

# Probability And Statistics

- fundamental area
- histogram, Pdf, cdf
- random variable

$$n = 2, 3, \dots$$

dice = {1, 2, 3, 4, 5, 6} Discrete Random  
rolled a fair dice  $\rightarrow$  equally likely

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

$$P(n=2) = \frac{1}{6}$$

$$P(X \text{ is even}) = \frac{1}{2} = P(X=2) + P(X=4) + P(X=6)$$

$$= \frac{1}{6} + \frac{1}{6} + \frac{1}{6}$$

$$P(X \text{ is odd}) = \frac{1}{2}$$

Height of a Random Student

Continuous Random Variable

Outlier

Y: Height of Student

[122.2, 146.4, 132.5, ...]

12.26 156.3

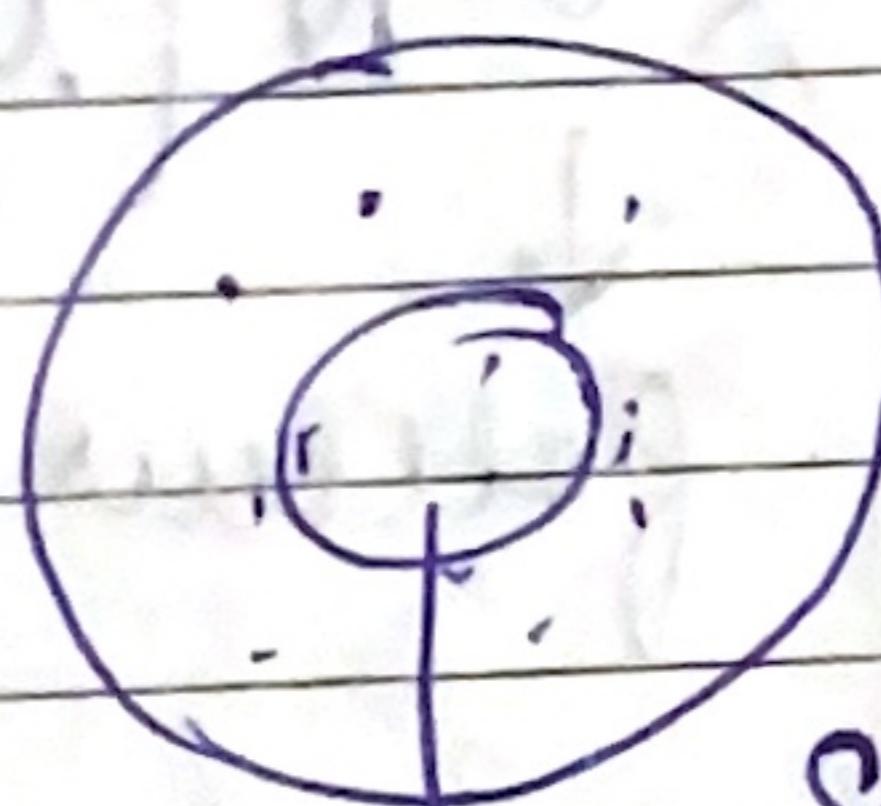
Outlier

## Population and Sample

→ estimate average height of human?

set of all people in world

Mean  $\bar{y} = \frac{1}{TB} + \sum_{i=1}^{TB} h_i$



Random Sample

Sample of size 1000

Subset

$$\bar{h} = \frac{1}{1000} + \sum_{i=1}^{1000} h_i$$

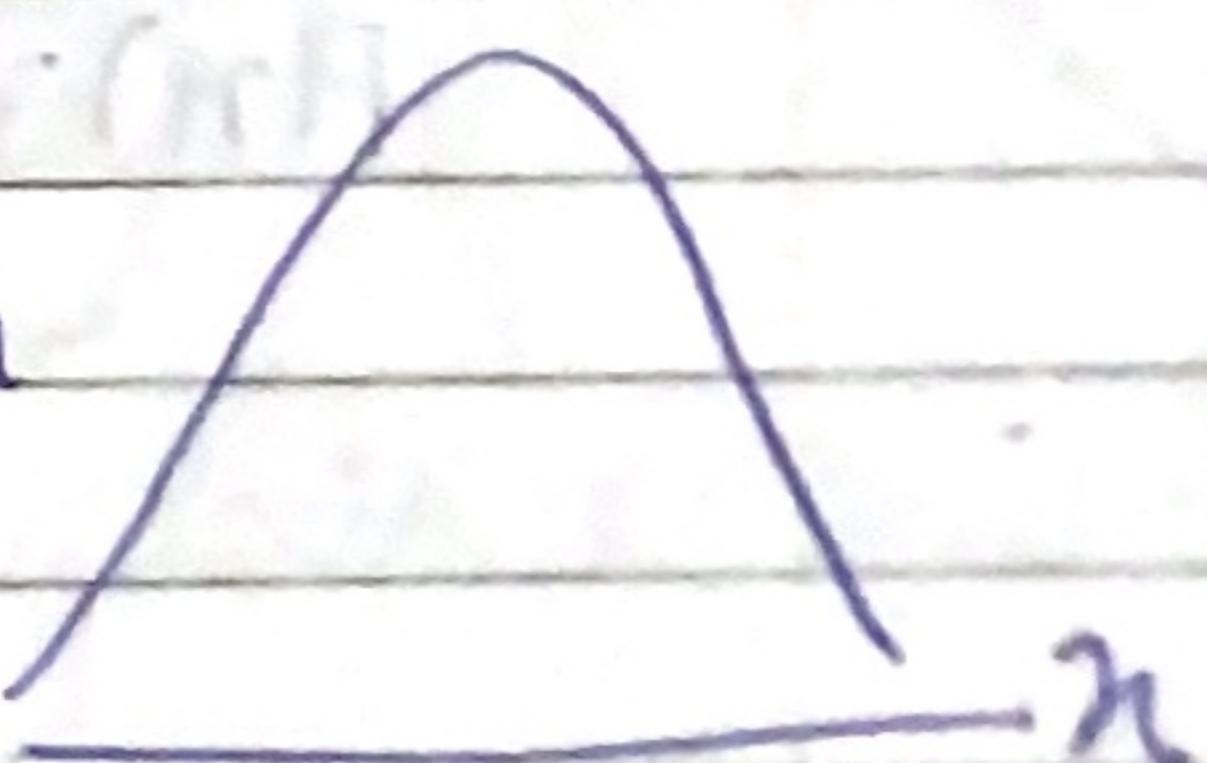
height in my sample

As sample size increases

$\bar{x} = \mu$   
 Sample mean      population Mean

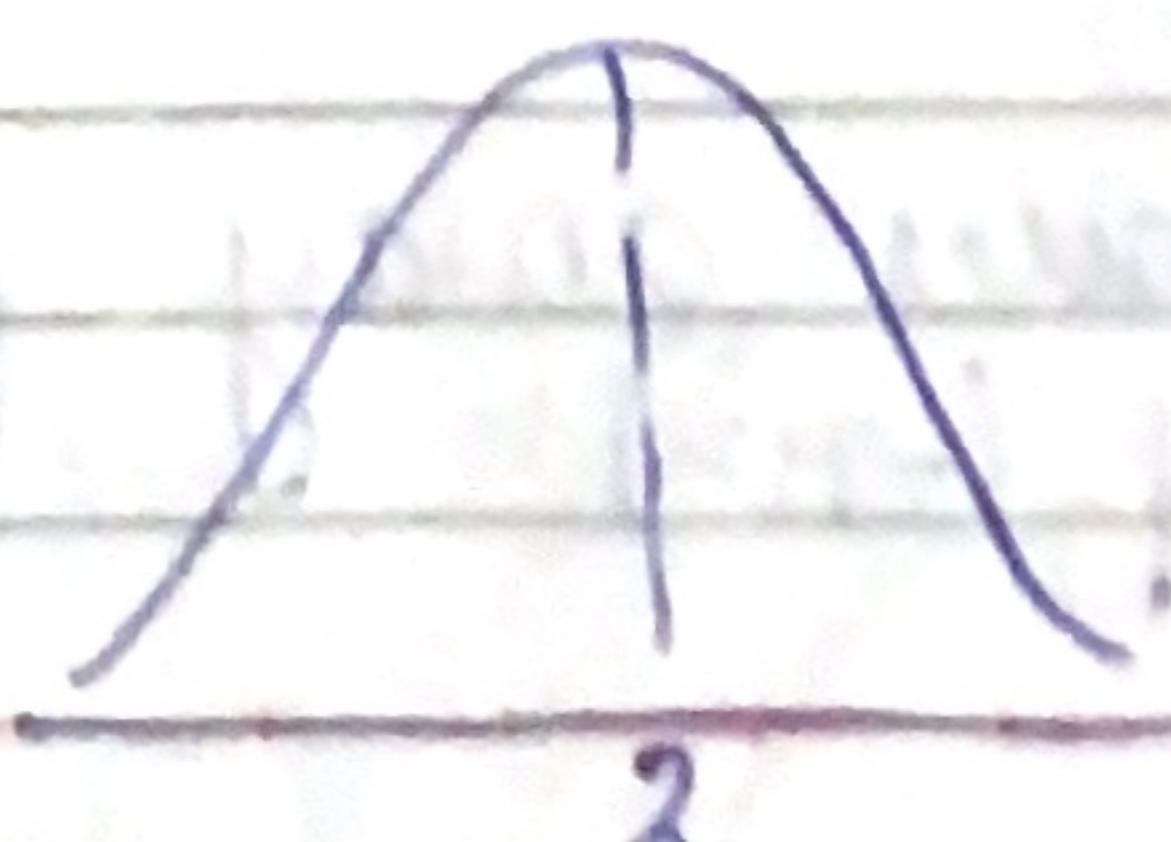
## Gaussian Distribution

PDF of a Gaussian distribution  
 random variable



$x$ : continuous random variable

Parameters of Gaussian dist =  $\mu, \sigma^2$



$\mu = 2$

# Normal/Gaussian

Date :

Page No.

$$X \sim N(\mu, \sigma^2)$$

↑  
↓

follows → mean  
Variance

$$P(X = \bar{x}) \xrightarrow{160} P(n) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(\bar{x}-\mu)^2}{2\sigma^2}\right\}$$

$$\text{let } \mu = 0 \quad \sigma^2 = 1$$

$$P(n) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}n^2\right\} = y \quad y = p(n)$$

Constant

$$P(n) = y = \exp(-n^2)$$

$$n \uparrow \quad (-n^2) \downarrow \\ \exp(-n^2) \downarrow$$

- ① as  $n$  moves away from  $\mu$ ,  $y \downarrow$
- ② symmetric
- ③  $n$  moves away from  $\mu$ ,  $y$  reduces  
 $\exp(-n^2)$

$$y = \exp(-x^2)$$

$$n=0 \quad y=1$$

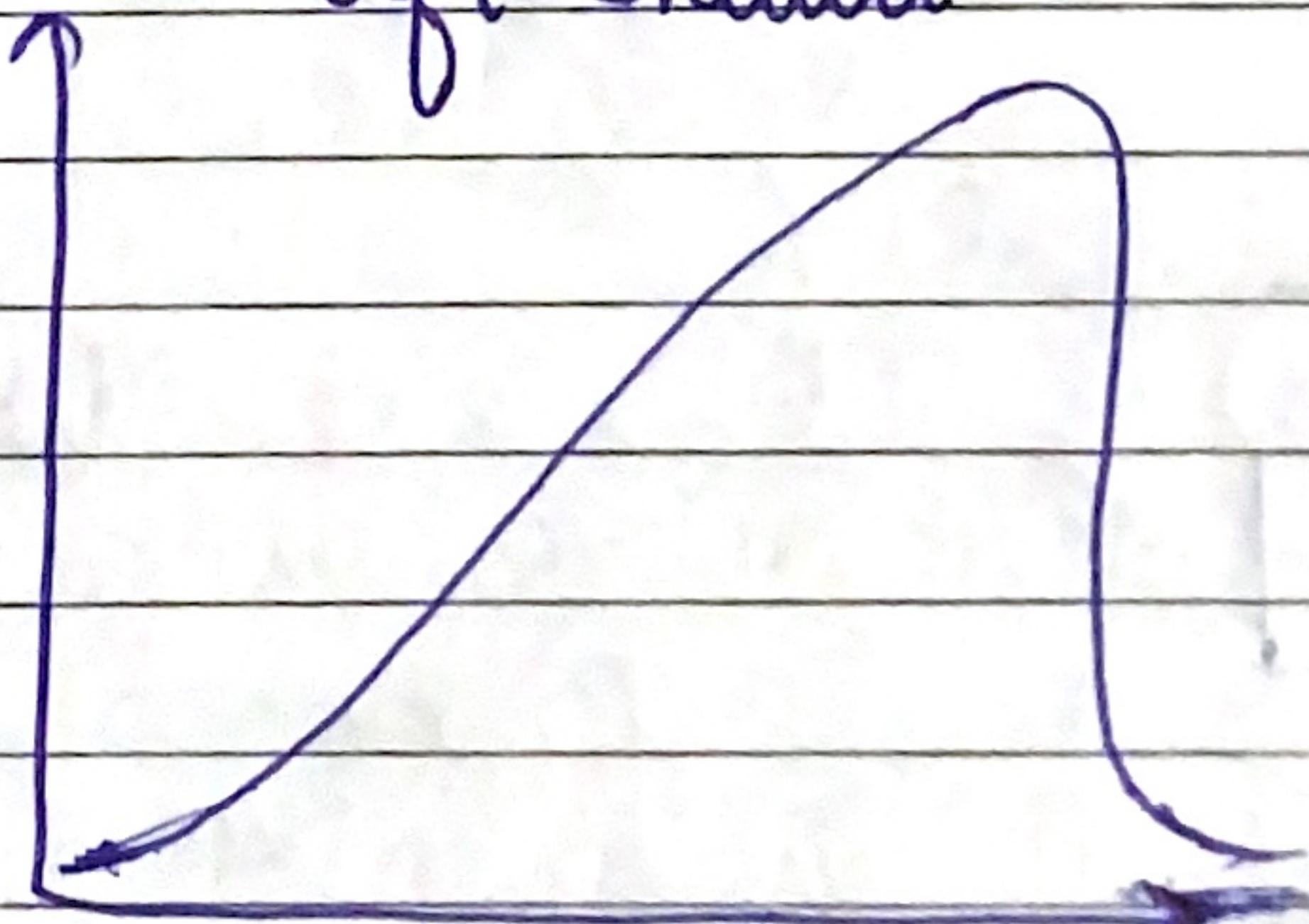
$$n=1 \quad y = \exp(-1) = \frac{1}{e^1} = 0.3678 \quad \left. \right\} 20x$$

$$n=2 \quad y = \exp(-4) = \frac{1}{e^4} = 0.018 \quad \left. \right\} 100x$$

$$n=3 \quad y = \exp(-9) = \frac{1}{e^9} = 0.000123$$

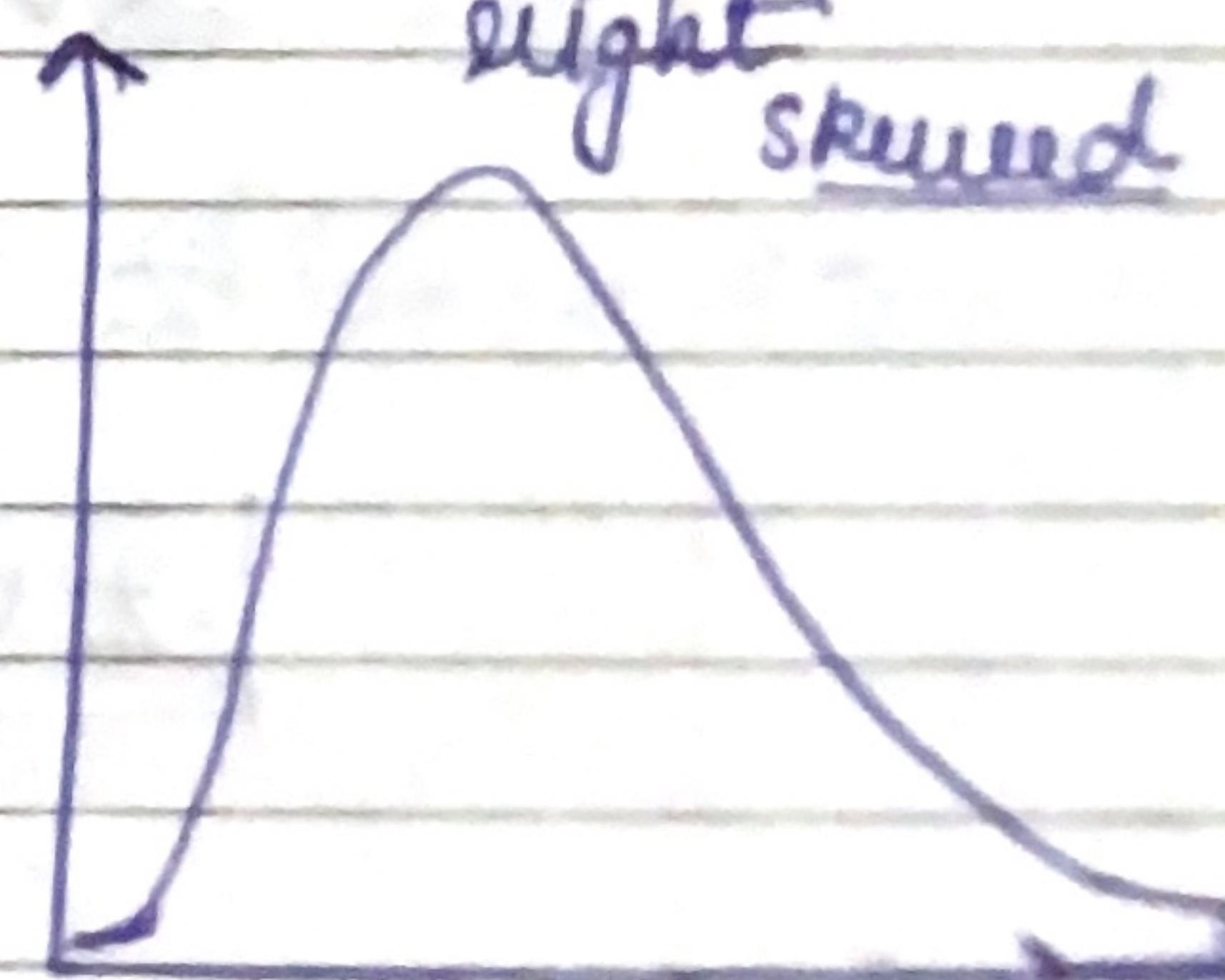
exponentially  
quadratic decay

left Skewed



Negative Skew

right skewed



Positive Skew

## Standard Normal Variate (2)

①  $Z \sim N(0, 1)$

$$\mu = 0$$

$$\sigma^2 = 1$$

② Let  $X \sim N(\mu, \sigma^2)$  | Standardization

$$\hookrightarrow [x_1, x_2, \dots, x_{50}]$$

$$x_i' = \frac{x_i - \mu}{\sigma}$$

$$x_i' \sim N(0, 1)$$

↑  
standard  
normal  
variant

Standardization of Data

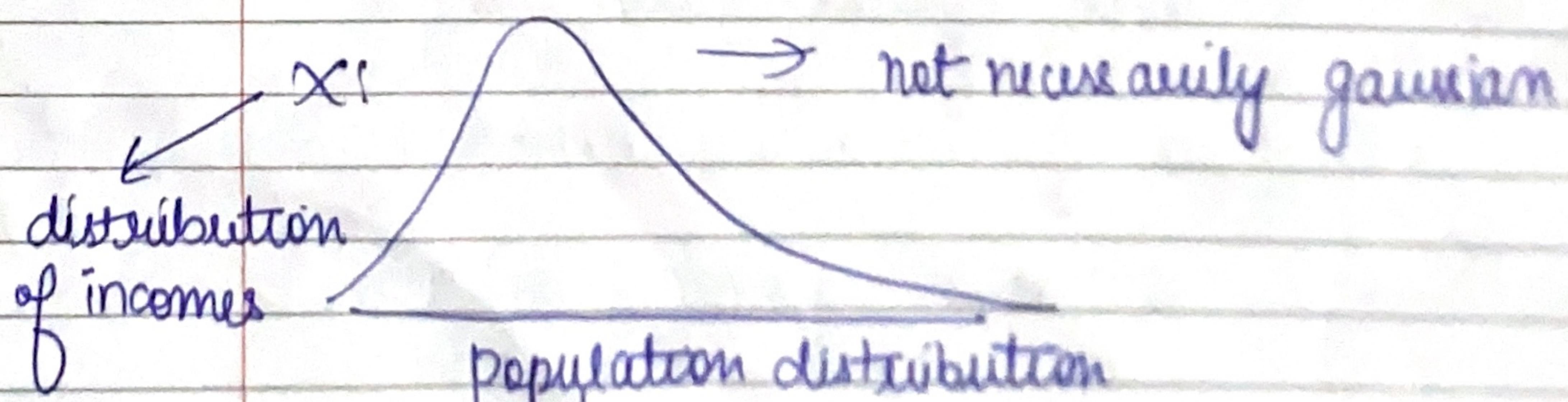
$$\bar{X} \sim N(\mu, \sigma^2)$$

$$Z = \frac{\bar{X} - \mu}{\sigma}$$

$$Z \sim N(0, 1)$$

Sampling Distribution & CLT

(Central Limit Theorem)



Population distribution

Sample of size ( $n$ )  $\rightarrow S_2 \rightarrow \frac{S_2}{n}$

3

$\rightarrow S_m \rightarrow \frac{S_m}{n}$

$\bar{x}_1, \bar{x}_2, \dots, \bar{x}_m$  :-  $m$  sample - means

$\bar{x}_i$   $\sim$  distribution

$\uparrow$  distribution of sample means

dist of  $\bar{x}_i$  = Sampling dist of sample - mean

CLT

need not gaussian

( $\mu$ ) ( $\sigma^2$ )

infinite  $\mu \& \sigma^2$

If  $X$ : finite mean and variance

$S_1, S_2, \dots, S_m$  (sample of size  $n$ )

(Pareto dist)

$\downarrow$

$\downarrow$

Sample  
mean

$\bar{x}_1, \bar{x}_2, \dots, \bar{x}_m$

dist of  $\bar{x}_i$  = Sampling dist of sample mean.

CLT:-

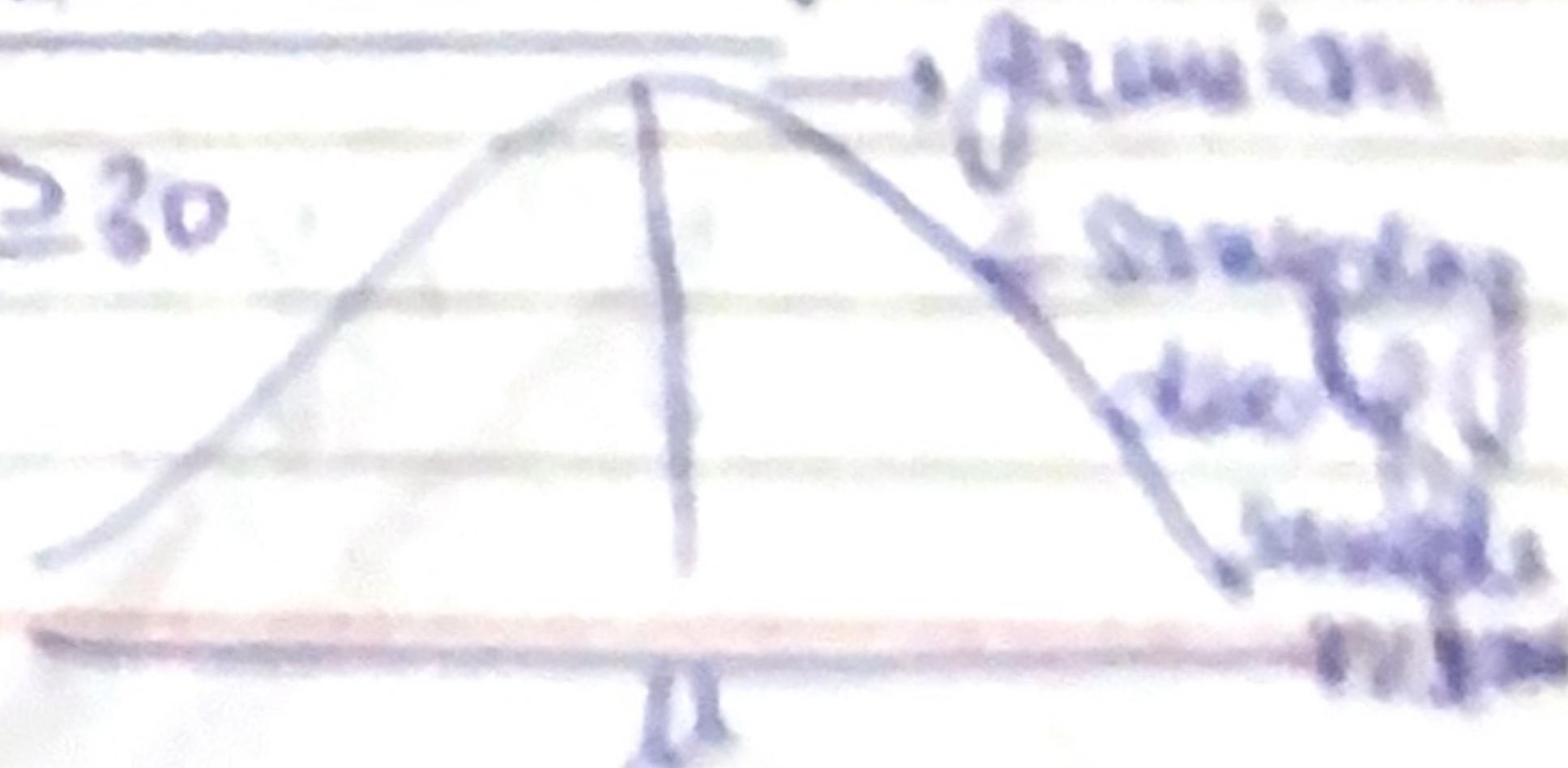
$\bar{x}_i \rightarrow N\left(\mu, \frac{\sigma^2}{n}\right) \text{ as } n \rightarrow \infty$

gaussian distribution

$\mu \rightsquigarrow$  mean of  $x_i$ , &

$n \geq 30$

$\frac{\sigma^2}{n} \rightsquigarrow$  var of  $\bar{x}_i$



# Quantile-Quantile plot: (Q-Q plot)



$x: x_1, x_2, x_3, \dots, x_{500}$

(Q) is  $x$  gaussian dist? (graphical Q-Q plot method)

statistical testing  
(KS, AD)

How

① Sort  $x_i$  and compute percentiles

$x_1, x_2, \dots, x_{500}$



Sort (asc)

$x'_1, x'_2, \dots, x'_{500}$  → Percentile ② (100)

$x'_5, x'_{10}, x'_{15} \dots x'_{500}$

$(x^{(1)}, x^{(2)}, x^{(3)})$

→ 1<sup>st</sup> percentile value of  $x'_1$

②

$y \sim N(0, 1)$  → Std Normal dist

$y_1, y_2, y_3, \dots, y_{1000}$

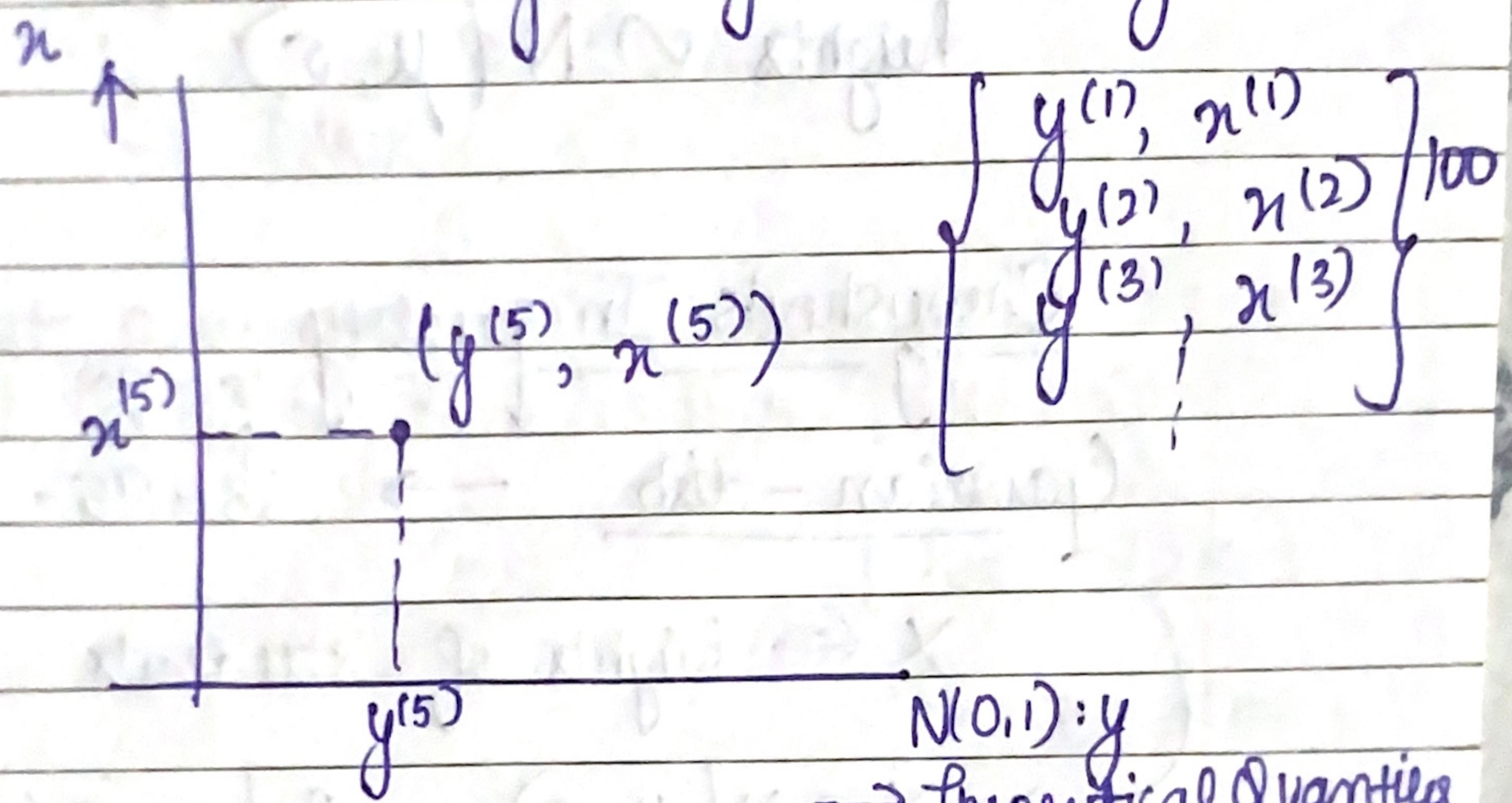
→ 1000 obsr from  $N(0, 1)$

$y'_1, y'_2, y'_3, \dots, y'_{1000}$   
↓ percentile

$y^{(1)}, y^{(2)}, y^{(3)}, y^{(4)}, \dots, y^{(100)}$

(3)

Plot Q-Q plot using  $x^{(1)}, x^{(2)} \dots x^{(100)}$   
 $y^{(1)}, y^{(2)} \dots y^{(100)}$



(100)

 $x_{500}$ 

If  $(y^{(i)}, n^{(i)})$   $i \rightarrow 100$  lie on straight line then  
 $\times 8 Y$  have similar distribution.

$N(0,1) : y$   
 $\rightarrow$  theoretical Quantiles

How / where to use distributions?

probability  $\rightarrow$  data analysis  $\rightarrow$  answering ques  
 about data

(Q1) Company - XYZ

Task: Order t-Shirts for all employees (100k)  
 S, M, L, XL

(a) How many XL t-shirts to order?

① collect data for all (100k)

height  $\rightarrow$  180cm  $\rightarrow$  XL  
 160  $\rightarrow$  L

collect heights of 500 random employees

mean, std dev

$$\text{heights} \sim N(\mu, \sigma)$$

Chebyshev's Inequality

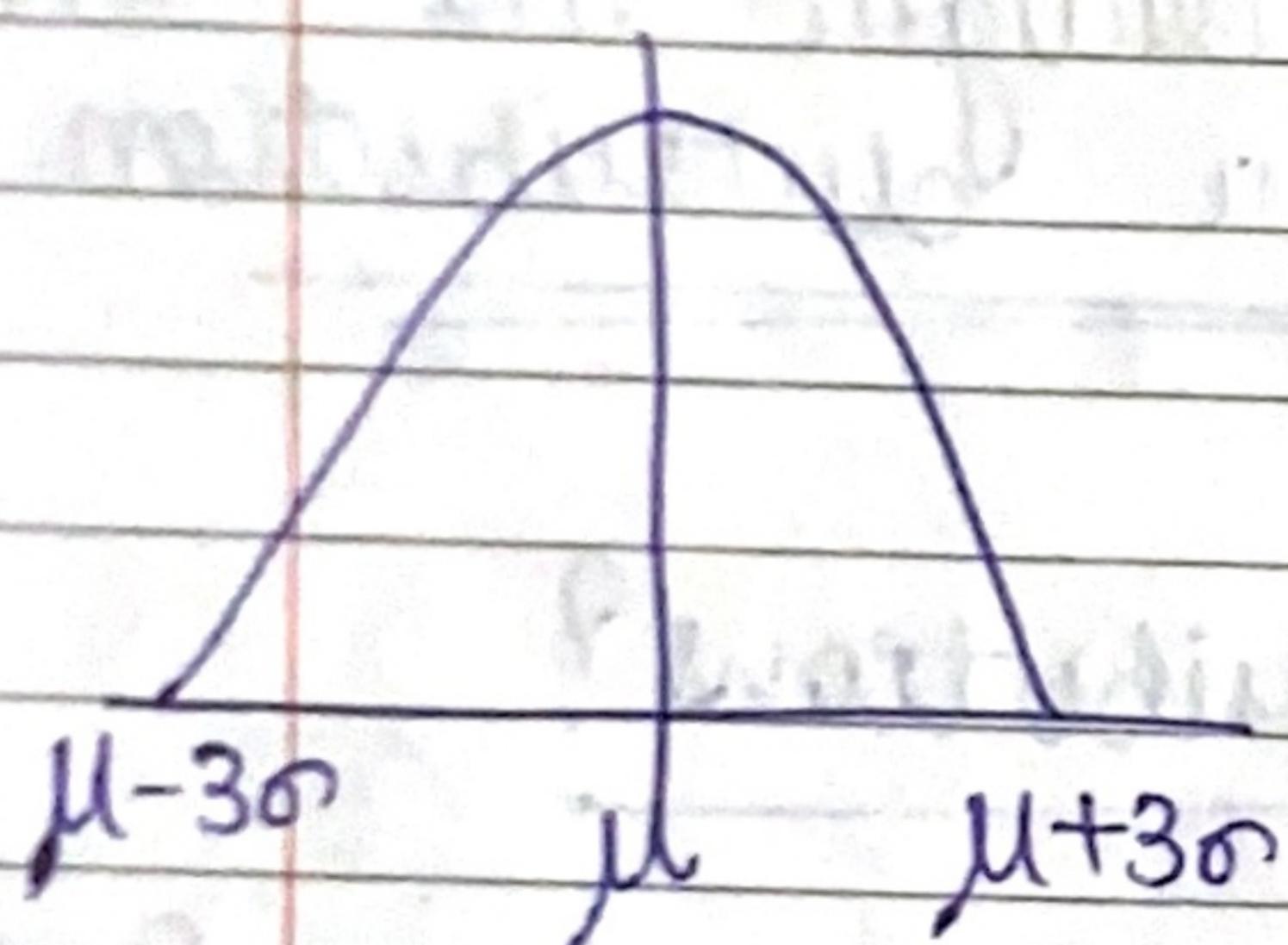
Gaussian - dist  $\rightarrow$  68-95-99.7 rule

1σ    2σ    3σ

$X \leftarrow \text{heights of students}$

$$X \sim N(\mu, \sigma) \quad P(X \in [\mu - 2\sigma, \mu + 2\sigma]) = 95\%$$

95% of students heights  
lie in the [130, 170]



Q) what if I don't know dist

know  $\mu, \sigma$

non zero & finite

finite

$n\%$  of data lies  $\mu - 2\sigma \leq X \leq \mu + 2\sigma$

$y\%$  of data lies  $\mu - 1.5\sigma \leq X \leq \mu + 1.5\sigma$

Why?

$X: \text{r.v}$

Salaries of people  $\leftarrow$  don't know disb

$\mu, \sigma \leftarrow \text{know}$   
 $\downarrow \quad \downarrow$   
 $\$10K$   
 $\$40K$

26 (Q) what percentage of individuals have  
 Salary in range  $[20K, 60K]$

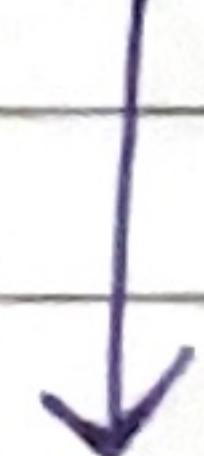
36 (Q) " " "  $[10K, 70K]$   
 $\downarrow$   
 $40 - 3(10K)$

Chebyshev's Inequality

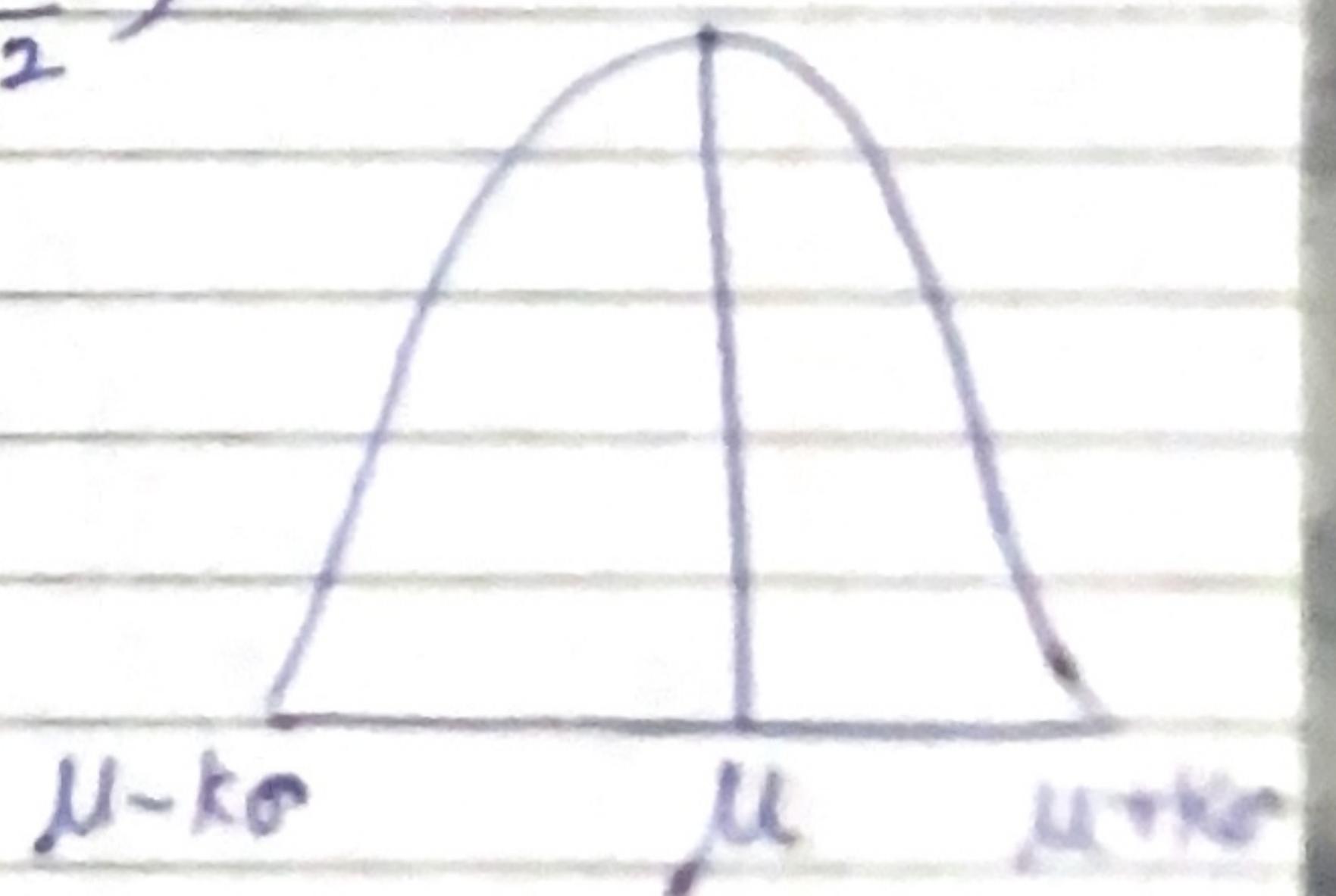
$X: \text{r.v}$        $\mu = \text{finite mean}$   
 $\sigma = \text{non zero & finite}$

P

$$P(|x-\mu| \geq k\sigma) \leq \frac{1}{k^2}$$



$$\begin{aligned} x - \mu &\geq k\sigma \\ \text{or } x &\leq \mu - k\sigma \end{aligned}$$



$$P(x \geq \mu + k\sigma) \leq \frac{1}{k^2}$$

Or

$$P(\mu - k\sigma < x < \mu + k\sigma) \geq 1 - \frac{1}{k^2}$$

all values equally probable

Date:

Page No.

## Uniform distribution

discrete (only take discrete value)  
continuous

PMF  $\rightarrow$  discrete r.v

PDF  $\rightarrow$  continuous r.v

Eg  $\rightarrow$  Throwing dice

$\{1, 2, 3, 4, 5, 6\}$   
equally likely  $[1/6]$

$$a = 1$$

$$b = 6$$

$$n = b - a + 1 = 6 - 1 + 1 = 6$$

$U(2, 6)$

$$a = 2$$

$$b = 6$$

$U[a, b]$

$a \in \{-\dots, -2, -1, 0, 1, 2, \dots\}$

$b \in \{-\dots, -2, -1, 0, 1, 2, \dots\}, b \geq a$

$$n = b - a + 1$$

$k \in \{a, a+1, \dots, b-1, b\}$

pmf  $\frac{1}{n}$

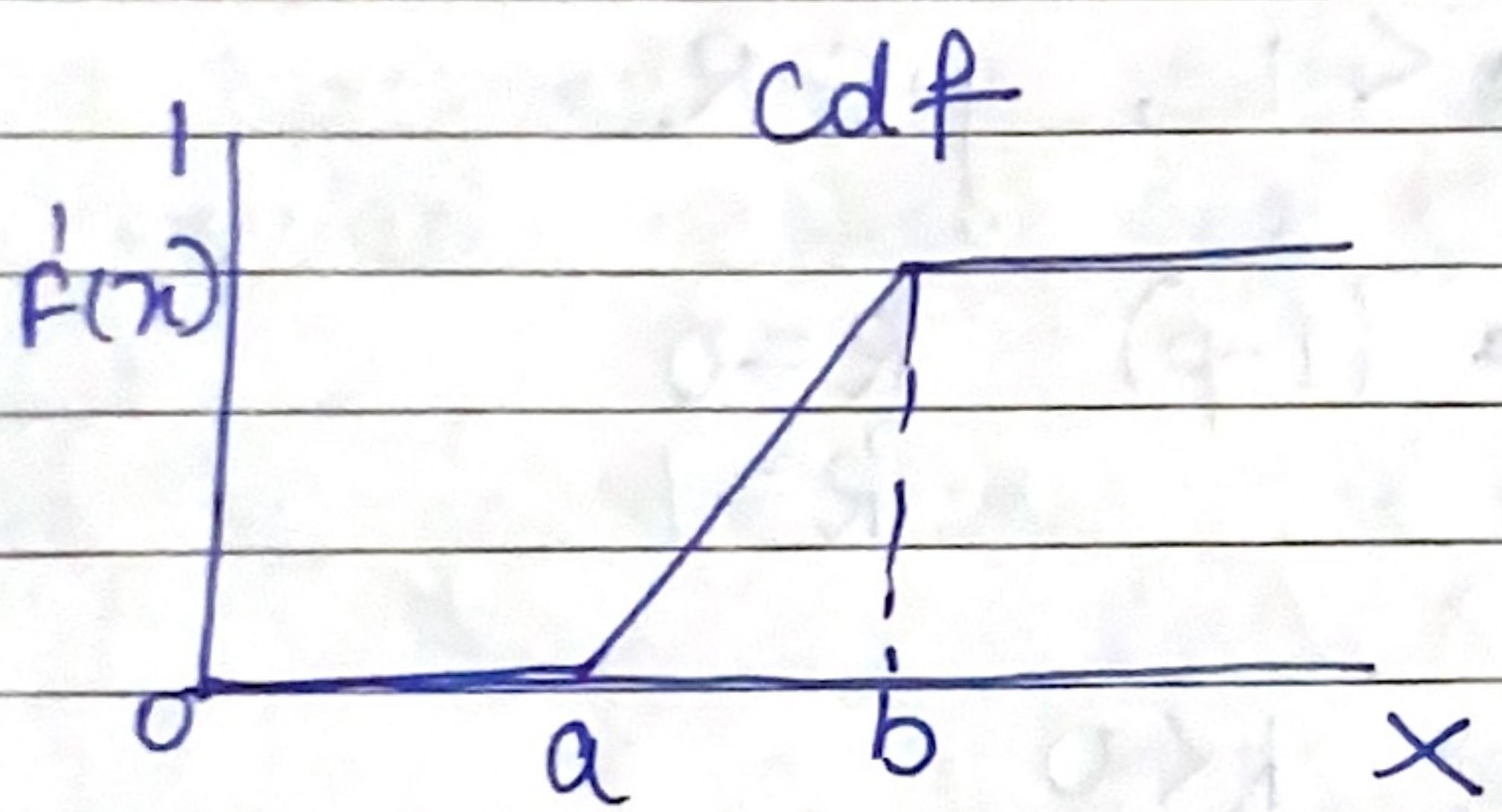
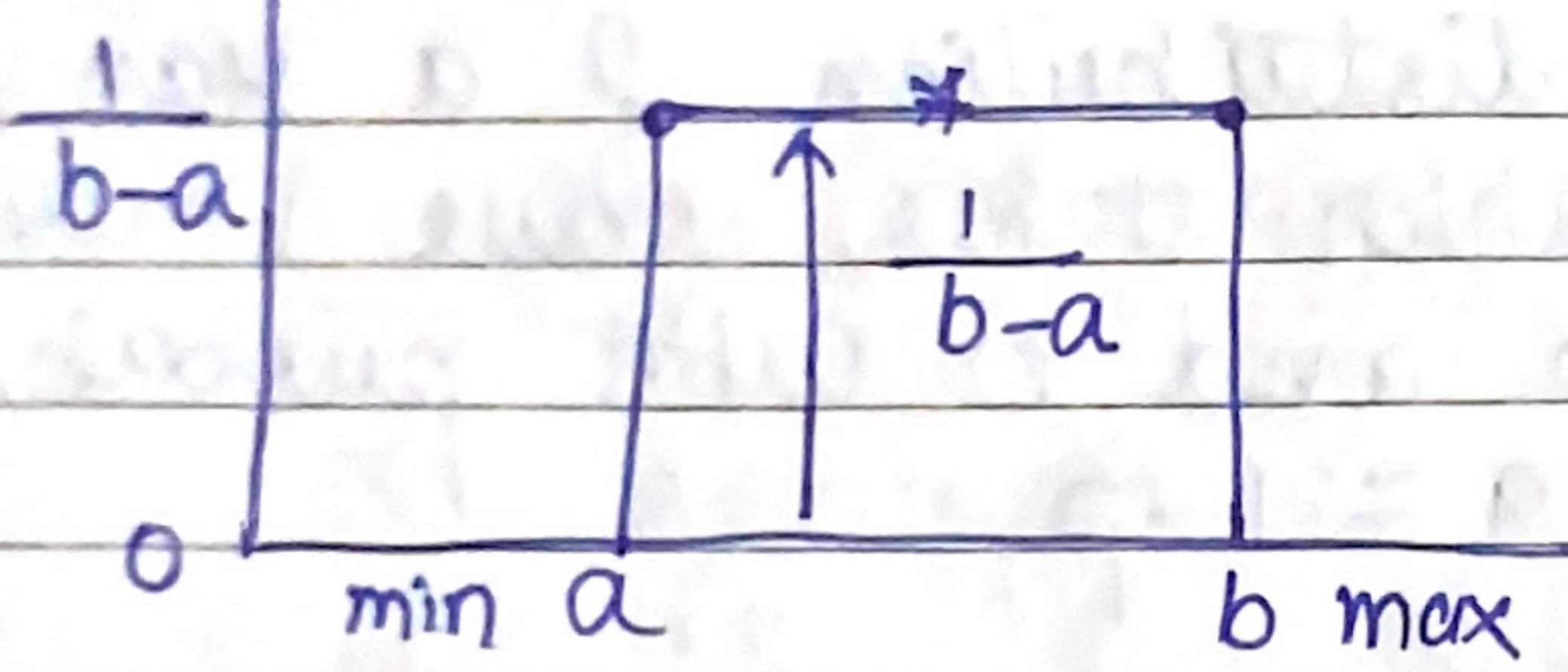
Mean  $a+b/2$

Skewness = 0

Median  $a+b/2$

Variance  $\frac{(b-a+1)^2}{12} = 1$

pdf

 $U(a, b)$  $-\infty < a < b < \infty$  $x \in [a, b]$ 

Pdf 
$$\begin{cases} \frac{1}{b-a} & x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

cdf 
$$\begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & x \in [a, b] \\ 1 & x \geq b \end{cases}$$

Mean  $\frac{1}{2}(a+b)$

Median  $\frac{1}{2}(a+b)$

## Bernoulli Distribution

- having only two outcomes
- probability distribution of a random variable which takes value 1 with probability  $p$  and 0 with probability  $q = 1-p$

$$0 < p < 1, p \in \mathbb{R}$$

pmf  $\begin{cases} q = (1-p) & k=0 \\ p & k=1 \end{cases}$

cdf  $\begin{cases} 0 & k < 0 \\ 1-p & 0 \leq k < 1 \\ 1 & k \geq 1 \end{cases}$

Mean  $p$

Variance  $p(1-p) = pq$

## Binomial Distribution

$Y \sim \text{Bin}(n, p) \rightarrow$  prob of getting head  
 ↓  
 no of trials

$$\text{B}(n, p)$$

$n \in \mathbb{N}_0$  — number of trials

$p \in [0, 1]$  — success probability in each trial

pmf  $\binom{n}{k} p^k (1-p)^{n-k}$

## log Normal Distribution

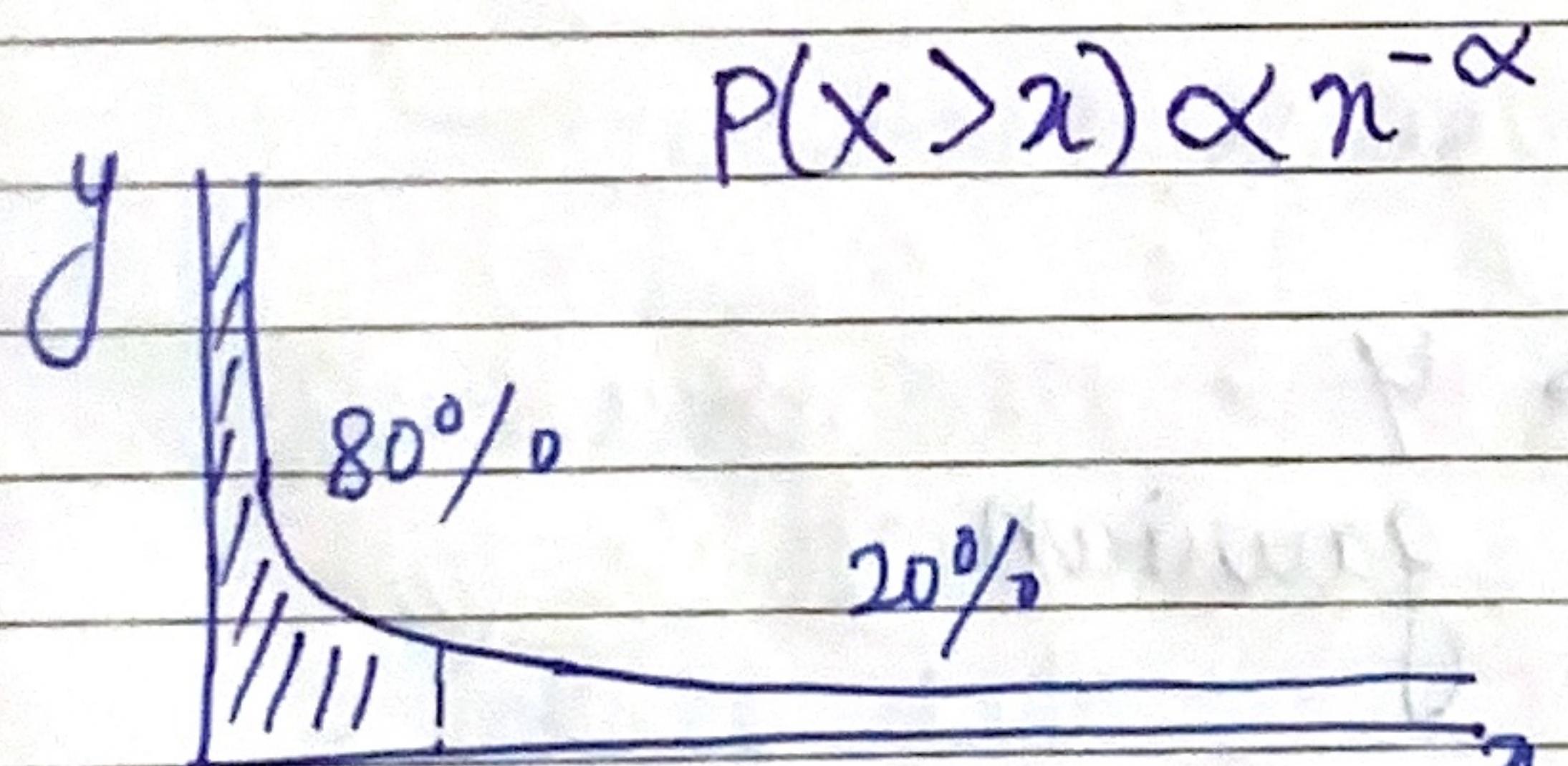
$X \sim \text{log normal}$

if  $\log(x)$  is normally distributed

log normal is a cont prob dist of random variable whose log is normally dist.

Thus if rand var  $x$  is log normally dist then  $y = \ln(x)$  has normal distribution

## Power law [called Pareto distribution]



[80-20 rule]

Power law is a relationship b/w two quantities where a relative change in one quantity results in proportional relative change in other quantity

$\therefore$  One quantity varies as power of another

log-log  $\rightarrow$  to determine whether dist follows power law

$\rightarrow$  Plot  $\log(n)$  vs  $\log(P(n))$

$\rightarrow$  straight line suggests power law

$$P(X > n) \propto n^{-\alpha}$$

- $P(X > n)$  is prob that variable  $X$  is greater than some value  $n$ .
- $\alpha > 1$  is power law exponent.
- implies that small events are common, large events are rare but not negligible.

## Power transform (Box-Cox transform)

$$\begin{array}{ccc} X & \xrightarrow{\ln} & Y \\ \downarrow & & \downarrow \\ \text{log normal} & & \text{gaussian} \end{array}$$

Power Law / pareto  $\xrightarrow{?}$  gaussian

Conversion  $\begin{cases} X: x_1, x_2, \dots, x_n \sim [\text{pareto}] \\ Y: y_1, y_2, \dots, y_n \sim [\text{gaussian}] \end{cases}$

① Box-Cox ( $x$ ) =  $\lambda^{\alpha}$   
 $(x_1, x_2, \dots, x_n)$

$$\textcircled{2} \quad y_i = \begin{cases} \frac{x_i - 1}{d} & \text{if } d \neq 0 \\ \lg(n_i) & \text{if } d = 0 \end{cases}$$

$A_i : 1 \rightarrow n$

## Application of non-gaussian dist

↳ 100s of dist  
↳ Why?

uniform  $\rightarrow$  random nos  
Bennelli, Binomial  
log normal, Pareto

## Weibull Distribution

↳ dams  $\rightarrow$  how high dam should be  
↳ water storage (max storage)

Civil engg:

Max one day rainfall  
Rainfall data  $\rightarrow$  30 yrs, 50 yrs

↳ few hundreds of data pts

$P(X > 20\text{cm})$  is very small  $\rightarrow$  before computation  
 $\leq 0.1\%$

$$P(X \geq 20\text{cm}) = ?$$

percentiles

$$\frac{1}{200} = 0.5\%$$

dist

100s of distb → gaussian - X  
 log normal - X  
 weibull - ✓

$x = 1 \rightarrow$  fit a known & well studied  
 $\quad \quad \quad 2 \quad$  distb / theoretical model  
 $\quad \quad \quad 3 \quad$

①  $x = \text{max daily rainfall}$

weibull (parameters) ②

③ CDF & PDF  $P(x \geq 20\text{cm}) = 1 - P(x < 20\text{cm})$   
 ↓  
 cdf of x

Bon-Cor Transform

$x \rightarrow x^{\frac{1}{n}}$  → all analysis  
 ↓  
 nongaussian → gaussian distb  
 max daily rainfall

## Relationship b/w Random Variables

X: heights  
 Y: weights

	X=h	Y=w
S <sub>1</sub>	160	62
S <sub>2</sub>	150	54
	1	1
S <sub>n</sub>	140	48

(Q) relation b/w X & Y

X↑, Y↑

X↑, Y↓

Covariance, Pearson correlation coeff,  
 Spearman rank corr. coeff

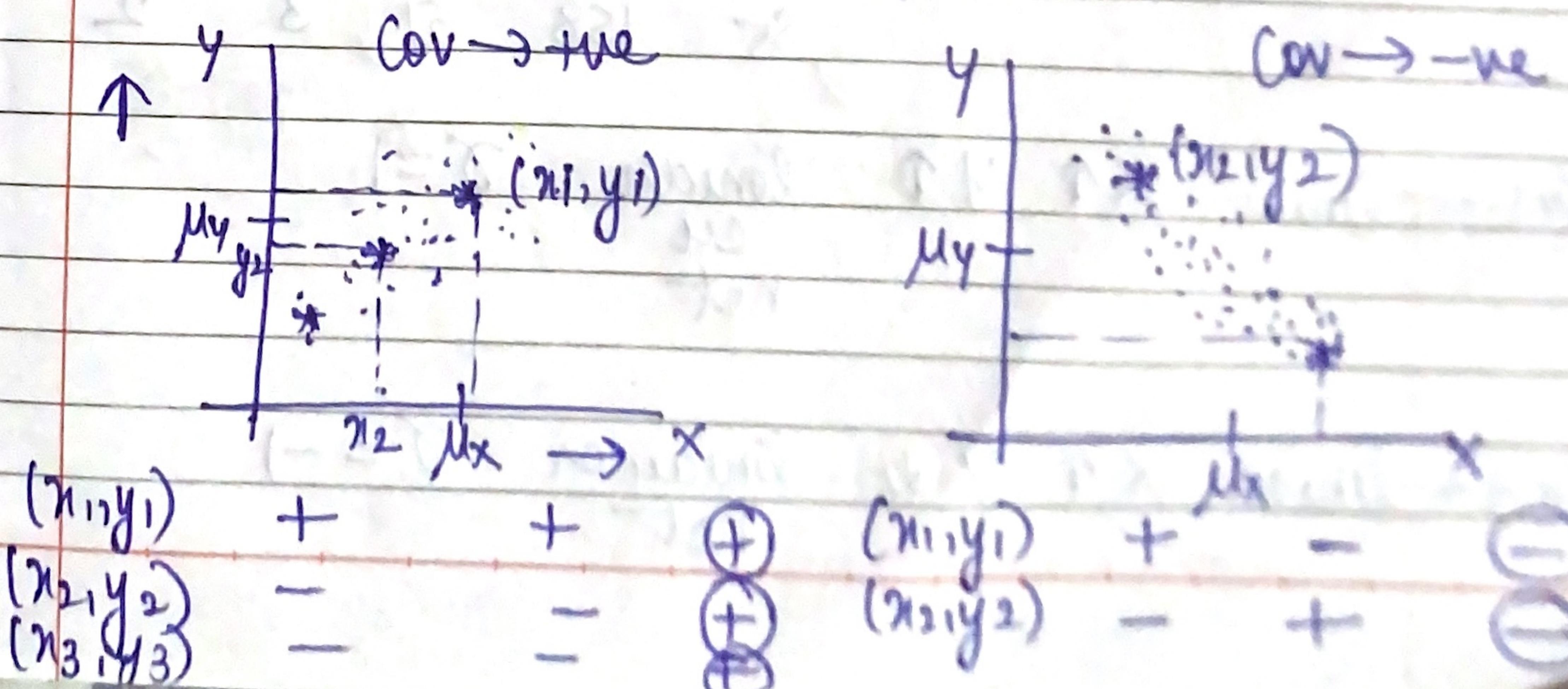
$$\text{Covariance } (X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x) * (y_i - \mu_y)$$

$$\text{Var}(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)^2$$

$$\text{cov}(X, X) = \text{Var}(X)$$

\*  $\text{cov}(X, Y) = \text{+ve}$       X↑, Y↑

\*  $\text{cov}(X, Y) = \text{-ve}$       X↑, Y↓



$\text{cov}(x, y)$

$\neq$

$\text{cov}(x_1, y_1)$   
feet lbs

### Pearson Correlation Coeff (PCC)

$$r_{x,y} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

$$\sigma_x = \sqrt{\text{var}(x)}$$

$$-1 \leq r \leq 1$$

Only studies linear relations

### Spearman Rank-Correl Coeff ( $\gamma$ )

$$\gamma = \rho_{\text{rank}}$$

	$\text{rank } x$	$\text{rank } y$
$s_1$	160	52
$s_2$	150	66
$s_3$	170	68
$s_4$	140	46
$s_5$	158	51

$\gamma = 1 \leftarrow$  linear  $X \uparrow Y \uparrow$  linear  $\gamma = 1$   
 or  
 not

$\gamma = -1 \leftarrow$  linear  $X \uparrow Y \downarrow$  linear or  $\gamma = -1$   
 not

Correlation does not imply causation

X causes Y } X  
Y causes X } X

Causal Models ↗ what causes what

Applications of Correlation → not causation

SRCC      PCC

- ① Is salary correlated with sq. feet of home?
- ② Is no of years of education correlated with income?
- ③ ecommerce: amazon  
Is time spent in 24 hrs → money spent in 24 hrs

# unique visitors in day Vs \$ salu in day

$$\begin{array}{ccc} 100K & \longrightarrow & \$1M \\ 120K & \longrightarrow & \$1.6M \\ | & & | \end{array}$$

- ④ Medicine  
(dosage of drug) (reduction in sugar)

1mg	z
2mg	y
3mg	x