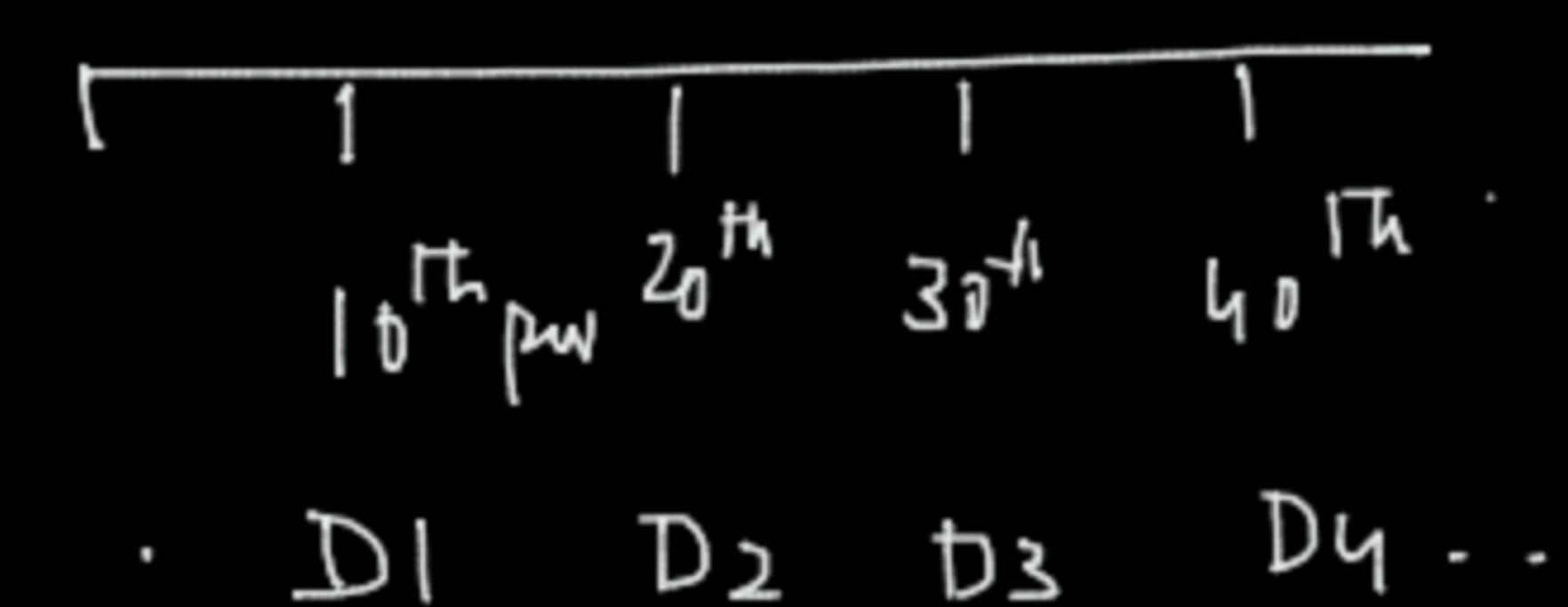
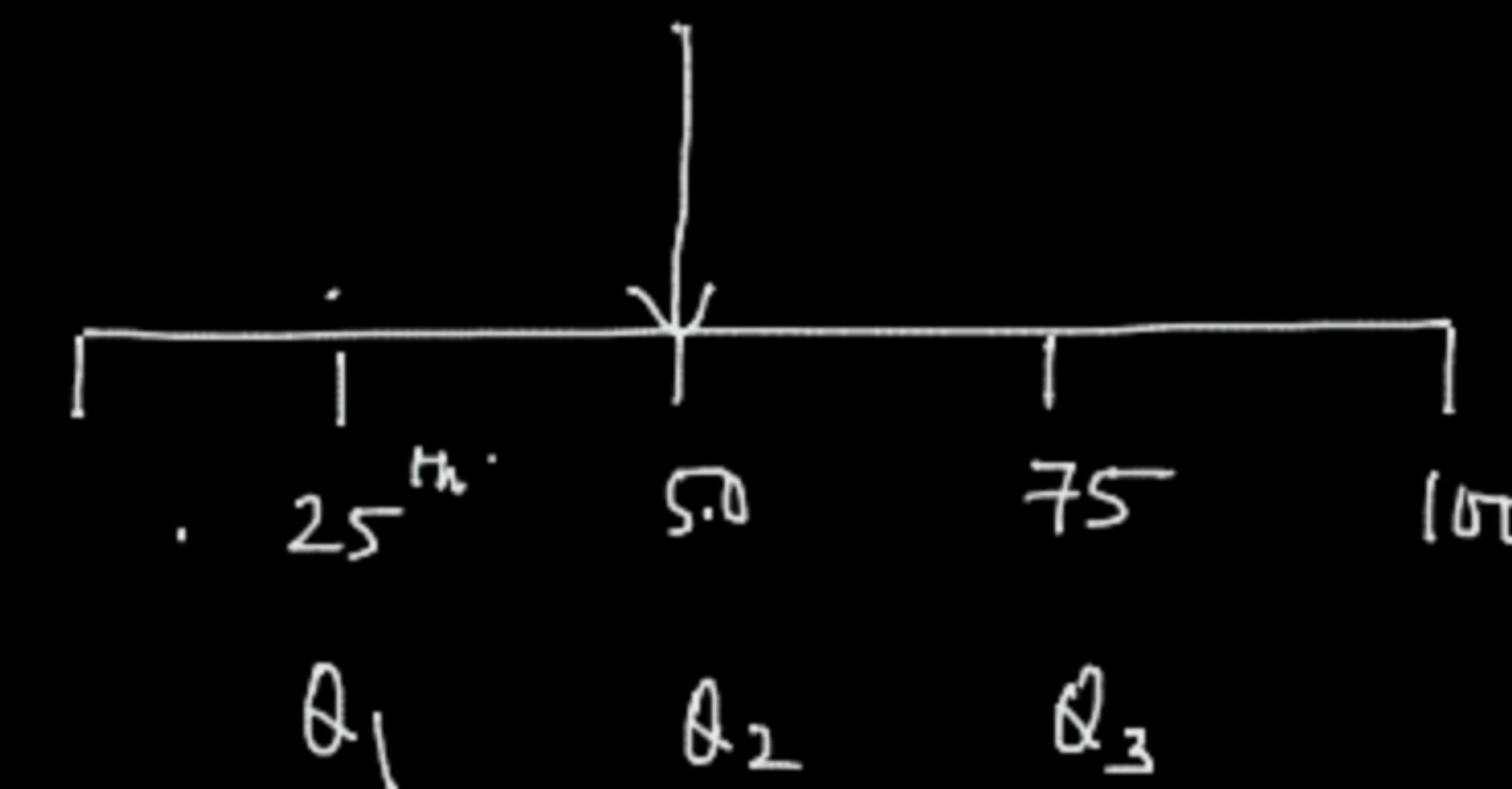
:

#L



Quartiles →  $25^{\text{th}}$  percentile →  $Q_1$

Deciles → 10<sup>th</sup> pw

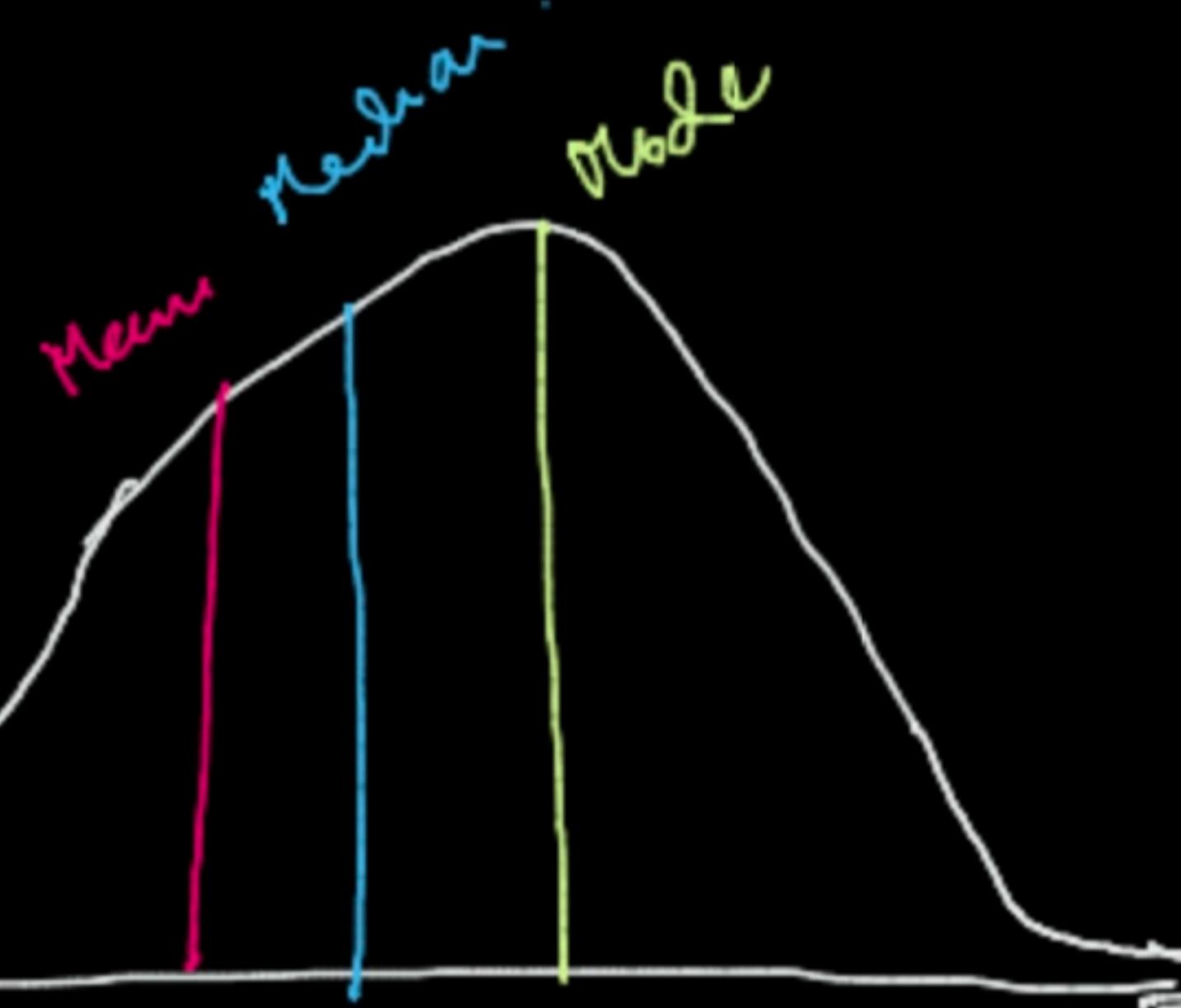


Quan**N**tile → Percentile  
 → Decile  
 → Qua**R**tile

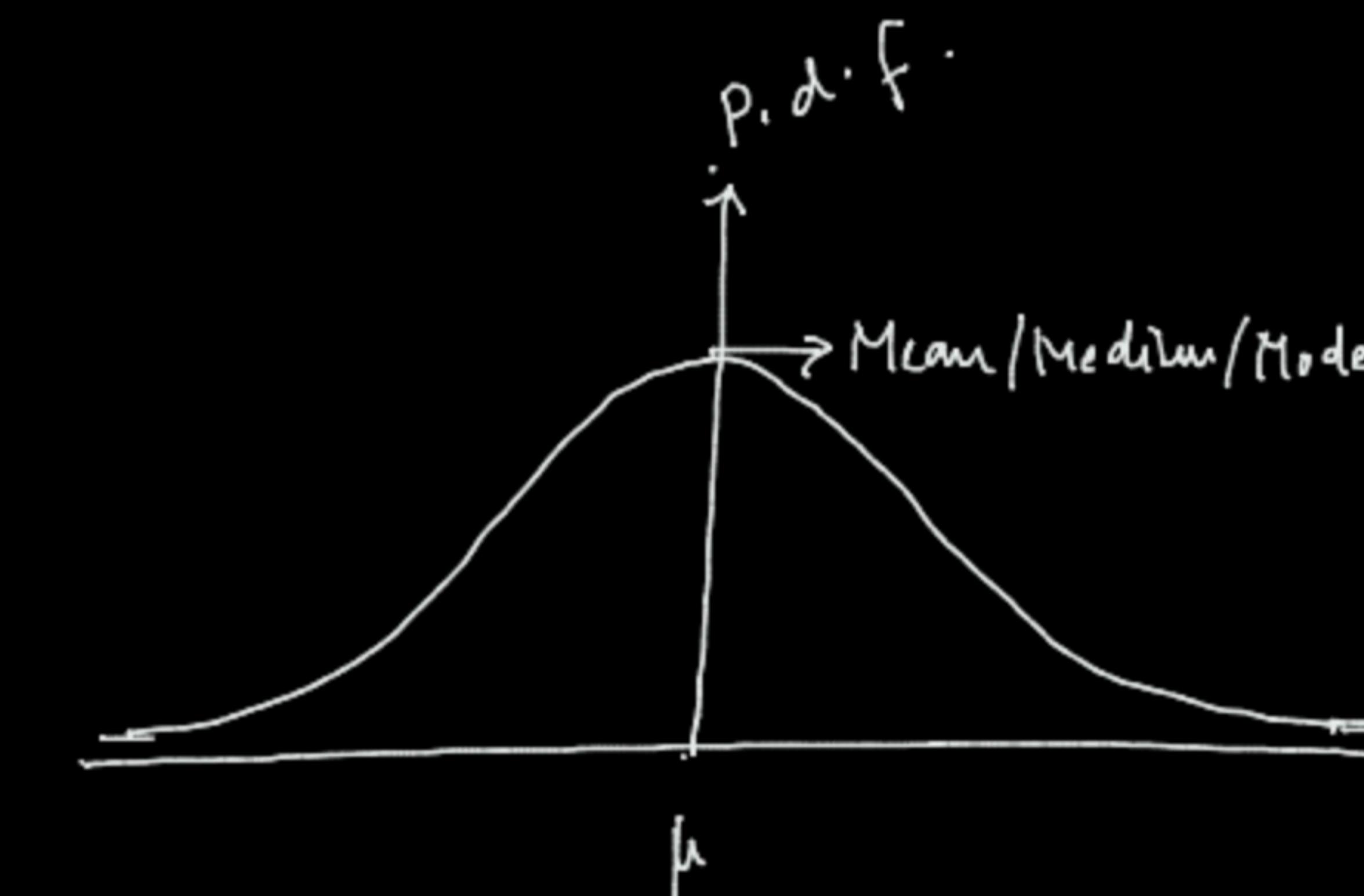
← Measures of Shape →

## (a) Skewness

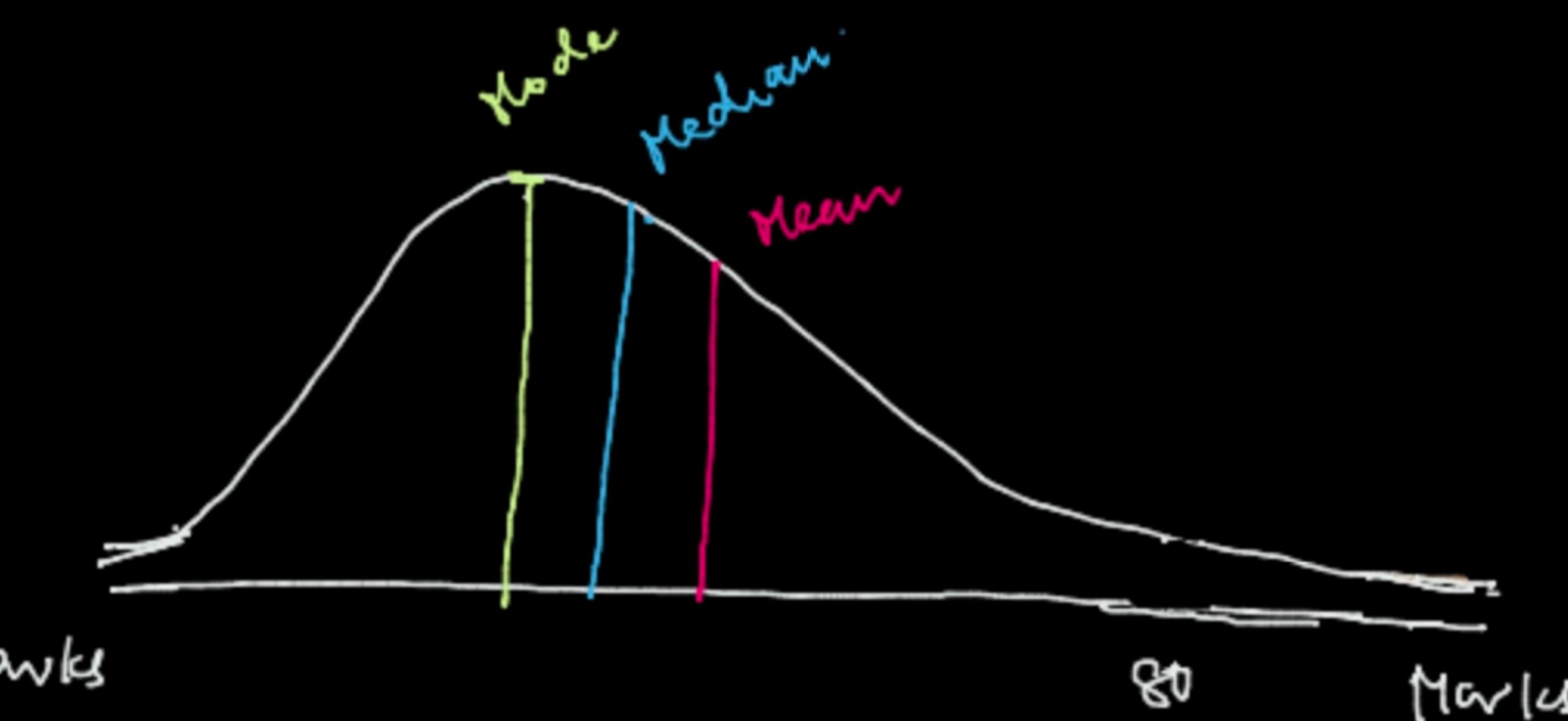
— How symmetric  
is the data dist.



Left skew (-ve)

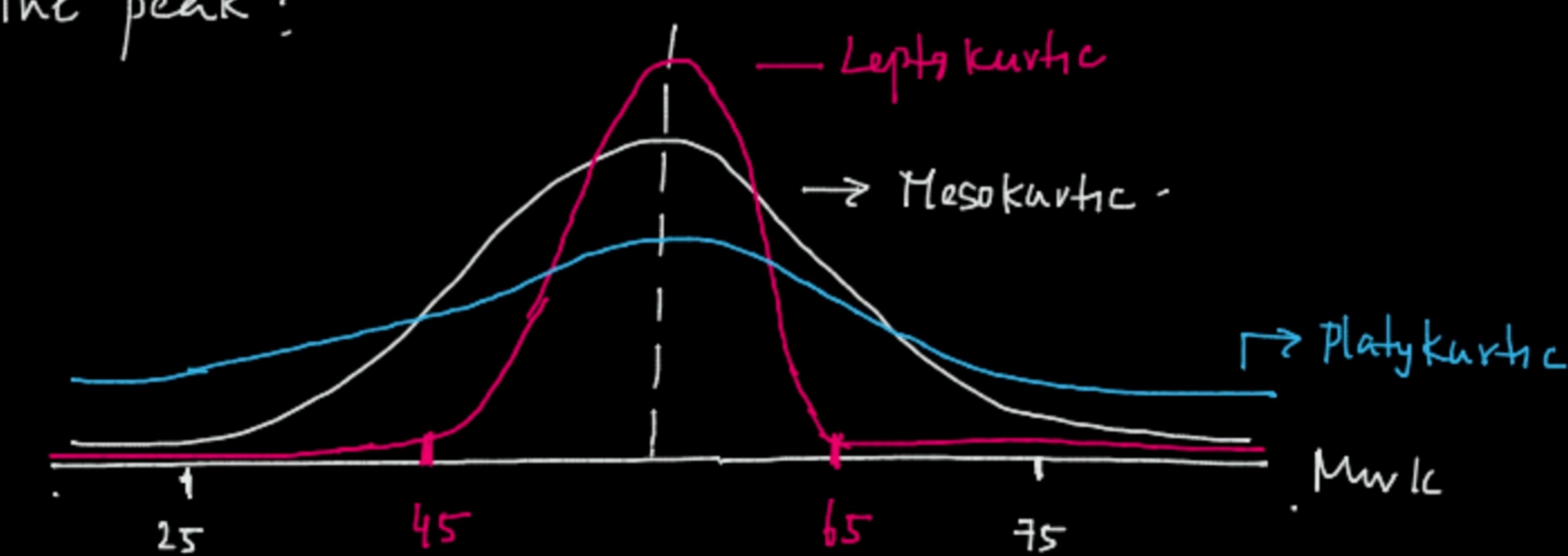


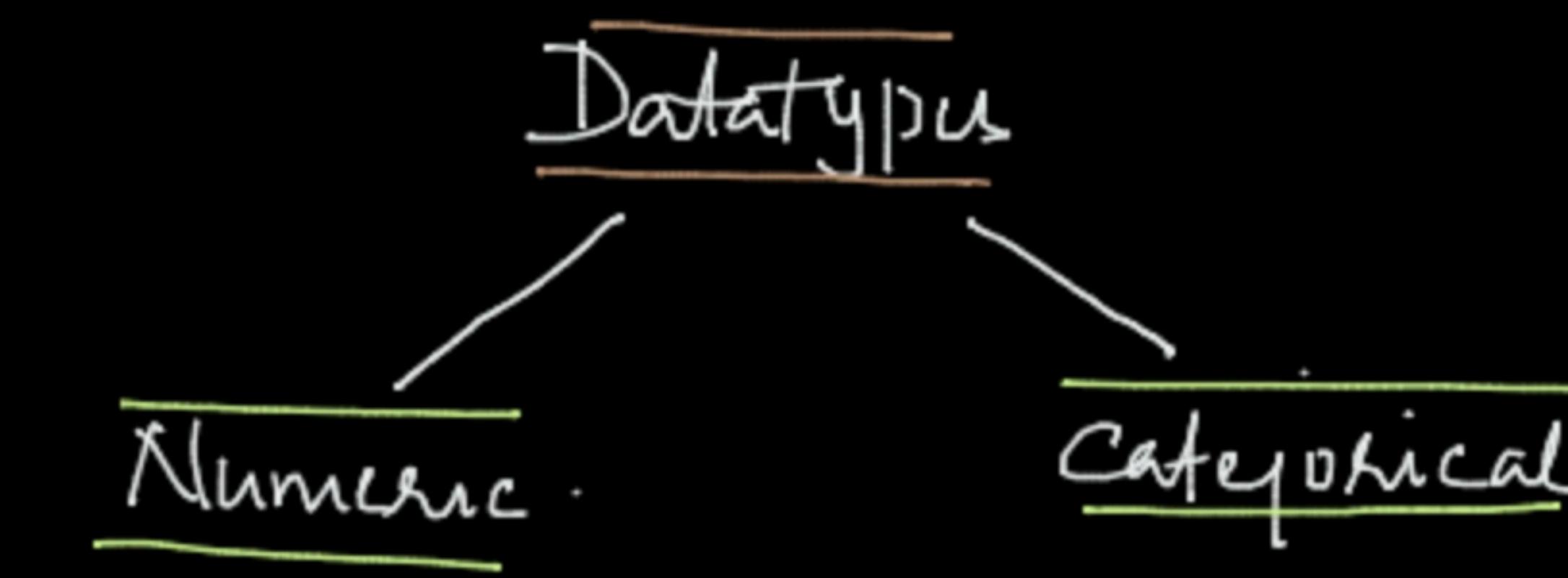
No skew



Right skew (+ve)

## (b) Kurtosis — How sharp is the peak?





- Counted

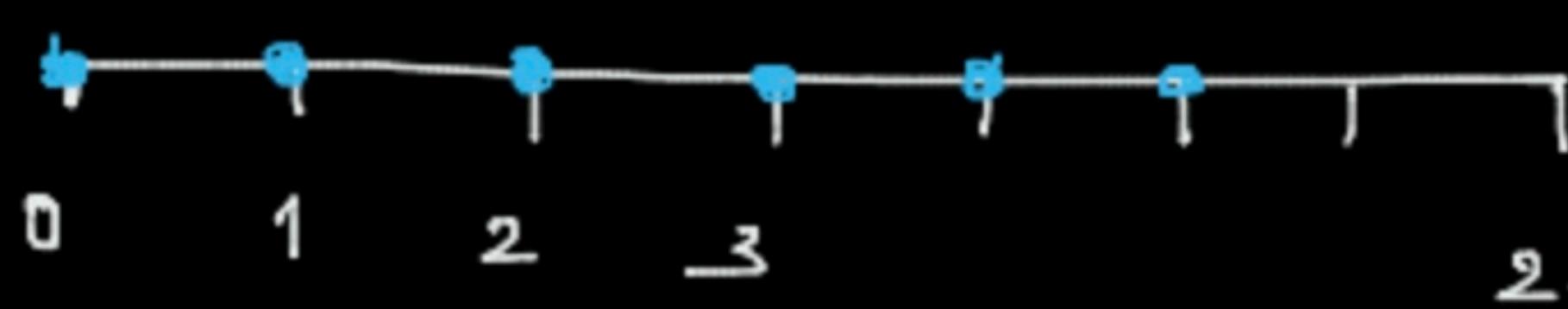
Eg: No. of students

No. of cars

No. of children

- Fractions don't make sense

- Only certain values are possible.



- Finite prob. of getting specific value

- Categorize the data
- Ex: Gender, Race, Religion

- only a label

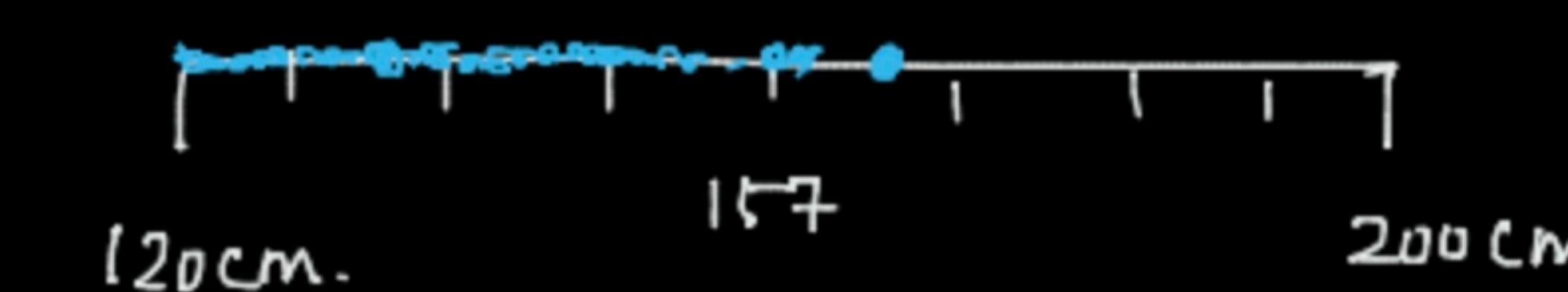
M/F . B/G .

0/1

150.75cm

85.3 kg

- Fractions are meaningful



- Any value is possible within a range

- Prob. of finding a particular value = 0

$\leftarrow$  Probabilities  $\rightarrow$ 

Random Expts: Result changes

Sample space  $\rightarrow$  [ All possible outcomes  
of a random expt. ]

Toss a coin  $\rightarrow$  [ H, T ]

Roll a dice  $\rightarrow$  [ 1, 2, 3, 4, 5, 6 ]

$$P(X) = \frac{\text{No. of Fav. events}}{\# \text{ Total possible outcomes}}$$

$$P(T) = \frac{1}{2}$$

$$P(3) = \frac{1}{6}$$

Draw a card  $\rightarrow$  [ 13 H, 13 C, 13 S, 13 D ]

$$P(\text{Heart}) \rightarrow \frac{13}{52} = \frac{1}{4}$$

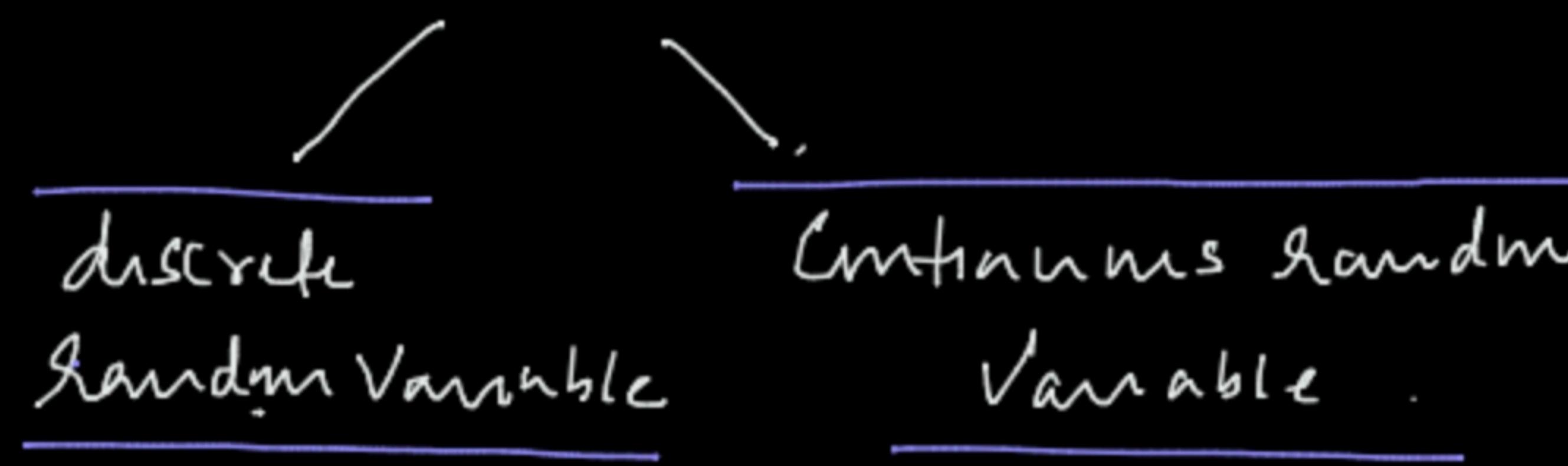
$P(X = 150) \rightarrow [120, -\infty \dots \dots , 200]$

$$\rightarrow \frac{10}{\infty} \rightarrow 0$$

$P(151 < x < 152) \rightarrow \text{Finite}$

## Random Variable:

Variable used to hold the result of a random experiment.



$X \rightarrow$  hold the result } -  
tossing a coin

$X \rightarrow$  prob. } finding  
a person whose height = 150 cms.

1. Normal Dist.
2. Z-dist (or) Std Normal dist.
3. t-dist.

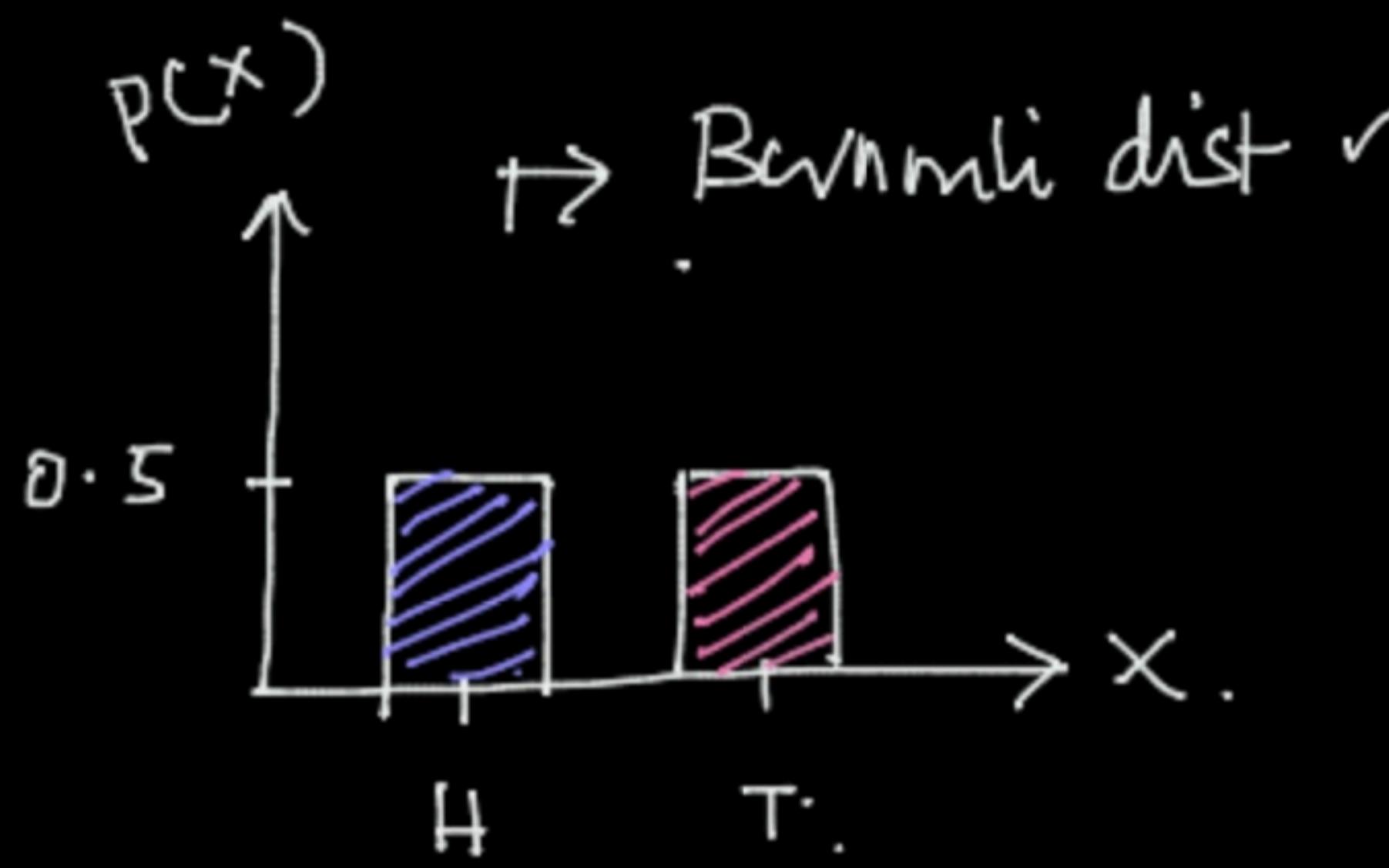
Probability distribution of a Random Variable

- The prob. associated with all possible values

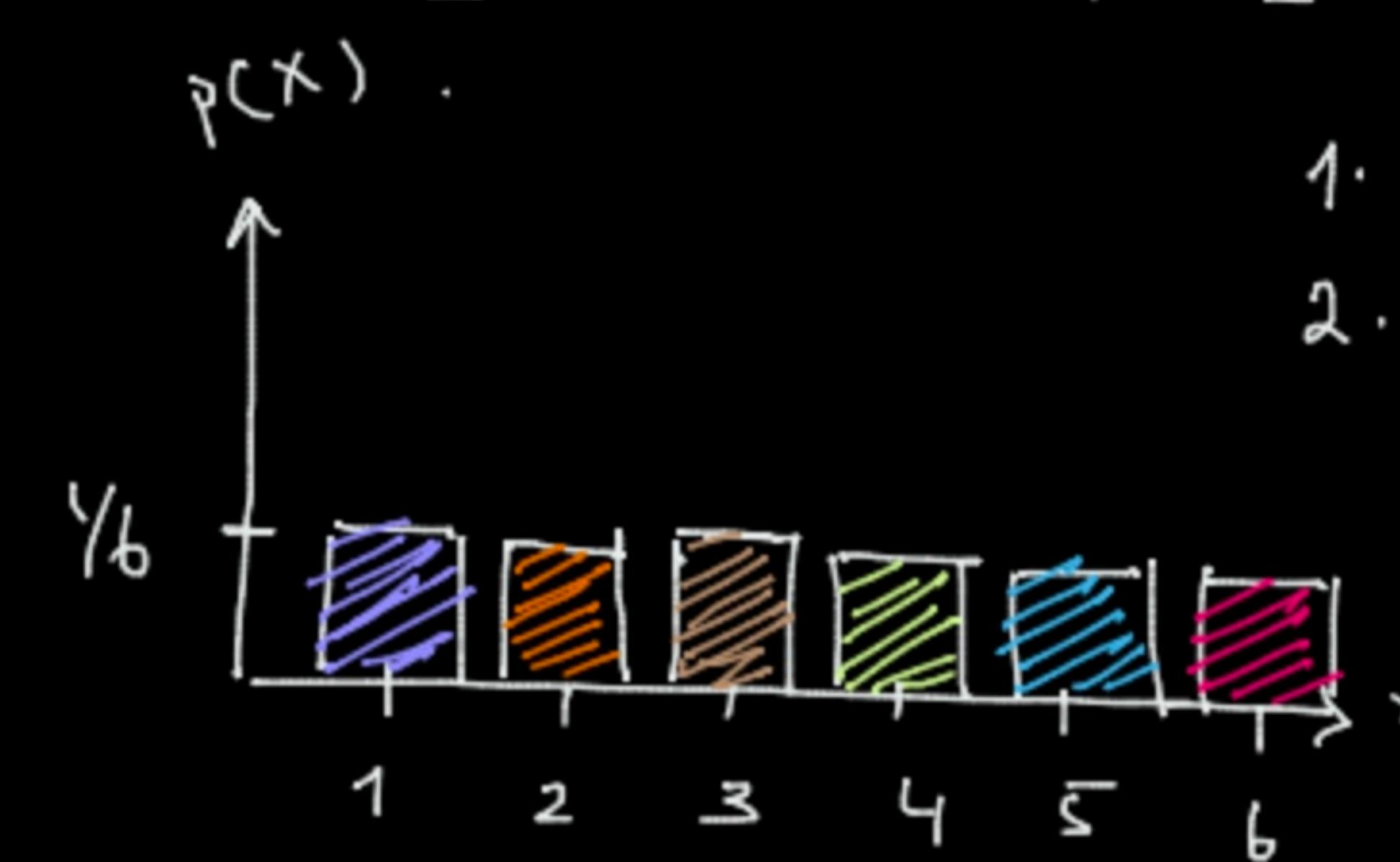
The Random Variable can hold -

1. Bernoulli
2. Binomial
3. Uniform

$$X \rightarrow [H, T]$$



$$X \rightarrow [1, 2, 3, 4, 5, 6] \quad (3) \text{ Bernoulli}$$



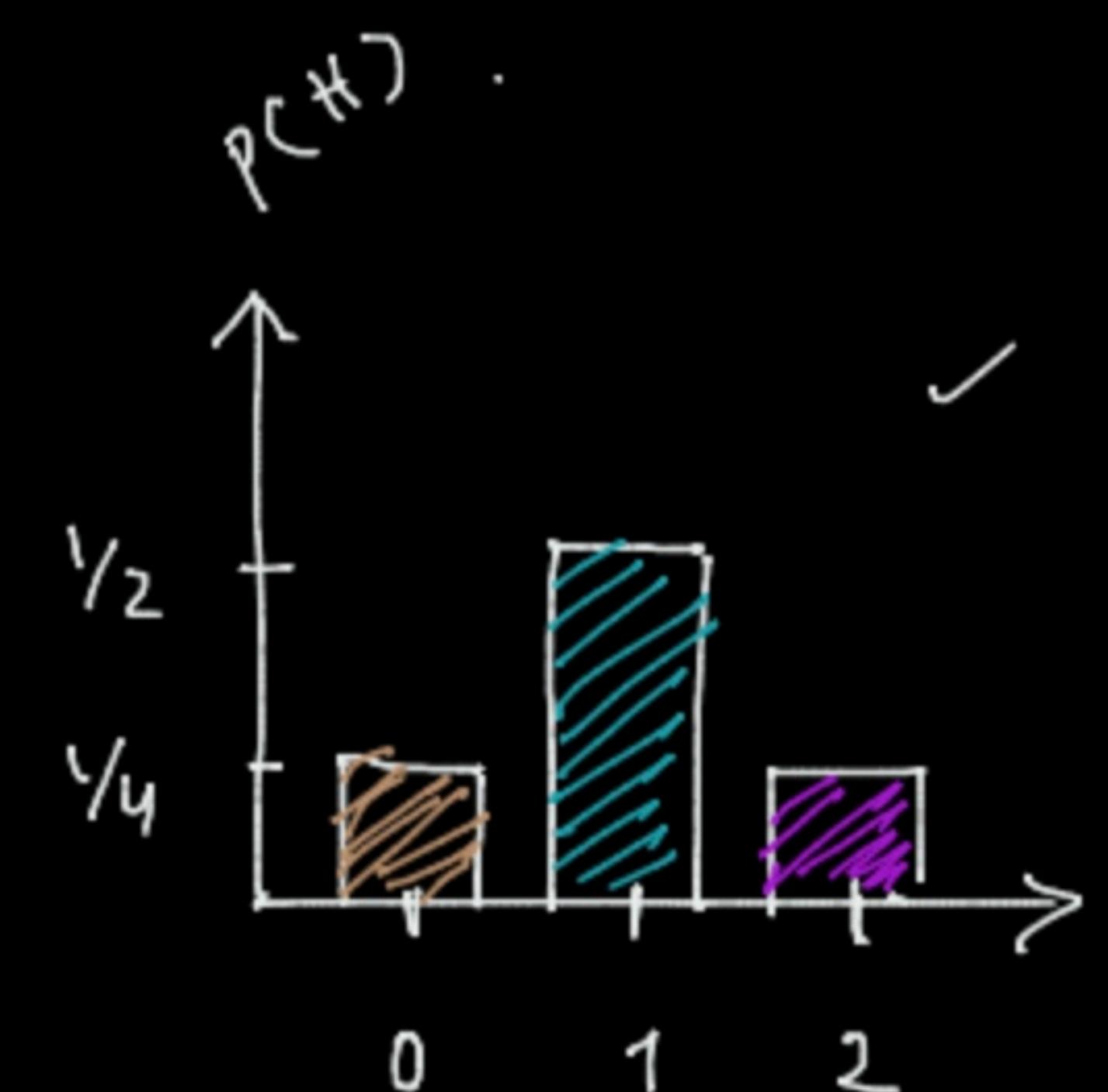
1. Expt. performed only once.
2. There should only be 2 possible outcomes.

### Binomial Distribution

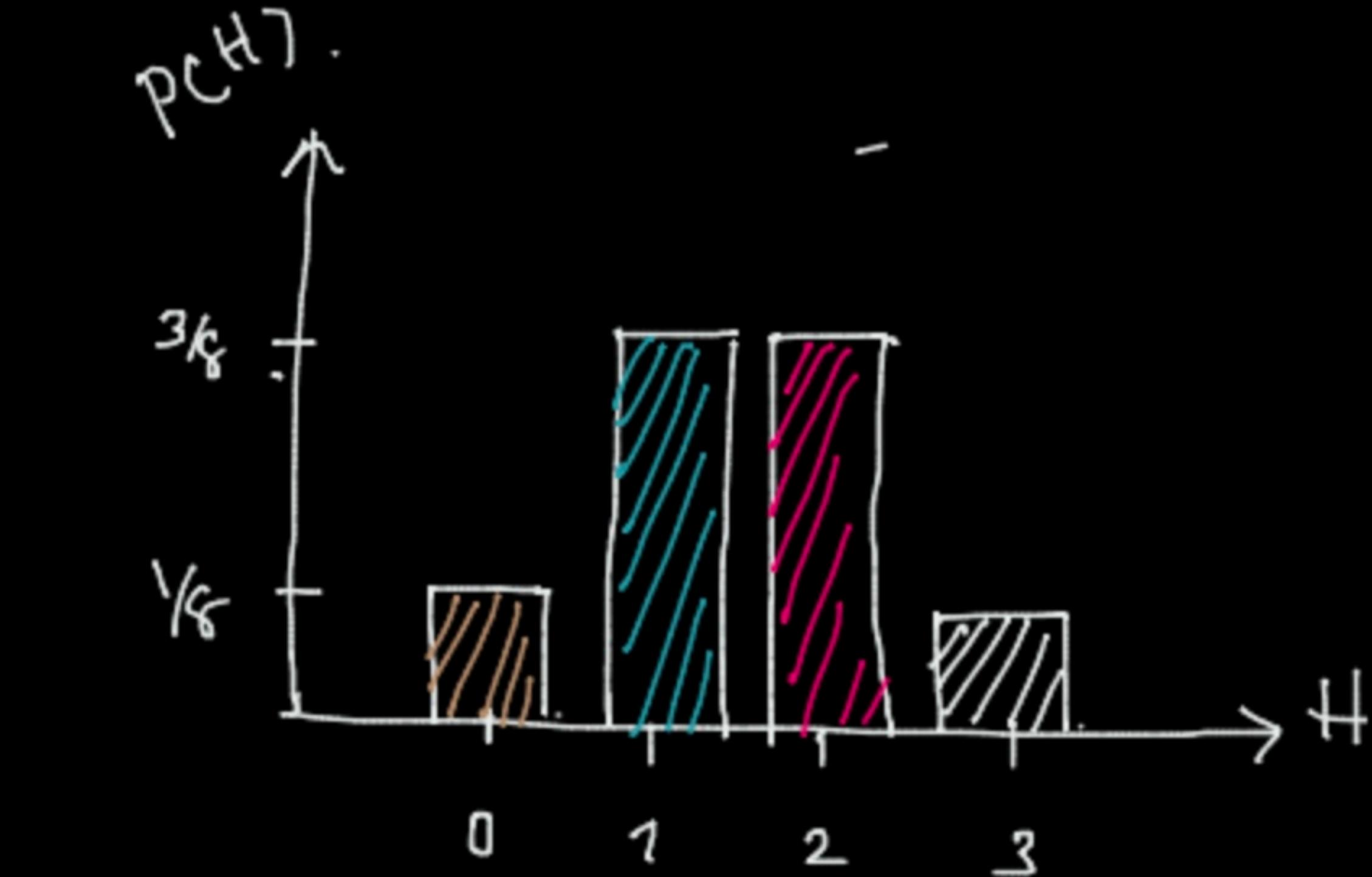
1. Expt. should be repeated more than once.
2. There should only be 2 possible outcomes.

Toss a coin twice

$$\rightarrow \begin{bmatrix} H & H \\ H & T \\ T & H \\ T & T \end{bmatrix}$$



T	T	T	-
H	H	H	-
H	H	T	-
H	T	H	-
H	T	T	-
T	H	H	-
T	H	T	-
T	T	H	-

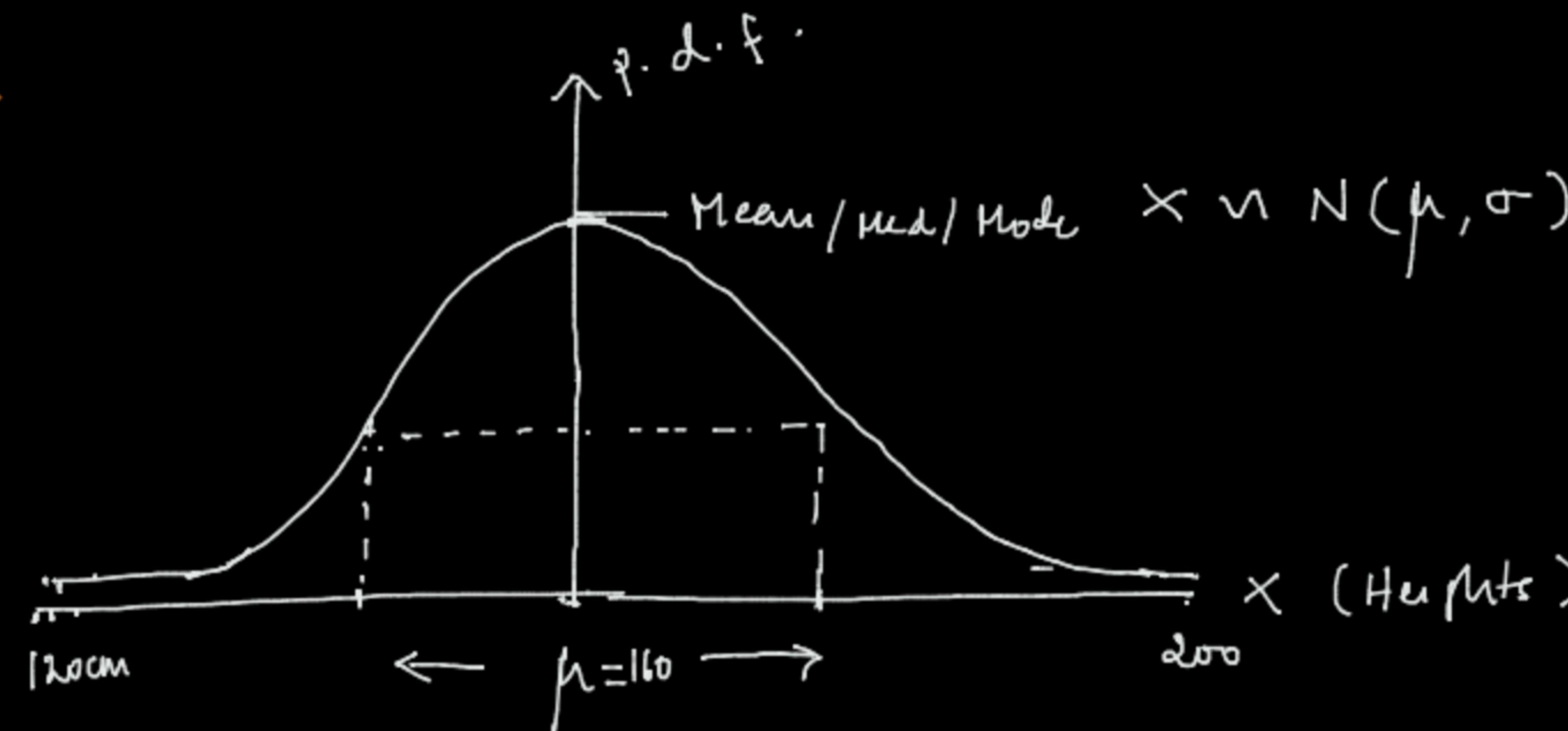


(a) Normal distribution $\leftarrow$  Continuous distributions  $\rightarrow$ 

1. Bell shaped .

2. Symmetric

$$3. \text{ p.d.f} = \frac{1}{2\pi\sigma^2} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$



4. Mean = Median = Mode

5. Empirical Rule



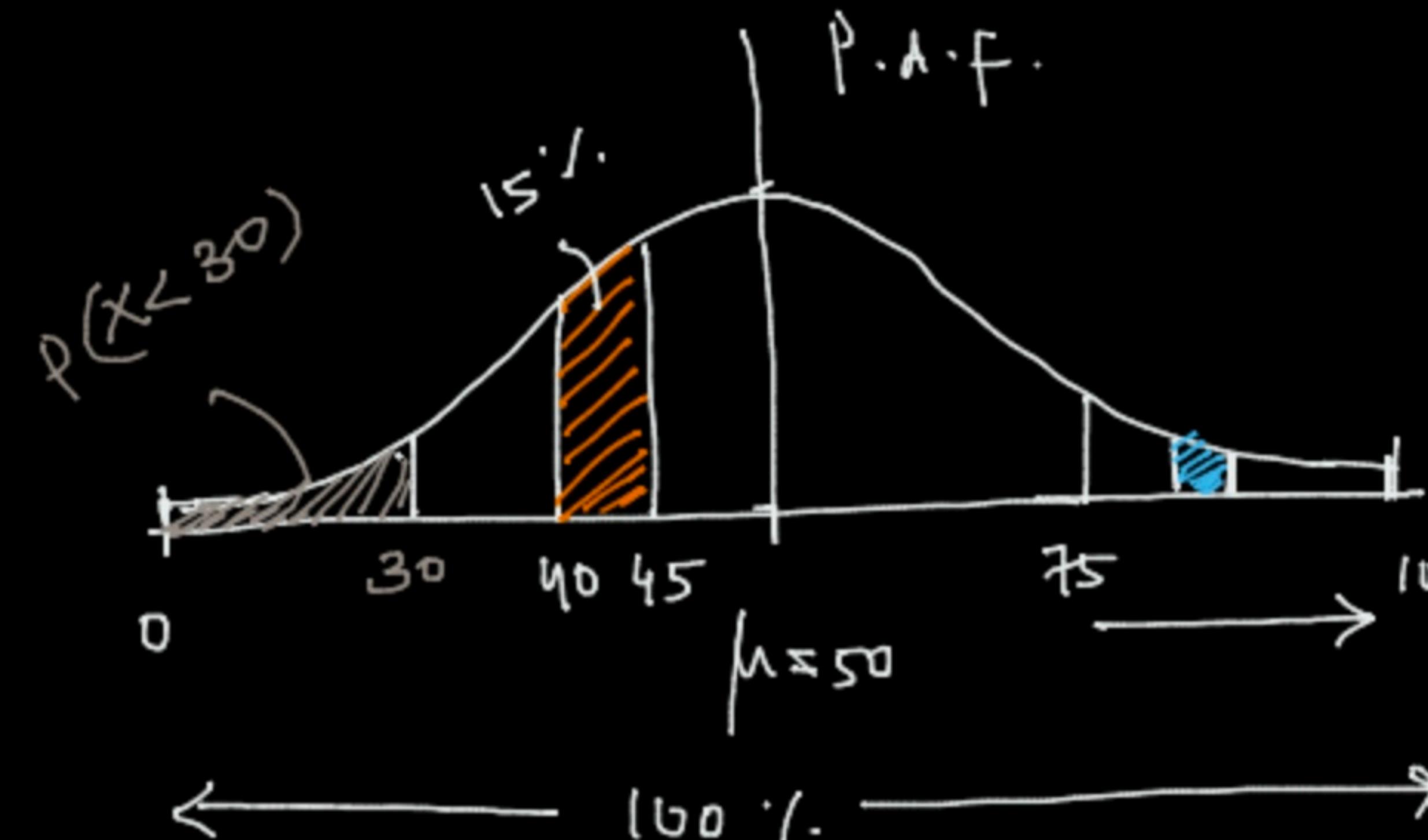
(x)  
heights

⋮  
⋮  
⋮  
⋮  
⋮

$$\mu = 160$$

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

$$\sigma = 5$$



p.d.f  
Area under the curve  $\Rightarrow$  Probability

0  $\rightarrow$  100

Total Area = 100 %

1.  $P(40 < x < 45) \rightarrow 15\%$

15% of students have marks  
between 40 & 45

$$\mathcal{N}(60, 5)$$

$$\mu = 60$$

68%  $\rightarrow$  55 & 65

$$\sigma = 5$$

95%  $\rightarrow$  50 & 70

$$(x)$$

99.7%  $\rightarrow$  45 & 75

$$\rightarrow \begin{cases} 45 \\ 70 \\ 65 \\ \vdots \end{cases} \quad \left| \begin{array}{l} \mu = 55 \\ \sigma = 7 \end{array} \right.$$

Standardize

$$Z_x = \left( \frac{x - \mu}{\sigma} \right)$$

Hawls

$x$	$\frac{x - \mu}{\sigma}$	$Z_x = \frac{x - \mu}{\sigma}$
45	-17	-1.54
70	+8	+0.72
54	-8	-0.72
66	+4	+0.36
75	+13	+1.18
$\mu = 62$		$\mu = 0$

$$\sigma = 11$$

Scaling  $\rightarrow$  Bring all features to same scale

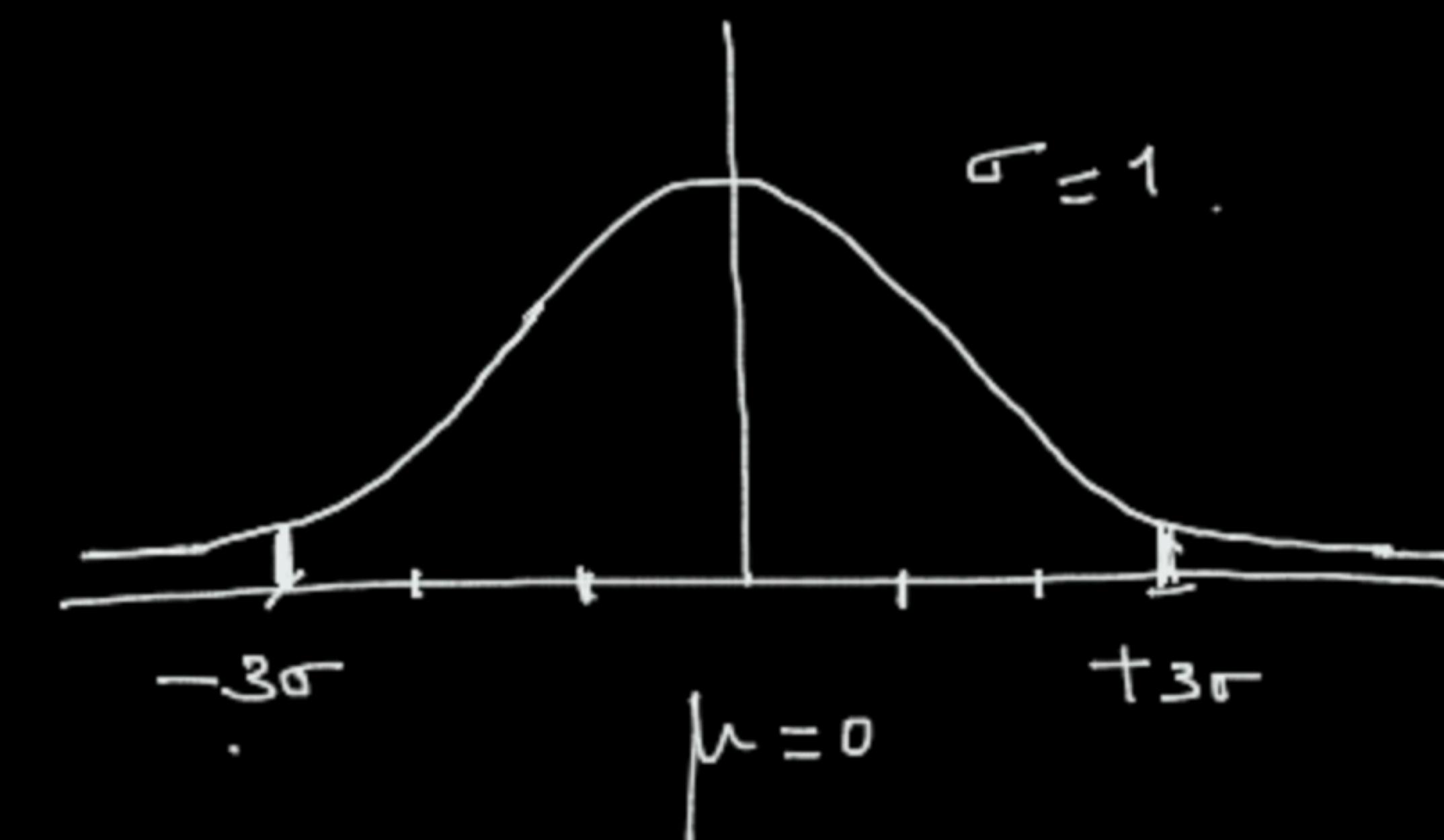
1. Normal dist  
 $\mu = 0; \sigma = 1$

$Z$ -distribution / Std. Normal distribution

$$\begin{array}{c} Z_x \\ \hline -1.54 \\ 0.72 \\ -0.72 \\ 0.36 \\ +1.18 \\ \hline \end{array}$$

$\mu - 3\sigma, \mu + 3\sigma \rightarrow 99.7\%$

$[-3, +3] \rightarrow 99.7\%$



$\longleftrightarrow 99.7\% \longrightarrow$



