

## Introduction to Statistics

Dfn: Statistics is the science of collecting, organizing and analyzing  
data ↓  
Decision Making Process

Data: "facts or pieces of information"

Eg: Height of students in classroom

{ 175cm, 180cm, 190cm, 160cm ... }

IQ

{ 100, 90, 95, 80 ... }

## ② Types of Statistics

### ① Descriptive Statistics

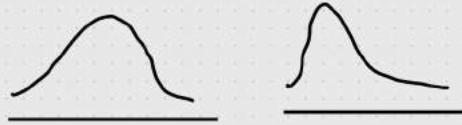
Defn: It consists of organizing And summarizing data

① Measure of Central Tendency [Mean, Median, Mode]

② Measure of Dispersion [Variance, Std]

③ Different type of Distribution of data.

Eg: Histogram, pdf, pmf



Eg: Let's say there are 20 statistics classes at your college. And you have collected the heights of students in the class.

Heights are recorded [175cm, 180cm, 140, 140, 135, 160, 135, 190cm]

### Descriptive Question

$$\frac{175 + 180 + 140 + 140 + 135 + 160 + 135 + 190}{8} = \text{Avg Height}$$

"What is the average height of the entire classroom"

### Inferential Question

↗ Sample

"Are the height of the students in classroom similar to what you expect in the entire college"

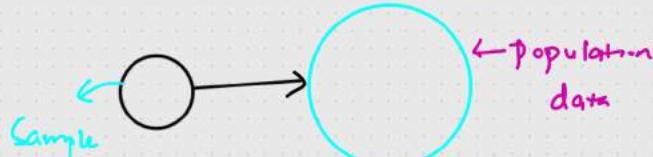
↳ population data

## Population And Sample Data

Population: The group you are interested in studying

### ② Inferential Statistics

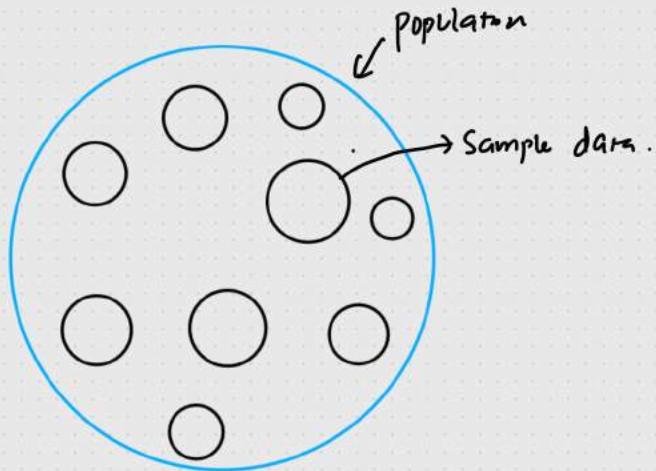
Defn: It consists of using data you have measured to form conclusion



① Z-test      ② t-test      ③ CHI SQUARE } Hypothesis Testing  
Mo, H, p value, Significance value

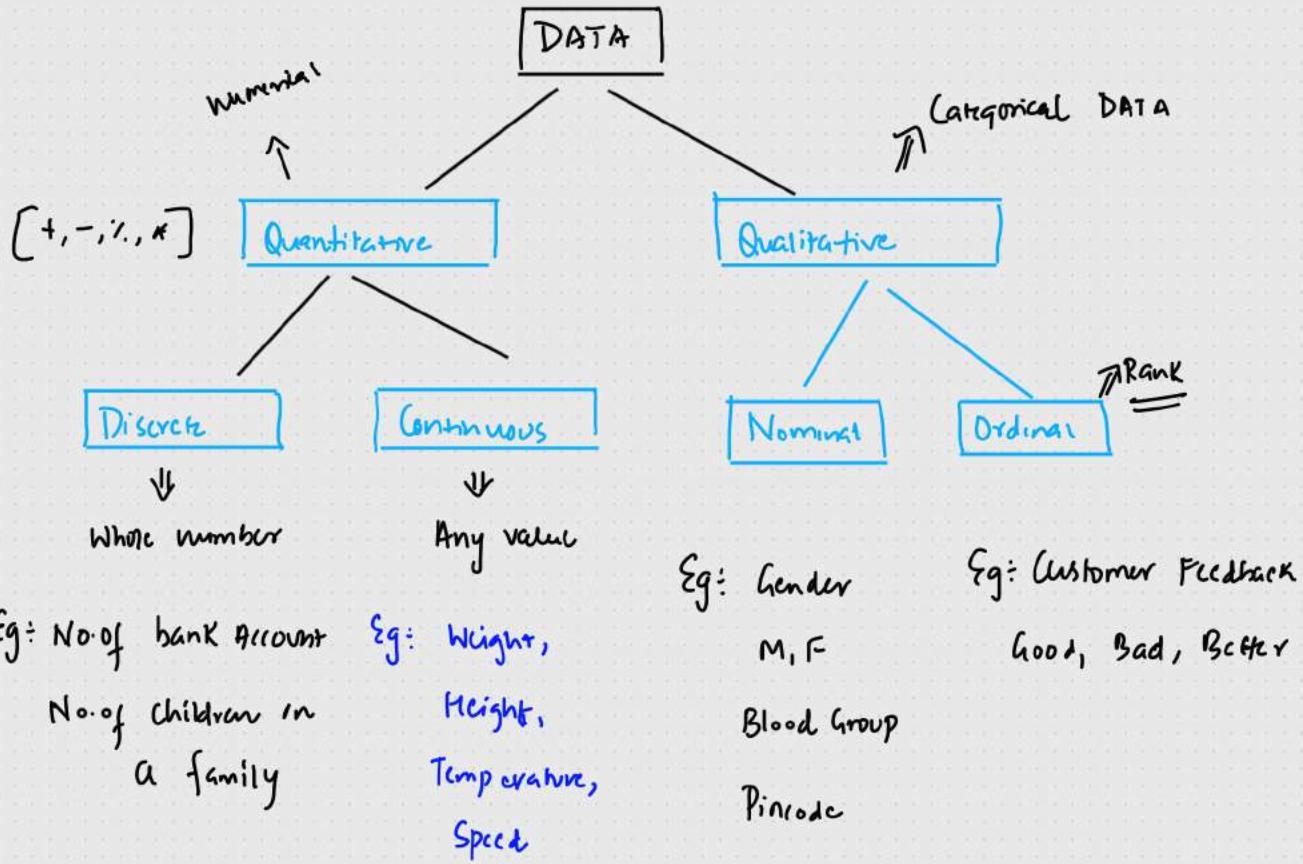
Sample : a subset of population

Eg:  $\frac{\text{Total Population}}{\downarrow}$



### ③ Types of Data

| DMC  | DC   | ISI | BUI  | FWI | Classes  | Region |
|------|------|-----|------|-----|----------|--------|
| 3.4  | 7.6  | 1.3 | 3.4  | 0.5 | not fire | 0      |
| 4.1  | 7.6  | 1   | 3.9  | 0.4 | not fire | 0      |
| 2.5  | 7.1  | 0.3 | 2.7  | 0.1 | not fire | 0      |
| 1.3  | 6.9  | 0   | 1.7  | 0   | not fire | 0      |
| 3    | 14.2 | 1.2 | 3.9  | 0.5 | not fire | 0      |
| 5.8  | 22.2 | 3.1 | 7    | 2.5 | fire     | 0      |
| 9.9  | 30.5 | 6.4 | 10.9 | 7.2 | fire     | 0      |
| 12.1 | 38.3 | 5.6 | 13.5 | 7.1 | fire     | 0      |



### ④ Scale of Measurement

- ① Nominal Scale Data
- ② Ordinal Scale Data
- ③ Interval Scale Data
- ④ Ratio. Scale Data.

## ① Nominal Scale Data

- Qualitative / Categorical
- Eg: Gender, Colors, Habits
- Order does not matter

Eg: Favorite color

Red → 5 → 50%  
 Blue → 3 → 30%  
 Orange → 2 → 20%  
 $\frac{10}{10}$

Gender

M

F

## ② Ordinal Scale Data

- Ranking is important
- Order matter
- Difference cannot be measured

Eg: 1 → Best

2 → Good

3 → Bad

Working Professional

M

T

W

Th

← F

Race: 1<sup>st</sup> [ 4:20  
 2<sup>nd</sup> [ 5:30  
 3<sup>rd</sup> ] 6:00

OF

## ③ Interval Scale Data

- The order matter
- Difference can be measured
- Ratio cannot be measured ✓
- No "0" Starting Point

Eg: Temperature Variable

-30° F

-15° F

[ 30 F ] 60 : 30 = 2 : 1 →  
 [ 60 F ] 90 - 60 = 30 F  
 [ 90 F ] 120 - 90 = 30 F

## ④ Ratio Scale Data

- The order matter ✓
- Differences are measurable (including ratio)
- Contains a "0" Starting point.

Eg: Students marks in a class

0, 90, 60, 30, 75, 40, 50

Avg = 30, 40, 50, 60, 75, 90

$$40 - 30 = 10$$

$$50 - 30 = 20$$

$$\text{Ratio} = \frac{90}{30} = 3 : 1$$

### Example

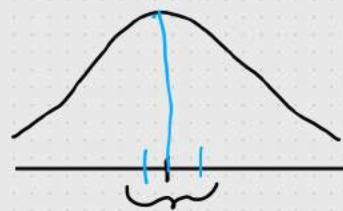
- ① length of Different Rivers In the World ?
- ② favorite food based on Gender ?
- ③ Marital Status ?
- ④ IQ measurement ?

## ① Measure of Central Tendency

① Mean or Average

② Median

③ Mode



### ① Mean

Population ( $N$ )

$$X = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$$

Sample ( $n$ )

$$\text{Population mean } (\mu) = \frac{\sum_{i=1}^n x_i}{N} \quad (\bar{x}) \quad \text{Sample mean} = \frac{\sum_{i=1}^n x_i}{n}$$

$$= \left[ \frac{1+1+2+2+3+3+4+5+5+6}{10} \right] =$$

$$= \frac{32}{10} = 3.2$$

### ② Median

Steps

$$X = \{4, 5, 2, 3, 2, 1\}.$$

① Sort the Random Variable  $\{1, 2, 2, 3, 4, 5\}$ .

② No. of elements Count = 6

③ if Count == even

$$\{1, 2, \boxed{2, 3}, 4, 5\}$$



$$\text{Median} = \frac{2+3}{2} = 2.5$$

④ if Count is odd

$$\{1, 2, 2, \boxed{3}, 4, 5, 6\}$$

$$\text{Median} = \underline{\underline{3}}$$

Why Median?

$$X = \{1, 2, 3, 4, 5\}$$

$$X = \{1, 2, 3, 4, 5, 100\}$$

Outlier  
↓

$$\bar{x} = \frac{1+2+3+4+5}{5}$$

$$\bar{x} = \frac{1+2+3+4+5+100}{6}$$

$$\bar{x} = \frac{15}{5} = 3 \quad \xrightarrow{\quad} \quad \frac{115}{6} \approx 19.$$

$$X = \{1, 2, \boxed{3, 4}, 5, 100\}$$

$$\text{Median} = \frac{3+4}{2} = 3.5$$

Median is used to find the central Tendency

When outliers is present.

③ Mode : Frequency Maximum

$$\{2, 1, 1, 1, 4, 5, 7, 8, 9, 9, 10\}$$

Modc = 1

## EEDA AND Feature Engineering

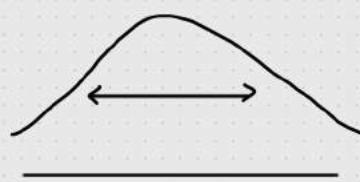
| Age  | Weight | Salary | Gender | Degree |
|------|--------|--------|--------|--------|
| 24   | 70     | 40K    | M      | BE     |
| 25   | 80     | 70K    | F      | -      |
| 21   | 95     | 45K    | F      | -      |
| → 24 | [-]    | 50K    | M      | PHD    |
| → 32 | -      | 60K    | -      | BE     |
| → -  | 60     | -      | -      | Master |
| → -  | 65     | 55K    | -      | BSC    |
| → 40 | 72     | -      | M      | BE     |

DATA  
IS  
MISSING

## ① Measure of Dispersion [Spread of the data]

① Variance

② Standard deviation



### ① Variance

Population Variance

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

$x_i \rightarrow$  DATA POINTS

$\mu \rightarrow$  Population mean

$N \rightarrow$  Population size

Sample Variance

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

$x_i \rightarrow$  DATA POINTS

$\bar{x} \rightarrow$  Sample Mean

$n \rightarrow$  Sample Size

Why we divide Sample Variance by  $n-1$ ?

Ans) The sample variance is divided by  $n-1$  so that  
 we can create an unbiased estimator of the  
 population variance

Bernie's Correction  
 ↑

Eg:  $\{1, 2, 3, 4, 5\}$

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

| $x$           | $\bar{x}$ | $(x_i - \bar{x})^2$ | $n=5$ |
|---------------|-----------|---------------------|-------|
| 1             | 3         | 4                   |       |
| 2             | 3         | 1                   |       |
| 3             | 3         | 0                   |       |
| 4             | 3         | 1                   |       |
| $\frac{5}{3}$ | 3         | $\frac{4}{10}$      |       |

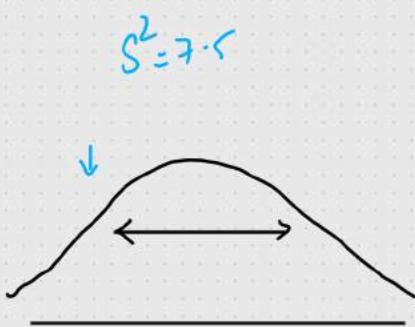
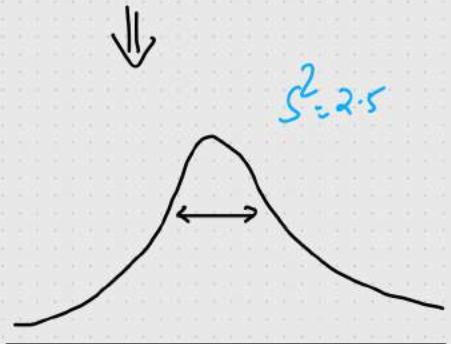
$$S^2 = \frac{10}{42} = 2.5$$

$$X = \{ \quad \} \quad Y = \{ \quad \}$$

$$S^2 = 2.5$$

$$S^2 = 7.5$$

Dispersion or Spread



## ② Standard Deviation

Population Standard Deviation

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

Sample Std

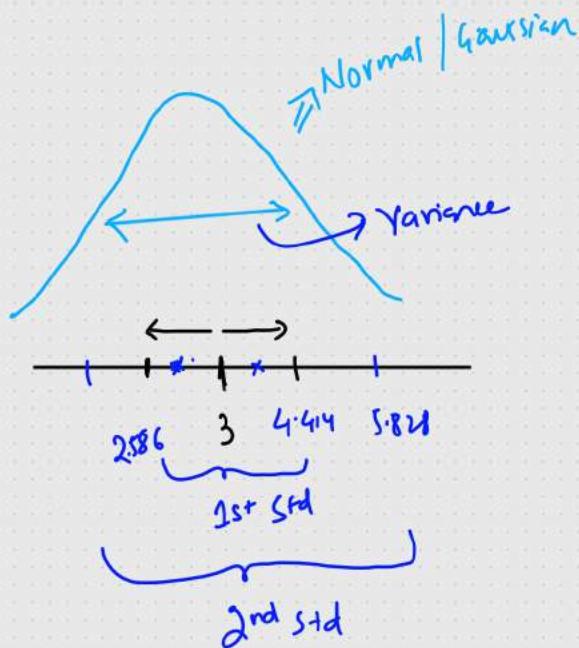
$$\text{Std} = \sqrt{s^2}$$

$$s^2 = \text{Sample Variance}$$

$$X = \{1, 2, 3, 4, 5\}$$

$$\bar{x} = 3$$

$$\sigma = 1.414$$



$$\begin{array}{r} 3.000 \\ 1.414 \\ \hline 2.586 \end{array}$$

$$\begin{array}{r} 4.414 \\ 1.414 \\ \hline 5.828 \end{array}$$

## ① Random Variable

$$x + 5 = 7$$

$x, y$  & variables?

$$x = 2$$

$$8 = y + x$$

$$y = 6$$

Random Variable is a process of mapping the output of a random process or experiments to a number.

Eg: Tossing a Coin

$$X = \begin{cases} 0 & \text{if Head} \\ 1 & \text{if Tails} \end{cases}$$

Rolling a dice

Measure the Temperature of the next day

$$Y = \begin{cases} \text{Sum of the rolling of dice 7 times} \\ \downarrow \end{cases}$$

$$Pr(Y > 15) \quad P(Y < 10)$$

$$X = \{-, -, -, -, -\} \Rightarrow \text{Ages}$$

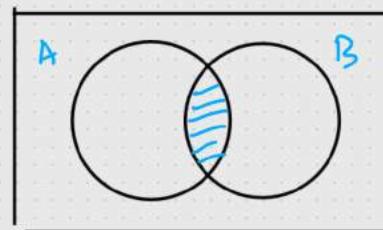
① Sets

$$A = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

$$B = \{3, 4, 5, 6, 7\}$$

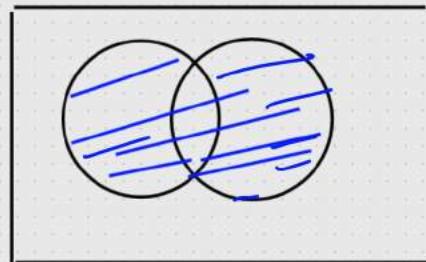
① Intersection

$$A \cap B = \{3, 4, 5, 6, 7\}.$$



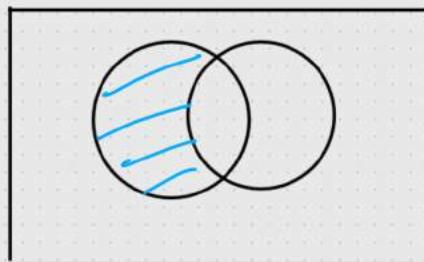
② Union

$$A \cup B = \{1, 2, 3, 4, 5, 6, 7, 8\}$$



③ Difference

$$A - B = \{1, 2, 8\}$$



④ Subset

$$A \rightarrow B \Rightarrow \text{False}$$

$$B \rightarrow A \Rightarrow \text{True}$$

## ⑤ Superset

$A \rightarrow B \Rightarrow \text{True}$

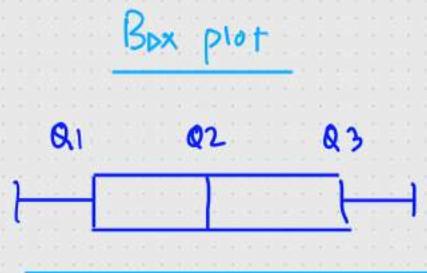
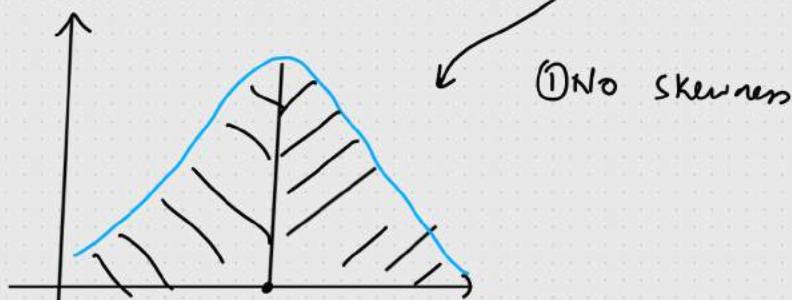
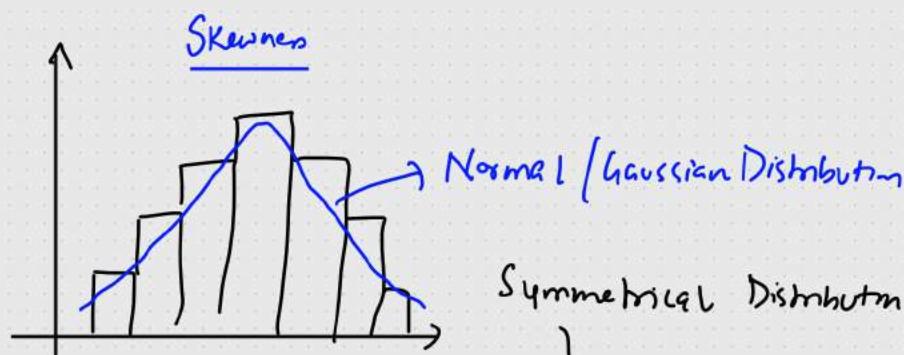
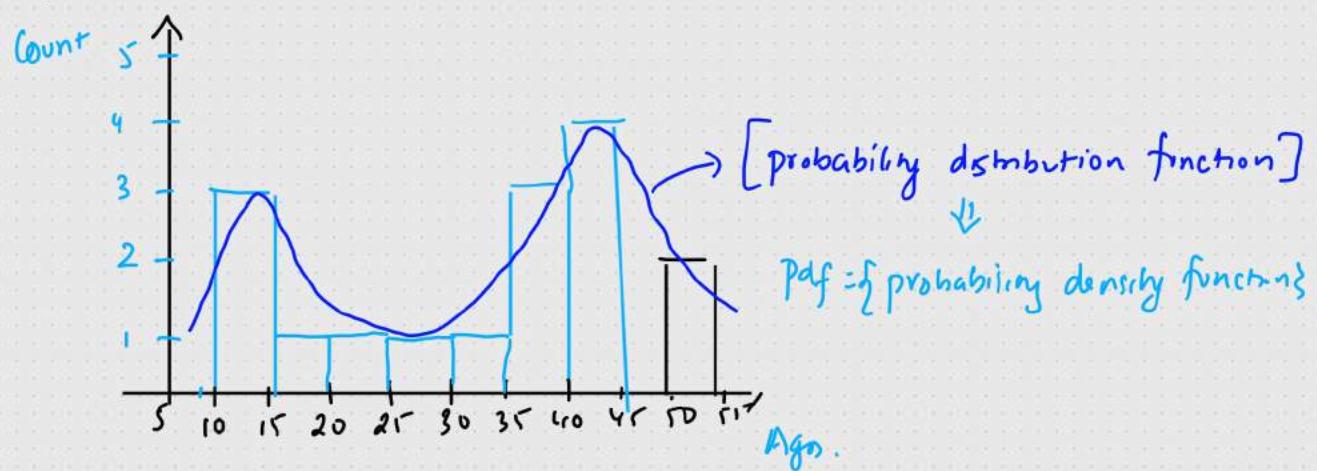
$B \rightarrow A \Rightarrow \text{False}$

## Histograms And Skewness → [Frequency]

Age = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 48, 50, 51}

$$\frac{50}{10} = \boxed{5} \rightarrow \text{bin size} \quad \left\{ \text{No. of bins} = 10 \right\}$$

$$\frac{50}{2.5} = 20 \rightarrow \text{bin size} \quad \left\{ \text{No. of bins} = 20 \right\}$$

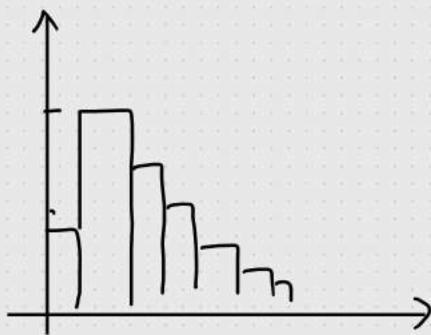


The mean, median, and mode all are perfectly at the center

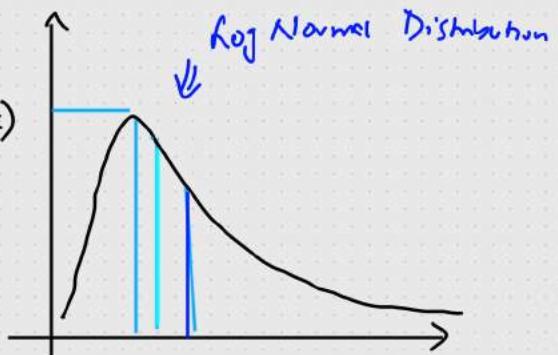
$$Q_3 - Q_2 \approx Q_2 - Q_1$$

mean = median = mode

② Right skewed



$\Rightarrow$  Positive Skewed  $\Rightarrow$



mean > median > mode

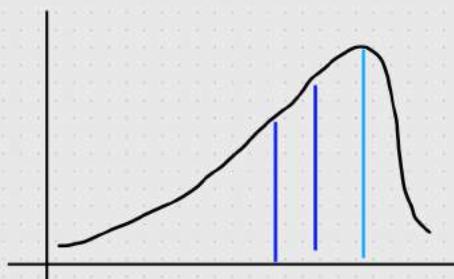
Relationship between Mean, Median, Mode

Box plot

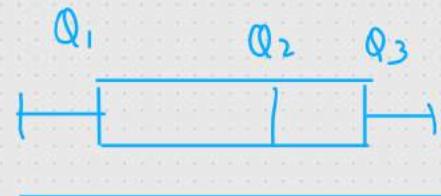


$$Q_3 - Q_2 \geq Q_2 - Q_1$$

③ Left Skewed Distribution



$\Rightarrow$  Negative Skewed



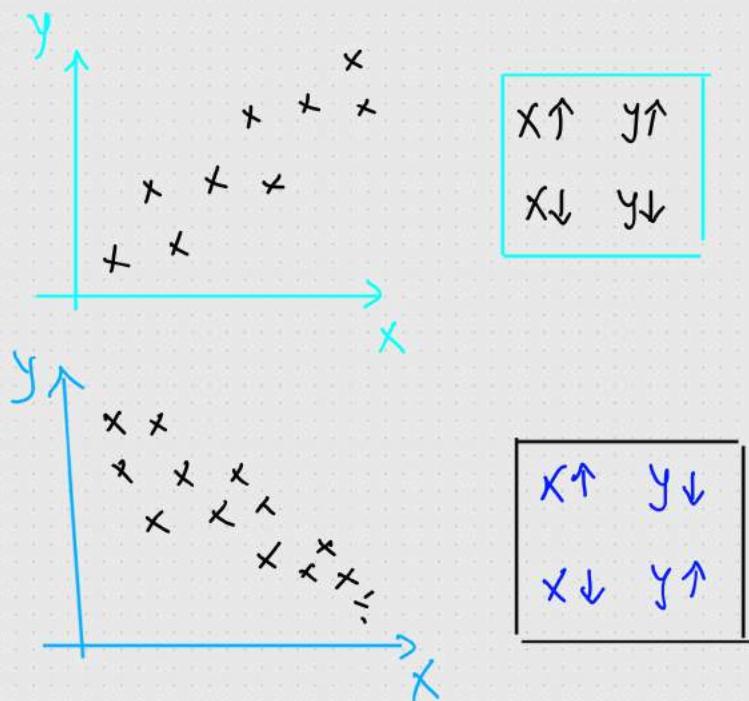
$$Q_2 - Q_1 \geq Q_3 - Q_2$$

Relationship : mean  $\leq$  median  $\leq$  mode

## Covariance And Correlation

[Relationship between X and Y]

| X | Y |                                   |  |
|---|---|-----------------------------------|--|
| 2 | 3 | $X \uparrow \quad Y \uparrow$     | $\downarrow \uparrow \text{Size of mouse} \rightarrow \text{Price of house}$ |
| 4 | 5 | $X \downarrow \quad Y \uparrow$   |  |
| 6 | 7 | $X \downarrow \quad Y \downarrow$ |  |
| 8 | 9 | $X \uparrow \quad Y \downarrow$   | $\downarrow \uparrow \text{Price of house}$                                  |



### Covariance

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$\begin{aligned} \text{Var}(x) &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x})(\bar{x} - x_i)}{n-1} \end{aligned}$$

$\approx \text{Cov}(x, x) \Rightarrow \text{Spread}$

$x_i \rightarrow$  Data points of x

$\bar{x} \rightarrow$  Sample mean of x

$y_i \rightarrow$  Data points of y

$\bar{y} \rightarrow$  Sample mean of y

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

|                |                |
|----------------|----------------|
| $x \uparrow$   | $y \uparrow$   |
| $x \downarrow$ | $y \downarrow$ |

+ve Covariance

|                |                |
|----------------|----------------|
| $x \uparrow$   | $y \downarrow$ |
| $x \downarrow$ | $y \uparrow$   |

-ve Covariance

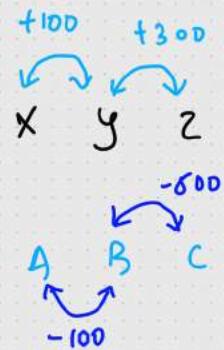
|                 |               |
|-----------------|---------------|
| $X$             | $y$           |
| $\rightarrow 2$ | 3             |
| $\rightarrow 4$ | 5             |
| $\rightarrow 6$ | 7             |
| $\bar{x} = 4$   | $\bar{y} = 5$ |

$$\text{Cov}(x,y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$= \frac{[(2-4)(3-4) + (4-4)(5-4) + (6-4)(7-4)]}{n-1}$$

$$= \frac{4+0+4}{2} = \frac{8}{2} = 4 \text{ tve value}$$

Positive  
Covariance



$x$  &  $y$  are having a positive Covariance

### Advantages

- ① Relationship between  $x$  and  $y$   
+ve or -ve value

### Disadvantages

- ① Covariance does not have a specific limit value

## ② Pearson Correlation Coefficient [-1 to 1]

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

① The more the value towards +1 the more +ve correlated it is  $(x,y)$

② The more the value towards -1 the more -ve correlated it is  $(x,y)$

## ③ Spearman Rank Correlation [-1 to 1]

$$\gamma_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) \cdot \sigma(R(y))}$$

| <u>x</u> | <u>y</u> | <u><math>R(x)</math></u> | <u><math>R(y)</math></u> |
|----------|----------|--------------------------|--------------------------|
| 1        | 2        | 5                        | 5                        |
| 3        | 4        | 4                        | 4                        |
| 5        | 6        | 3                        | 3                        |
| 7        | 8        | 2                        | 1                        |
| 0        | 7        | 6                        | 2                        |
| 8        | 1        | 1                        | 6                        |

## Feature Screen

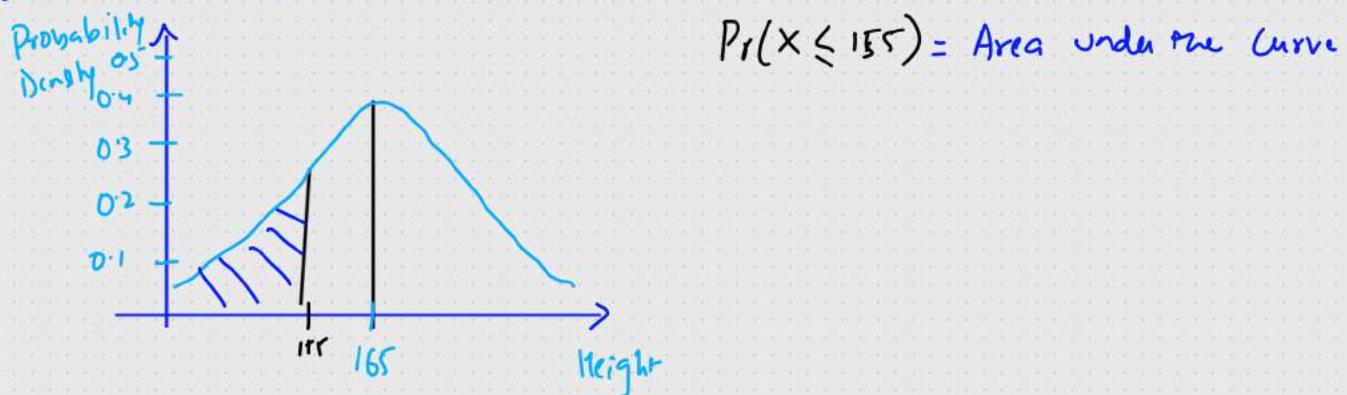
|                        |                       |                   |   |                |         |
|------------------------|-----------------------|-------------------|---|----------------|---------|
| +ve<br>Size of House ↑ | +ve<br>No. of Rooms ↑ | +ve<br>Location ↑ | ≈ 0<br><del>No. of people staying</del> | -ve<br>Haunted | Price ↑ |
|------------------------|-----------------------|-------------------|---|----------------|---------|

# Probability Distribution Function / Density Function (pdf) (pmf)

## ① Probability Density Function (pdf) :

### ① Continuous Random Variable

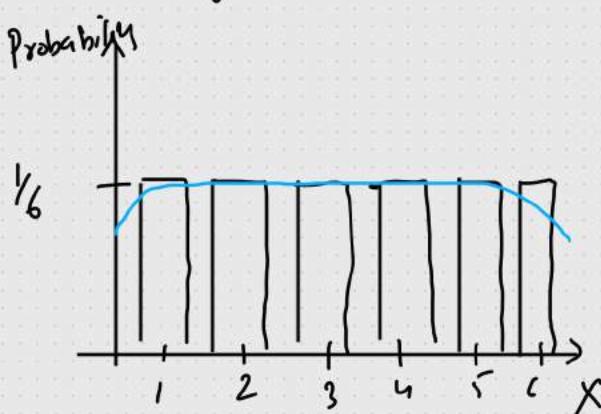
Eg: Heights of students in classroom [0 - 1]



## ② Probability Mass Function (pmf)

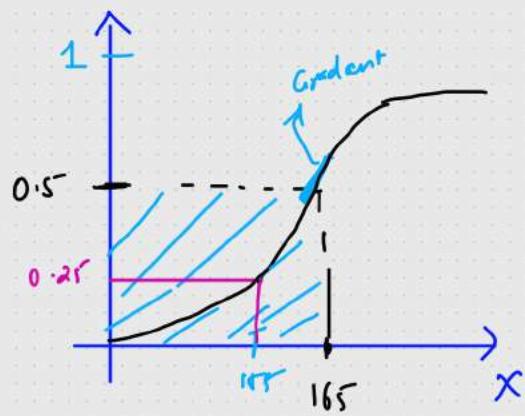
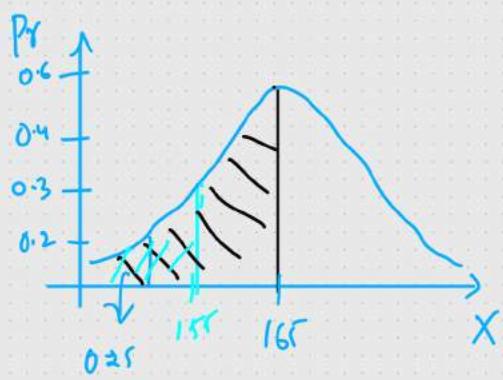
Variable  $\rightarrow$  Discrete Random Variable

Eg: Rolling a Dice  $\{1, 2, 3, 4, 5, 6\}$



$$\begin{aligned} \Pr(X \leq 4) &= \Pr(X=1) + \Pr(X=2) \\ &\quad + \Pr(X=3) + \Pr(X=4) \\ \Pr(1) &= \frac{1}{6} \\ \Pr(2) &= \frac{1}{6} \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{4}{6} = \frac{2}{3} \end{aligned}$$

### ③ Cumulative Distribution Function (cdf) Cumulative probability



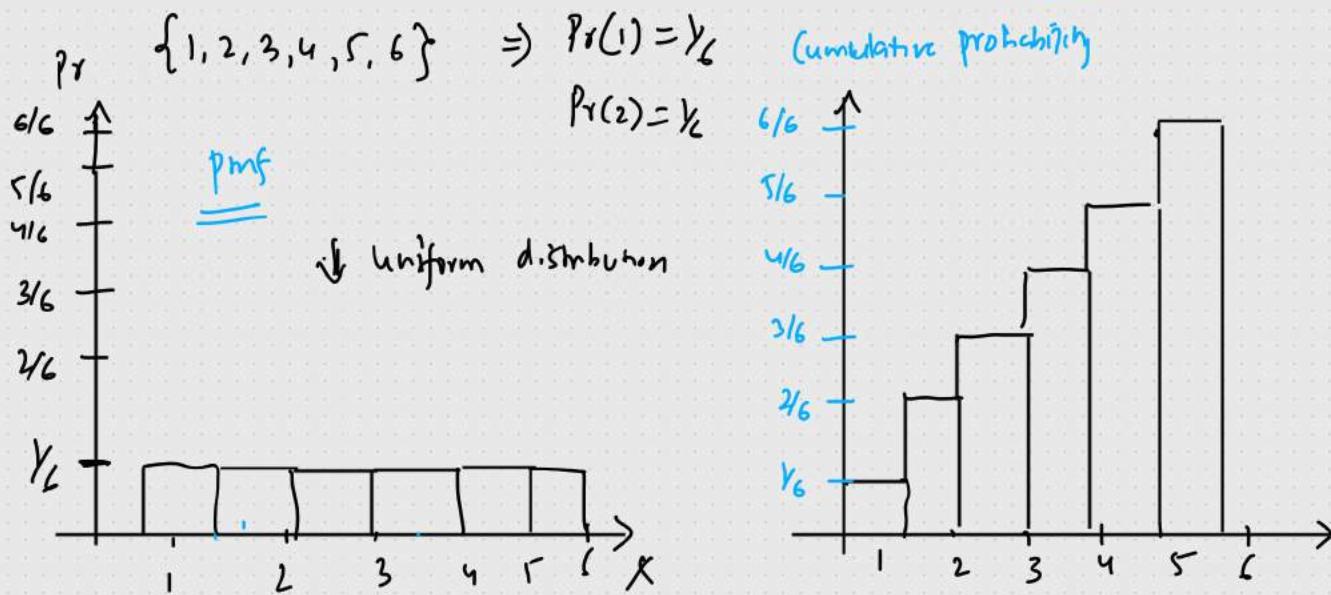
# Probability Density Function & Probability Mass Function

## Cumulative Distribution Function

### ① PMF

#### ① Discrete Random Variable

Eg: Rolling a dice

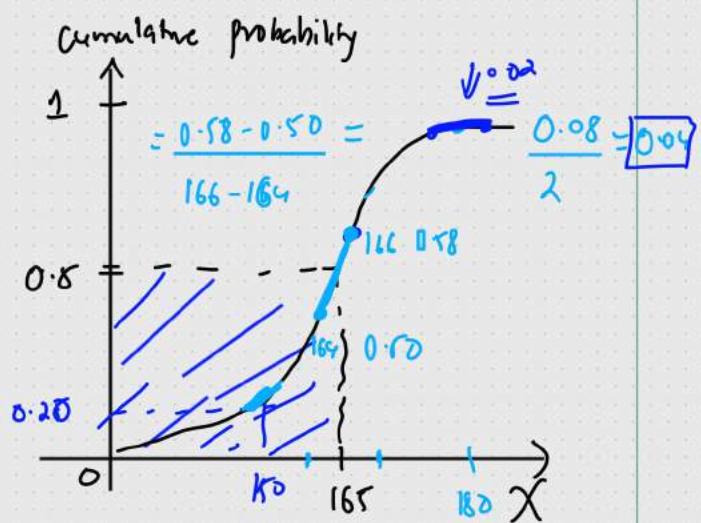
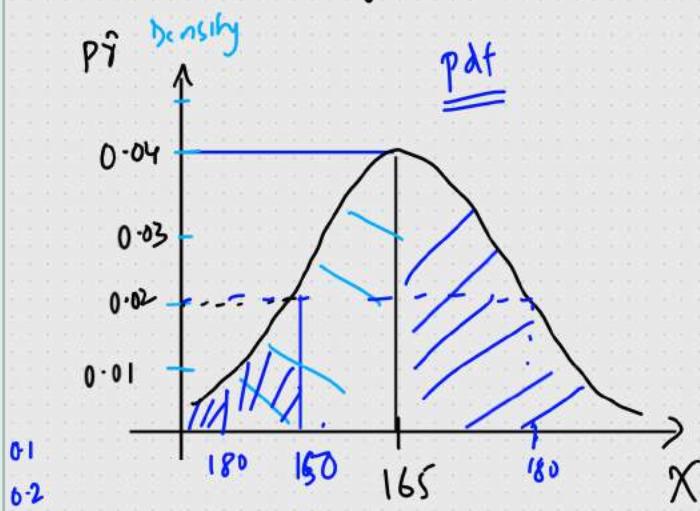


$$\begin{aligned} Pr(X \leq 2) &= Pr(X=1) + Pr(X=2) \\ &= \frac{1}{6} + \frac{1}{6} = \frac{2}{6} \end{aligned}$$

$$Pr(X \leq 6) = 1$$

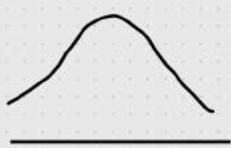
## ② Probability Density Function (pdf)

### ① Distribution of Continuous Random Variable



Probability Density  $\rightarrow$  Gradient of Cumulative Curve

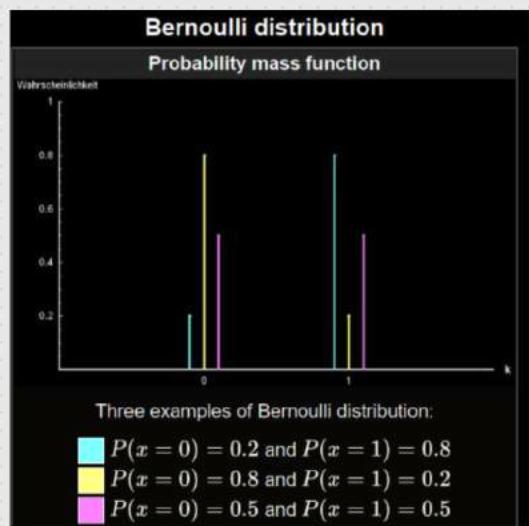
## Types of Probability Distribution



- ① Normal / Gaussian Distribution (pdf)
- ② Bernoulli Distribution (pmf) → Outcomes are binary only
- ③ Uniform Distribution (pmf)
- ④ Poisson Distribution (pmf)
- ⑤ Log Normal Distribution (pdf)
- ⑥ Binomial Distribution (pmf)

# Bernoulli distribution

In probability theory and statistics, the Bernoulli distribution, named after Swiss mathematician Jacob Bernoulli, is the discrete probability distribution of a random variable which takes the value 1 with probability p and the value 0 with probability  $q=1-p$ . Less formally, it can be thought of as a model for the set of possible outcomes of any single experiment that asks a yes–no question.



① Discrete Random Variable {pmf}

② Outcomes are Binary

Eg: ① Tossing a coin {H,T}

$$Pr(H) = 0.5 = p$$

$$Pr(T) = 0.5 = 1-p = q$$

② Whether the person will Pass/Fail

$$Pr(Pass) = 0.7 = p$$

$$Pr(Fail) = 1 - 0.7 = 0.3 = q.$$

## Parameters

$$0 \leq p \leq 1$$

$$q = 1-p$$

$$K = \{0, 1\} \rightarrow 2 \text{ outcomes}$$

## ① PMF

$$\text{PMF} = p^k * (1-p)^{1-k} \quad k \in \{0, 1\}$$

$$\text{if } K=1$$

$$Pr(K=1) = p^1 (1-p)^{1-1} \\ = p$$

Simplified

$$\text{pmf} \quad \begin{cases} q = 1-p & \text{if } K=0 \\ p & \text{if } K=1 \end{cases}$$

if  $K=0$

$$\Pr(K=0) = p^0 + (1-p)^{1-0}$$
$$= (1-p) = q_{11}.$$

② Mean of Bernoulli Distribution

$$E(K) = \sum_{k=0}^K k \cdot p(k)$$
$$\Pr(K=1) = 0.6 \Rightarrow p$$
$$\Pr(K=0) = 0.4 \Rightarrow 1-p$$
$$= [0 * 0.4 + 1 * 0.6]$$
$$= 0.6 = p_{11}.$$

③ Median of Bernoulli Distribution

Median

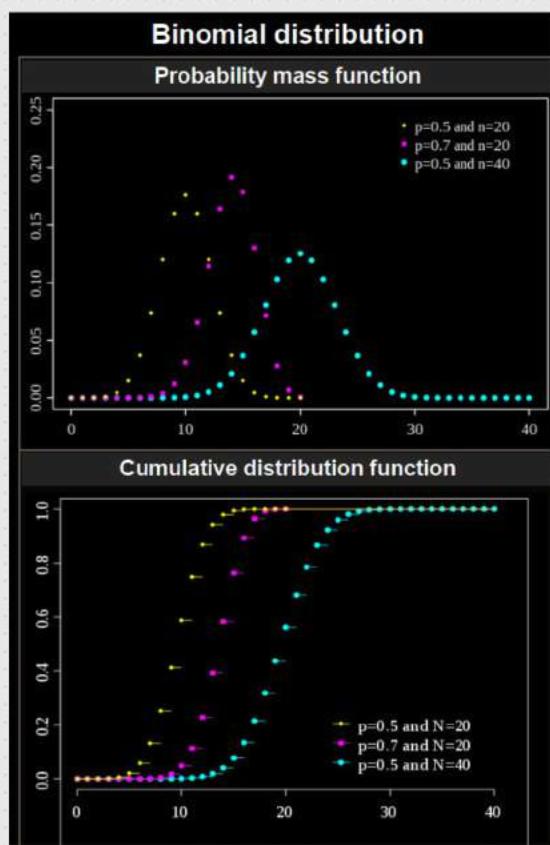
$$\begin{cases} 0 & \text{if } p < \frac{1}{2} \\ [0,1] & \text{if } p = \frac{1}{2} \\ 1 & \text{if } p > \frac{1}{2} \end{cases}$$

④ Variance Std

$$\text{Var} = p * (1-p)$$
$$= pq_{11}$$
$$\text{Std} = \sqrt{pq_{11}}.$$

# Binomial Distribution

In probability theory and statistics, the binomial distribution with parameters  $n$  and  $p$  is the discrete probability distribution of the number of successes in a sequence of  $n$  independent experiments, each asking a yes–no question, and each with its own Boolean-valued outcome: success (with probability  $p$ ) or failure  $q=1-p$ . A single success/failure experiment is also called a Bernoulli trial or Bernoulli experiment, and a sequence of outcomes is called a Bernoulli process; for a single trial, i.e.,  $n = 1$ , the binomial distribution is a Bernoulli distribution. The binomial distribution is the basis for the popular binomial test of statistical significance.



① Discrete Random Variable:

① Every experiment outcome is binary

② These experiment is performed  
for  $n$  times.

Eg: Tossing a coin 10 times

Notation  $\hat{=} \text{B}(n, p)$

Parameters  $\hat{=} n \in \{0, 1, 2, \dots\} \rightarrow$  no. of trials

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

$$q = 1-p$$

Support  $\hat{=} K \in \{0, 1, 2, \dots, n\} \rightarrow$  Number of success

PMF :

$$Pr(X, n, p) = {}^n C_k p^k (1-p)^{n-k}$$

for  $k=0, 1, 2, \dots, n$  where

$$\boxed{{}^n C_k = \frac{n!}{k!(n-k)!}}$$

Mean

$$\text{Mean} = np$$

Variance

$$\text{Var} = npq$$

$$\text{Std} = \sqrt{npq}.$$

## Poisson Distribution

① Discrete Random Variable (pmf)

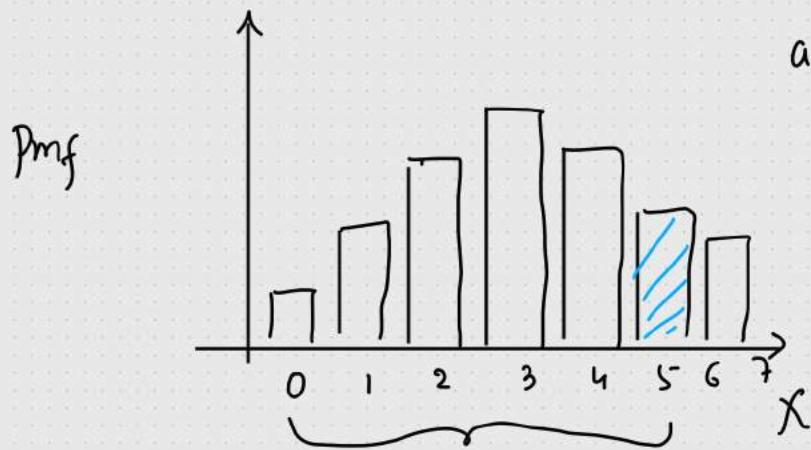
② Describes the number of events occurring in a fixed time interval

Eg: No. of people visiting hospital every hour

No. of people visiting banks every hour

$\lambda=3 \Rightarrow$  Expected no. of event occur

at every time interval



PMF

$$P(X=5) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$\lambda=3$

$$= \frac{e^{-3} 3^5}{5!}$$

$\lambda$

$$P(X=4) + P(X=5)$$

## Mean of Poisson Distribution

$$\text{Mean} = E(X) = \mu = \lambda * t$$

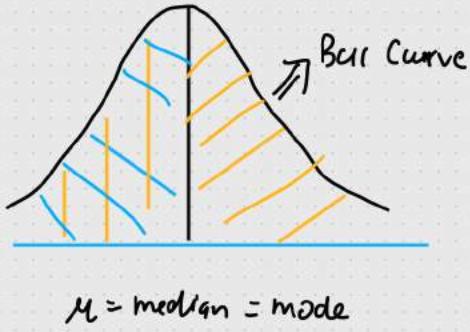
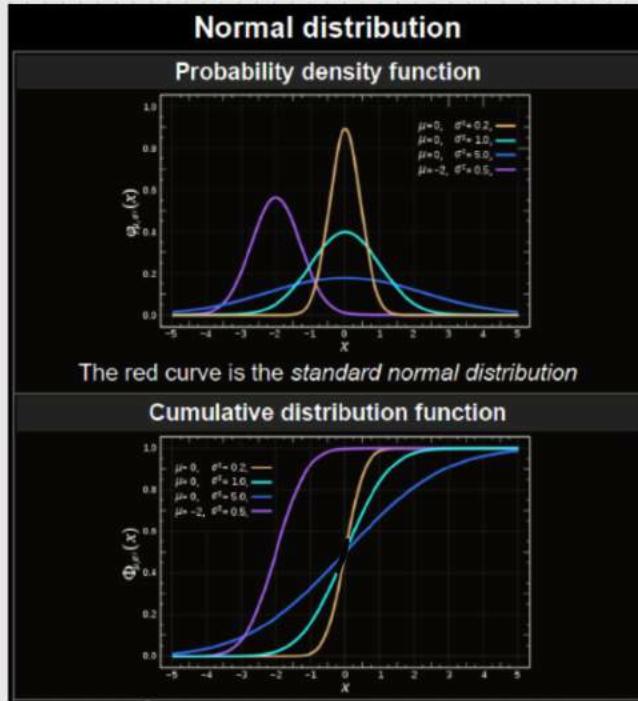
$$\text{Variance} =$$

$\lambda$  = Expected No. of events to occur  
at every time interval

$t$  = Time interval

# Normal/Gaussian Distribution (pdf)

In statistics, a normal distribution or Gaussian distribution is a type of continuous probability distribution for a real-valued random variable.



Notation  $N(\mu, \sigma^2)$

Parameters :  $\mu \in \mathbb{R}$  = mean  
 $\sigma^2 \in \mathbb{R} > 0$  = variance

$x \in \mathbb{R}$

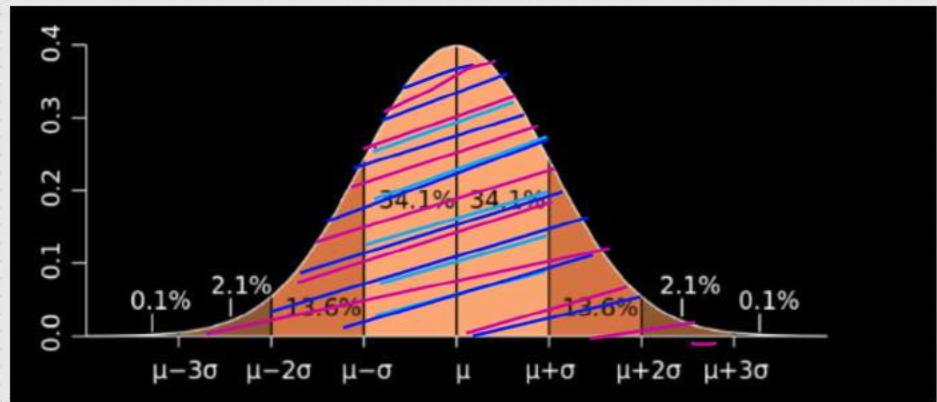
$$\text{PDF} = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Mean of Normal Distribution

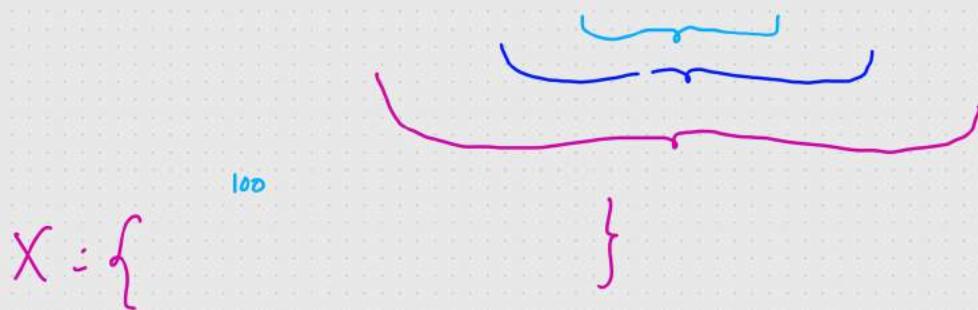
Mean =  $\mu$  = Average

Variance & Std  $\text{Var} = \sigma^2$   $\sigma = \sqrt{\text{Var}}$

## Empirical Rule of Normal Distribution



68-95-99.7%  
=



## Probability

$$\Pr(\mu - \sigma \leq X \leq \mu + \sigma) \approx 68\%$$

$$\Pr(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \approx 95\%$$

$$\Pr(\mu - 3\sigma \leq X \leq \mu + 3\sigma) \approx 99.7\%$$

- Example : ① Weight of the student in the class  
② Height of the " " " "  
③ IRIS DATASET {sepal width}

## Q-Q plot { Quantile Quantile Plot }

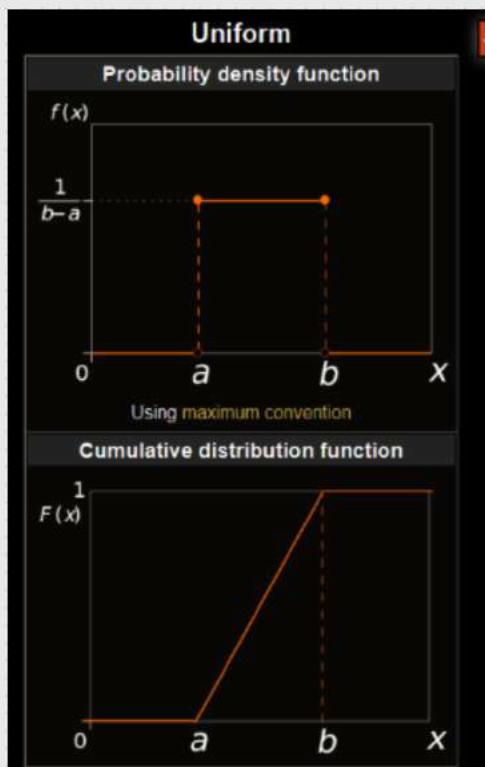
## Uniform Distribution

① Continuous Uniform Distribution (pdf)

② Discrete Uniform Distribution (pmf)

① Continuous Uniform Distribution [Continuous Random Variable]

In probability theory and statistics, the continuous uniform distribution or rectangular distribution is a family of symmetric probability distributions. The distribution describes an experiment where there is an arbitrary outcome that lies between certain bounds. The bounds are defined by the parameters,  $a$  and  $b$ , which are the minimum and maximum values.



Notation :  $U(a,b)$

Parameters :  $-\infty < a < b < \infty$

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

$$\text{Cdf} = \begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } x \in [a,b] \\ 1 & \text{for } x > b \end{cases}$$

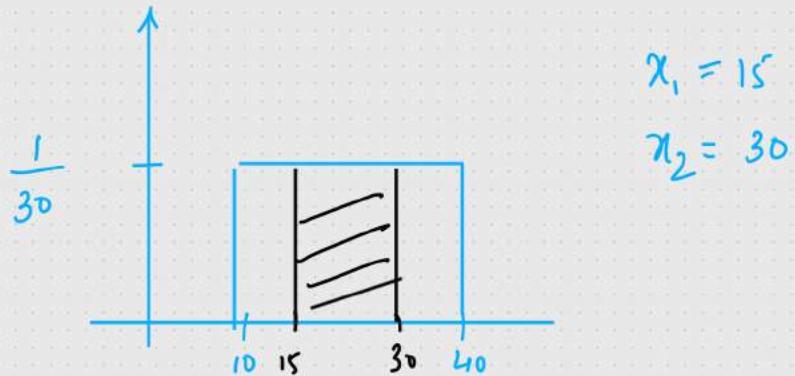
$$\text{Mean} = \frac{1}{2}(a+b) \quad \text{Variance} = \frac{1}{12}(b-a)^2$$

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

Eg: The number of candies sold daily at a shop is uniformly distributed with a maximum of 40 and a minimum of 10

i) Probability of daily sales to fall between 15 and 30?

Ans)

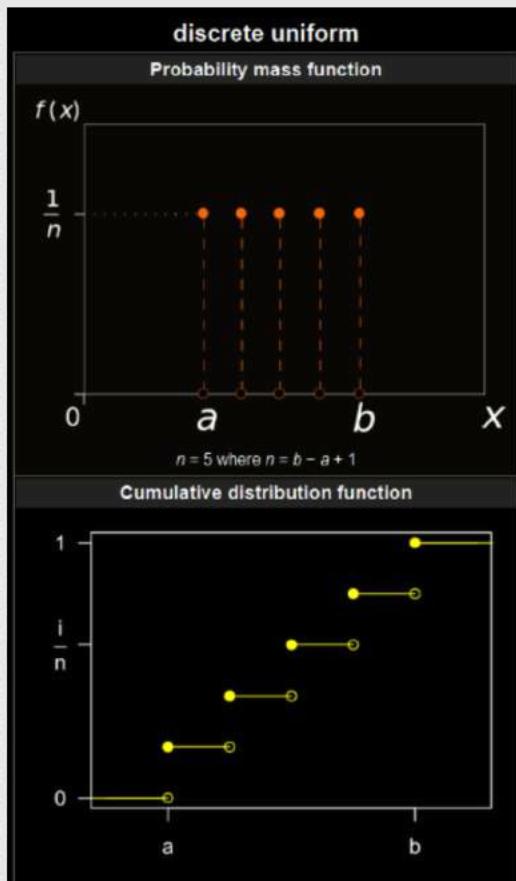


$$\begin{aligned} \Pr(15 \leq x \leq 30) &= (x_2 - x_1) * \frac{1}{b-a} \\ &= 15 * \frac{1}{30-2} \\ &= 0.5 \end{aligned}$$

$$\begin{aligned} \Pr(x \geq 20) &= (40-20) * \frac{1}{30} \\ &= 0.66 = 66\% \end{aligned}$$

## ② Discrete Uniform Distribution (pmf)

In probability theory and statistics, the discrete uniform distribution is a symmetric probability distribution wherein a finite number of values are equally likely to be observed; every one of  $n$  values has equal probability  $1/n$ . Another way of saying "discrete uniform distribution" would be "a known, finite number of outcomes equally likely to happen".



Eg: Rolling a dice

$$\{1, 2, 3, 4, 5, 6\}$$

$$a=1 \quad b=6$$

$$\Pr(1) = \frac{1}{6}$$

$$\Pr(2) = \frac{1}{6}$$

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

$$\frac{1}{n} \Rightarrow [n = b - a + 1]$$

Notation  $\mathcal{U}(a, b)$

Parameters  $a, b$  with  $b \geq a$

PMF

$$\frac{1}{n}$$

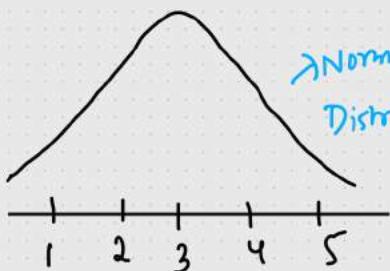
Mean  $\frac{a+b}{2}$

Median

## Standard Normal Distribution And Z-Score

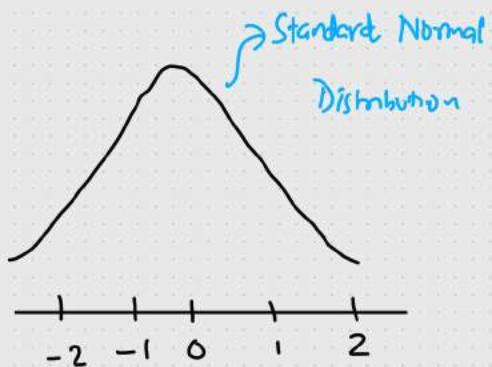
Z-stats

$$X = \{1, 2, 3, 4, 5\} \quad \text{Normally Distributed}$$



$$\sigma = 1.414 \approx 1$$

$$\Rightarrow \begin{aligned} \mu &= 0 \\ \sigma &= 1 \end{aligned}$$



$$Z\text{-score} = \frac{x_i - \mu}{\sigma}$$

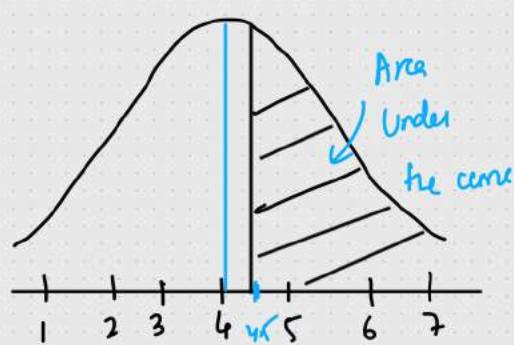
$$X \sim \text{SND}(\mu=0, \sigma=1)$$

$$= \frac{1-3}{1} = -2 \quad \frac{4-3}{1} = 1$$

$$= \frac{2-3}{1} = -1 \quad \frac{5-3}{1} = 2$$

$$= \frac{3-3}{1} = 0$$

①



$$\begin{aligned} \mu &= 4 \\ \sigma &= 1 \end{aligned}$$

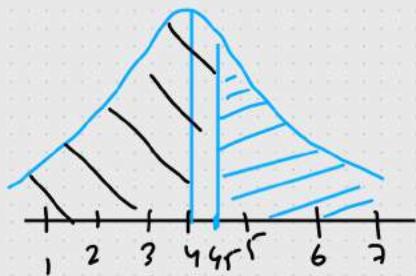
① How many Standard deviation 4.5 is away from mean?

$$x_i = 4.5$$

$$Z\text{-score} = \frac{4.5 - 4}{1} = 0.5$$

Question :  $\mu = 4$   $\sigma = 1$

What percentage of data is falling above 4.5?



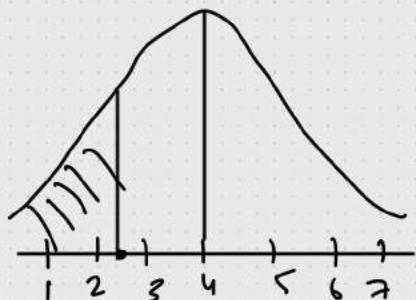
$$\begin{aligned} Z\text{-score} &= \frac{4.5 - 4}{1} \\ &= 0.5 \\ &= 0.69146 \\ &= 0.30854 \end{aligned}$$

$$\text{Area under the curve } ( \geq 4.5 ) = 1 - 0.69146$$

$$= 0.30854$$

$$= 30.85\%$$

④ What percentage of data is falling below 2.5?



$$Z\text{-score} = \frac{2.5 - 4}{1} = -1.5$$

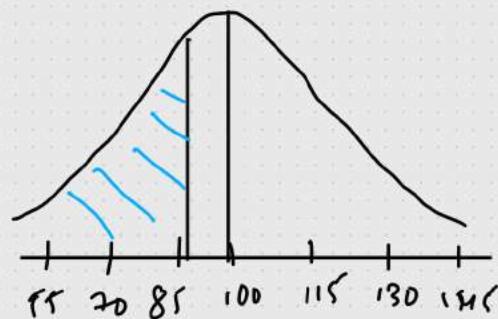
Area under the curve ( $\leq 2.5$ ) = 0.06681

= 6.6%

=====

Prob) In India the average IQ is 100, with a standard deviation of 15. What is the percentage of the population would you expect to have an IQ lower than 85?

Ans)  $\mu = 100 \quad \sigma = 15 \quad x_i = 85$



$$\begin{aligned} \textcircled{1} \quad Z \text{ score} &= \frac{x_i - \mu}{\sigma} \\ &= \frac{85 - 100}{15} \\ &= -1.11 \end{aligned}$$

Area under the curve = 0.15866

= 15.866

Area under the curve ( $> 85$ ) =  $1 - 0.15866$

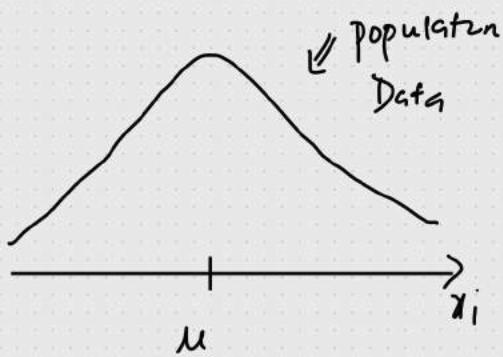
=  $\approx 84\%$

## Central Limit Theorem (CLT)

The central limit theorem says that the sampling distribution of the mean will always be normally distributed, as long as the sample size is large enough. Regardless of whether the population has a normal, Poisson, binomial, or any other distribution, the sampling distribution of the mean will be normal.

$$n = 20$$

$$\textcircled{1} \quad X \sim N(\mu, \sigma)$$



$$S_1 = \{x_1, x_2, x_3, \dots, x_{20}\} = \bar{x}_1$$

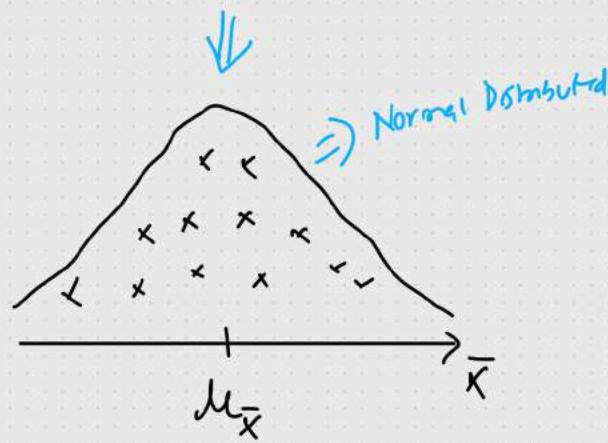
$$S_2 = \{x_2, x_3, x_5, \dots, x_{20}\} = \bar{x}_2$$

$$S_3 = \{ \dots \} = \bar{x}_3$$

$$\vdots \quad \vdots \quad \vdots$$

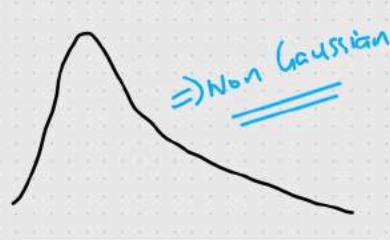
$$S_m = \{ \dots \} = \bar{x}_m$$

$$\bar{X} = \{\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_m\}$$



$$\rightarrow n \geq 30$$

$$\textcircled{2} \quad X \not\sim N(\mu, \sigma)$$



$\Rightarrow$

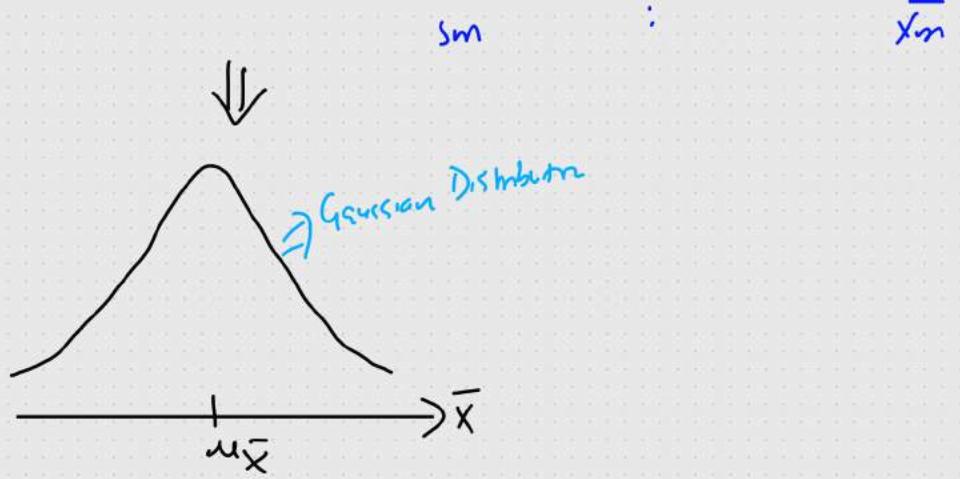
$$S_1 = \{x_1, x_2, \dots, x_{30}\} = \bar{x}_1$$

$$S_2 = \{x_2, x_3, \dots, x_{30}\} = \bar{x}_2$$

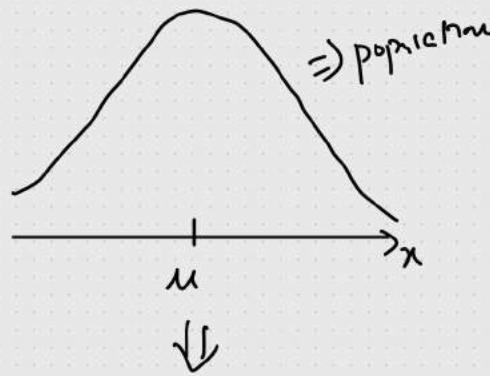
$$S_3 = \{ \dots \} = \bar{x}_3$$

$$S_4 = \{ \dots \} = \bar{x}_4$$

$$\vdots \quad \vdots \quad \vdots$$



## ① Important for Interview



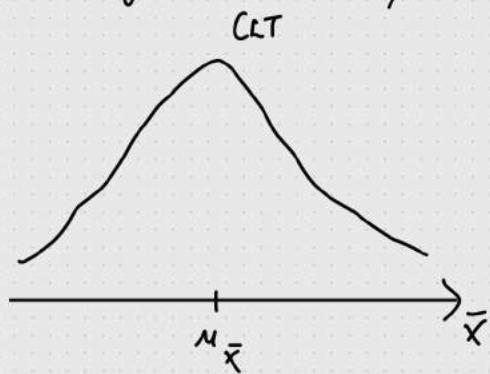
$$X \sim N(\mu, \sigma)$$

$\sigma$  = population Std

$\mu$  = population mean

$n$  = sample size

Sampling distribution of mean



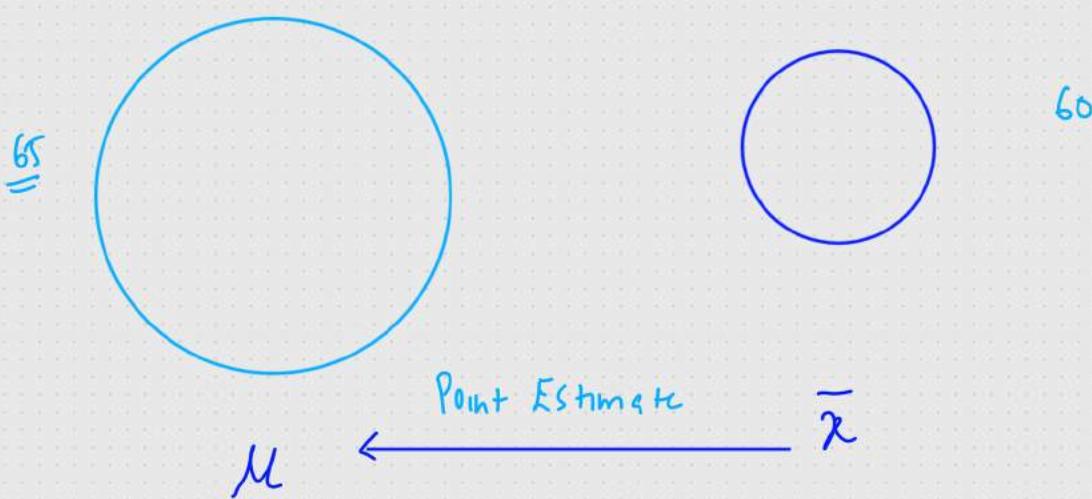
$$\bar{X} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

$n$  can be any  
Value

Estimate: It is an observed numerical value used to estimate an unknown population parameter

(1) Point Estimate : Single numerical value used to estimate the unknown population parameter

Sample mean is a point estimate of a population mean

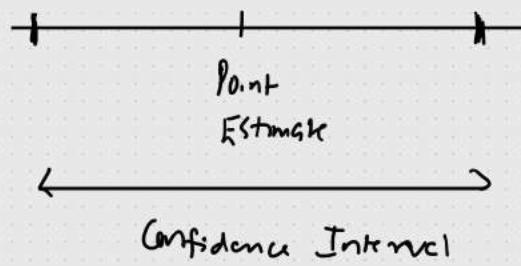


(2) Interval Estimate : Range of values used to estimate the unknown population parameters

Interval estimates of population parameters

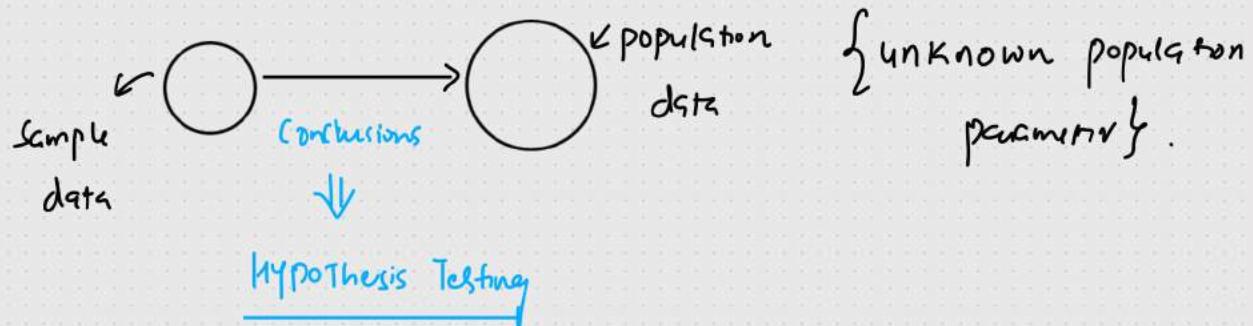
are called Confidence Intervals

$$[55 - 65]$$



## A Hypothesis And Hypothesis Testing Mechanism

Inferential Stats : Conclusion or Inferences



### Hypothesis Testing Mechanism

Person Crime → Court

① Null Hypothesis ( $H_0$ ) = The Person is not guilty

- The assumption you are beginning with

② Alternate Hypothesis ( $H_1$ ) = The Person is guilty.

- Opposite of Null Hypothesis

{ p value }  $\underline{=}$

③ Experiments → Proof collect (DNA Test, Finger Print)  $\Rightarrow$  Statistical Analysis

④ Accept the Null Hypothesis or Reject the Null Hypothesis

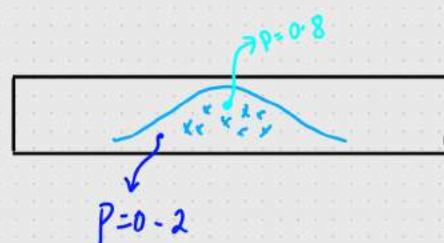
Eg: Colleges at District A states its <sup>Average</sup> passed percentage of Students are 85%. A new college opened in the district And it was found that a sample of student 100 have a pass percentage of 90% with a standard deviation of 4%. Does this school have a different passed percentage.

Ans) Null Hypothesis  $H_0 = \mu = 85\%$ .

Alternate Hypothesis  $H_1 = \mu \neq 85\%$ .

## P value

The p value is a number, calculated from a statistical test, that describes how likely you are to have found a particular set of observations if the null hypothesis were true. P values are used in hypothesis testing to help decide whether to reject the null hypothesis.



Out of 100 touches in this key the probability of touching in this region 20

## Hypothesis Testing

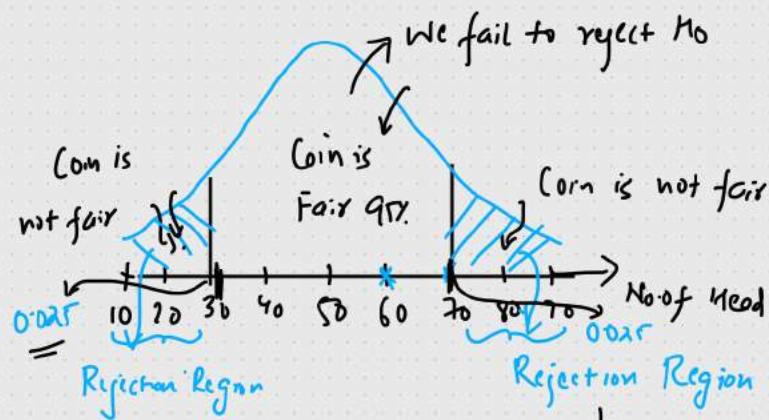
Exp = Coin is fair or Not {100 Tosses}

$$\{H, T\}$$

① Null Hypothesis  $H_0$  : Coin is fair  $P(H) = 0.5 \quad P(T) = 0.5$

② Alternate Hypothesis  $H_1$  : Coin is not fair  $P(H) = 0.6 \quad P(T) = 0.4$   
 $\rightarrow P(H) = 0.7 \quad P(T) = 0.3$

## Experiment



We Reject the Null

Hypothesis

We reject the  
null Hypothesis

Significance value ;  $\alpha = 0.05 \quad CI = 1 - 0.05$   
 $= 0.95$

## Conclusion

$P = 0.01 < \text{Significance}$

We reject the Null Hypothesis

Clue

We fail to reject Null hypothesis

## Hypothesis Testing And Statistical Analysis

- ① Z Test }  $\Rightarrow$  Average  $\Rightarrow$  Z table  $\rightarrow$  Z score And p value
- ② T Test  $\Rightarrow$  Average  $\Rightarrow$  t table
- ③ CHI SQUARE  $\Rightarrow$  Categorical Data
- ④ ANNOVA  $\Rightarrow$  Variance

Z test : i) Population std  $\Downarrow$  ii)  $n \geq 30$

With a  $\sigma = 3.9$

i) The average heights of all residents in a city is 168cm. A doctor believes the mean to be different. He measured the height of 36 individuals and found the average height to be 169.5 cm.

(a) State null and Alternate Hypothesis

(b) At a 95% confidence level, is there enough evidence to reject the null hypothesis.

Ans)  $M = 168\text{cm}$ ,  $\sigma = 3.9$ ,  $n = 36$   $\bar{x} = 169.5\text{cm}$

a) Null hypothesis  $H_0: M = 168\text{cm}$

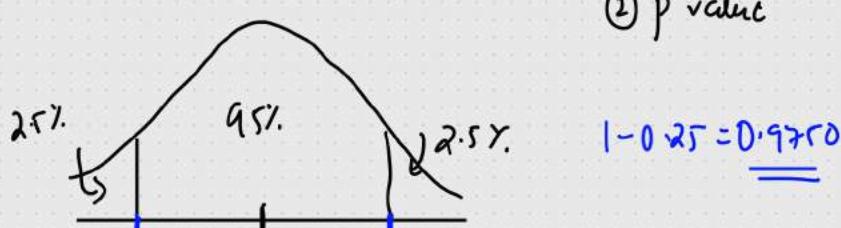
b) Alternate hypothesis  $H_1: M \neq 168\text{cm}$  { 2 Tail Test }

c)  $C.I = 0.95$   $\alpha = 1 - 0.95 = 0.05$

Decision Boundary

① Z test

② P value



-1.96      168      +1.96

$$Z\text{-score} = \frac{X_i - \mu}{\sigma}$$

## a) Statistical Analysis

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

$$\boxed{\sigma / \sqrt{n}}$$

$$= \frac{169.5 - 168}{3.9 / \sqrt{36}} = 2.31$$

If Z-test value is less than -1.96 or greater +1.96 We Reject the

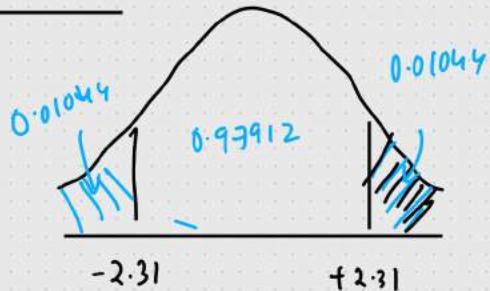
Null Hypothesis

Else

We Accept Null Hypothesis

$2.31 > +1.96$  {We Reject the Null Hypothesis}.

## ② P-value



$$1 - \text{Area under the curve of } +2.31 \\ 1 - 0.98956$$

=

$$p\text{-value} = 0.01044 + 0.01044 = 0.02088\%.$$

if P-value < Significance

0.02088 < 0.05 {Reject the Null Hypothesis}

② A factory manufactures bulbs with an average warranty of 5 years with standard deviation of 0.50. A worker believes that the bulb will malfunction in less than 5 years. He tests a sample of 40 bulbs and finds the average time to be 4.8 years.

- (a) State null and alternate hypothesis
- (b) At a 2% significance level, is there enough evidence to support the idea that the warranty should be revised?

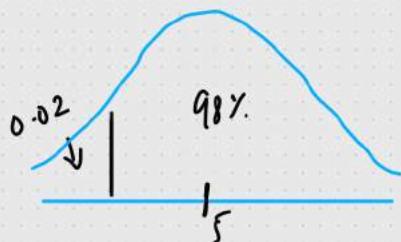
$$\text{Ans) } \mu = 5 \quad \sigma = 0.50 \quad n = 40 \quad \bar{x} = 4.8$$

① Null Hypothesis  $H_0 : \mu = 5$

② Alternate Hypothesis  $H_1 : \mu < 5$  {1 Tail Test}

③ Decision Boundary  $C.I = 0.98$

$$\alpha = 1 - 0.98 = 0.02$$

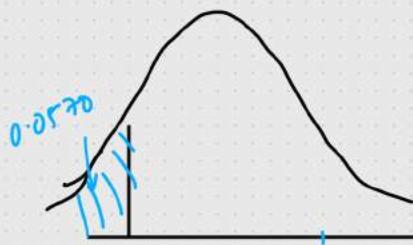


④ P value

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

$$= \frac{4.8 - 5}{0.50 / \sqrt{40}}$$

$$Z_{\text{test}} = -2.53$$



Area under the curve of  $-2.53$

$$Z_{\text{value}} \text{ is } \underline{\underline{0.0570}}$$

$P\text{-value} = 0.0570$ .

if P-value < Significance

$$0.0570 < 0.02 \Rightarrow \text{False}$$

{We accept the Null Hypothesis}

Conclusion : The Warranty needs to be revised.

## Student t distribution

Statistical Analysis using Z-score we need population standard deviation ( $\sigma$ )

How do we perform an analysis when we don't know the population standard deviation?

Soln: Student t distribution

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

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

s = sample std

## Degree of freedom

$$dof = \underline{n - 1}$$

n = sample size

3 people



<https://www.sjsu.edu/faculty/gerstman/StatPrimer/t-table.pdf>

T-stats  $\div$  J test  $\rightarrow$  One Sample t-test.

① In the population the average IQ is 100. A team of researchers want to test a new medication to see if it has either a positive or negative effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a standard deviation of 20. Did the medication affect intelligence? ( $\cdot I = 95\%$ )

$$\text{Ans) } \mu = 100 \quad n = 30 \quad \bar{x} = 140 \quad \sigma = 20 \quad (\cdot I = 0.95)$$

$$\alpha = 0.05$$

① Null Hypothesis  $H_0: \mu = 100$

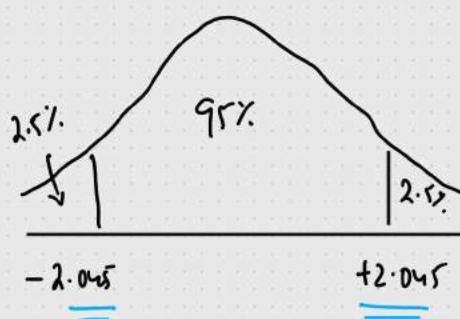
Alternate "  $H_1: \mu \neq 100$  {2 Tail Test}.

$$② \alpha = 0.05$$

③ Degree of freedom

$$df = n - 1 = 30 - 1 = 29$$

④ Decision Rule



If t test is less than  $-2.045$  and greater than  $2.045$ , Reject the Null Hypothesis

## ⑤ Calculate t test statistics

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{140 - 100}{20/\sqrt{30}} = \frac{40}{3.65} = 10.96$$

## ⑥ Conclusion

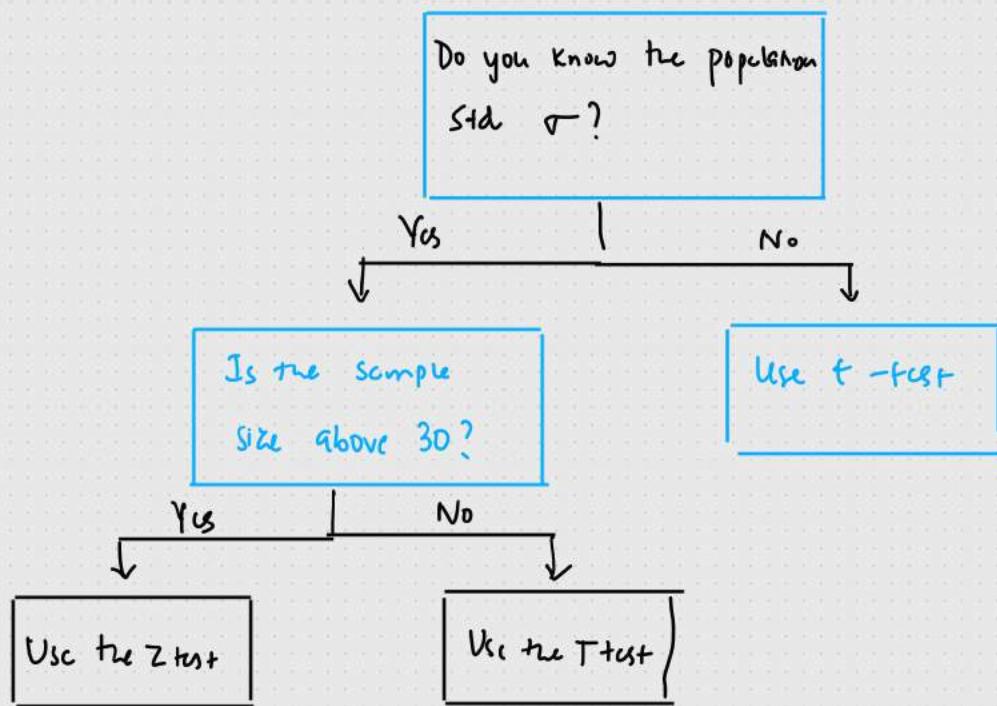
Decision Rule: If  $t$  is less than  $-2.0452$  and greater than  $2.0452$ , reject the Null Hypothesis

$t = 10.96 > 2.0452 \Rightarrow$  Rejecting the Null Hypothesis

Answer

Conclusion: Medication has increased the Intelligence.

## When To Use T-test Vs Z-test



## Type 1 and Type 2 Errors

Reality : Null Hypothesis is True or Null Hypothesis is False

Decision : Null Hypothesis is True or Null Hypothesis is False



Conclusion

Outcome 1 : We reject the Null Hypothesis When in reality it is False → Good

Outcome 2 : We reject the Null Hypothesis When in reality it is True → Type 1 Error

Outcome 3 : We retain the Null Hypothesis, when in reality it is False → Type 2 Error

Outcome 4 : We retain the Null Hypothesis when in reality it is True → Good

## Bayes Statistics (Bayes Theorem)

Bayesian statistics is an approach to data analysis and parameter estimation based on Bayes' theorem.

### Bayes' Theorem

Probability → Independent Events  
→ Dependent Events

#### ① Independent Events

Eg: Rolling a Dice

$$\{1, 2, 3, 4, 5, 6\}$$

$$Pr(1) = \frac{1}{6} \quad Pr(2) = \frac{1}{6} \quad \dots$$

#### Tossing a Coin

$$Pr(H) = 0.5 \quad Pr(T) = 0.5$$

#### ② Dependent Events

Red  $Pr(R) = \frac{2}{5}$  Yellow  $Pr(Y) = \frac{3}{4}$

|   |   |   |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 0 |   |

$$Pr(R \text{ and } Y) = Pr(R) * Pr(Y|R)$$

$$= \frac{2}{5} * \frac{3}{4} = \frac{6}{20}$$

$$Pr(A \text{ and } B) = Pr(B \text{ and } A)$$

$$Pr(A) * Pr(B|A) = Pr(B) * Pr(A|B)$$

$$Pr(B|A) = \frac{Pr(B) * Pr(A|B)}{Pr(A)}$$

Bayes theorem

conditional probability



$$\Pr(A|B) = \frac{\Pr(A) * \Pr(B|A)}{\Pr(B)}$$

$A, B$  = counts

$\Pr(A|B)$  = probability of  $A$  given  $B$  is true

$\Pr(B|A)$  = probability of  $B$  given  $A$  is true

$\Pr(A), \Pr(B)$  = The independent probability of  $A$  and  $B$ .

### Dataset

| Size of movie | No. of Room | Location | Price |
|---------------|-------------|----------|-------|
| $x_1$         | $x_2$       | $x_3$    | $y$   |

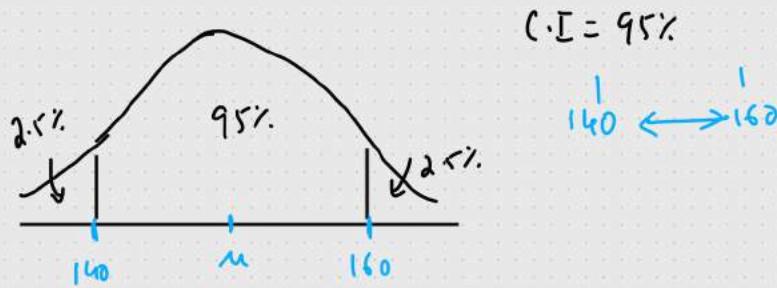
$$\Pr(y/x_1, x_2, x_3) = \frac{\Pr(y) * \Pr(x_1, x_2, x_3/y)}{\Pr(x_1, x_2, x_3)}$$



Bayes theorem

## Confidence Intervals and Margin of Error

$$\mu = 160$$



Point Estimate : A value of any statistics that estimates the value of an unknown population parameter is called Point Estimate

$$\bar{x} \longrightarrow \mu$$

$$\bar{x} = 2.95 \quad \mu = 3$$

## Confidence Interval

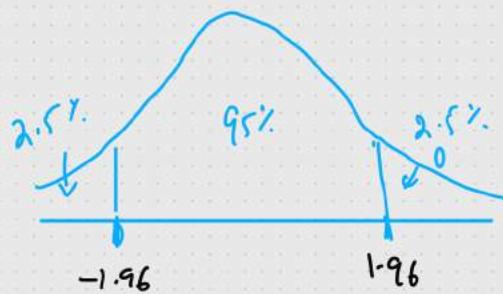
We construct a confidence interval to help estimate what the actual value of the unknown population mean is.

Point Estimate  $\pm$  Margin of Error

$$Z_{\text{MS}} + \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

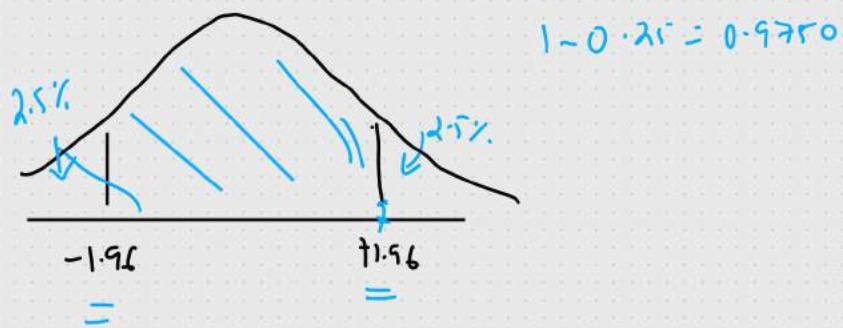
$\alpha = 0.05$

$$Z_{0.05/2} \Rightarrow Z_{0.025}$$



- ① On the Verbal section of CAT Exam, the standard deviation is known to be 100. A sample of 25 test takers has a mean of 520. Construct a 95% CI about the mean?

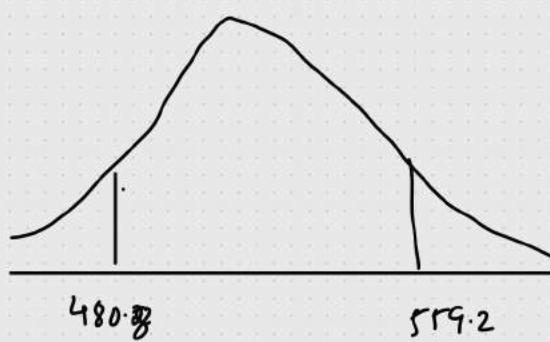
Ans)  $\bar{x} = 520$        $\sigma = 100$        $n = 25$        $CI = 0.95 \quad \alpha = 0.05$



$$\bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$\text{Lower CI} = 520 - (1.96) \times \frac{100}{\sqrt{25}} = 480.8$$

$$\text{Higher CI} = 520 + (1.96) \times \frac{100}{\sqrt{25}} = 559.2$$



I am 95% confident that the mean CAT score lies  
between 480.8 and 559.2

## CHI SQUARE TEST

The Chi Square Test for Goodness of fit test claims about population proportions. {Categorical variables}

It is a non parametric test that is performed on Categorical [ordinal, nominal] data

Eg: There is a population of Male who likes different color of bikes

|             | <u>Theory</u> | <u>Sample</u> | ① Goodness of fit. |
|-------------|---------------|---------------|--------------------|
| Yellow Bike | $\frac{1}{3}$ | 22            |                    |
| Orange Bike | $\frac{1}{3}$ | 17            |                    |
| Red Bike    | $\frac{1}{3}$ | 59            |                    |

↓      ↳ observe categorical distribution.

Theory      categorical distribution.

## Goodness of fit test

- \* In a student class of 100 students, 30 are Right handed. Does this class fit the theory 12% of people are right handed

|              | <u>O</u> | <u>E</u> | $\frac{12}{100} \times 100 = 12$       |
|--------------|----------|----------|--|
| Right handed | 30       | 12       |  |
| Left handed  | 70       | 88       | ⇒ Theory      categorical distribution |

↓  
Sample ↑

## CHI SQUARE For Goodness of Fit

In 2010 Census of the city, the weight of the individuals in a small city were found to be the following

| <50kg | 50 - 75 | >75 |
|-------|---------|-----|
| 20%   | 30%     | 50% |

In 2020, weight of  $n=500$  individuals were sampled. Below are the results

| <50 | 50-75 | >75 |
|-----|-------|-----|
| 140 | 160   | 200 |

Using  $\alpha=0.05$ , would you conclude the population difference of weights has changed in the last 10 years?

Ans)

2010  
Expected

| <50kg | 50 - 75 | >75 |
|-------|---------|-----|
| 20%   | 30%     | 50% |

2020  
 $n=500$   
Observed

| <50 | 50-75 | >75 |
|-----|-------|-----|
| 140 | 160   | 200 |

| Expected | <50                       | 50-75                     | >75                       |
|----------|---------------------------|---------------------------|---------------------------|
|          | $0.2 \times 500$<br>= 100 | $0.3 \times 500$<br>= 150 | $0.5 \times 500$<br>= 250 |

① Null Hypothesis :  $H_0$  : The data meets the expectation

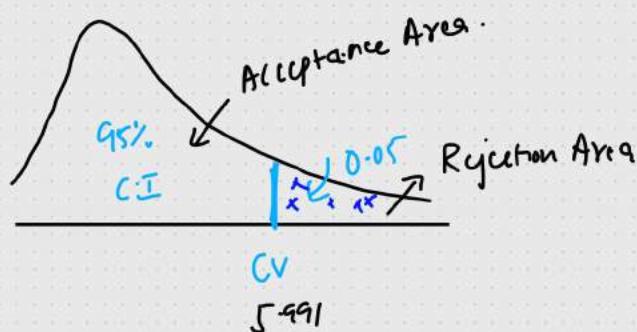
Alternate Hyp :  $H_1$  : The data does not meet the expectation

②  $\alpha = 0.05$  ( $I = 95\%$ )

③ Degree of freedom

$$df = K - 1 = 3 - 1 = 2$$

④ Decision Boundary



If  $\chi^2$  is greater than 5.99, Reject  $H_0$   
else

We fail to reject the Null Hypothesis

⑤ Calculate Chi-Square Test Statistic

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

$$= \frac{(140 - 100)^2}{100} + \frac{(160 - 150)^2}{150} + \frac{(200 - 250)^2}{250}$$

2020  
 $n=500$   
Observed

| <50 | 50-75 | >75 |
|-----|-------|-----|
| 140 | 160   | 200 |

Expected

| <50                    | 50-75                  | >75                    |
|------------------------|------------------------|------------------------|
| $0.2 \times 500 = 100$ | $0.3 \times 500 = 150$ | $0.5 \times 500 = 250$ |

$$= \frac{1600}{100} + \frac{100}{150} + \frac{2500}{250}$$

$$= 16 + 0.66 + 10$$

$$= 26.66$$

$$\chi^2 = 26.66$$

If  $\chi^2$  is greater than 5.99, Reject  $H_0$   
else

We fail to reject the Null Hypothesis

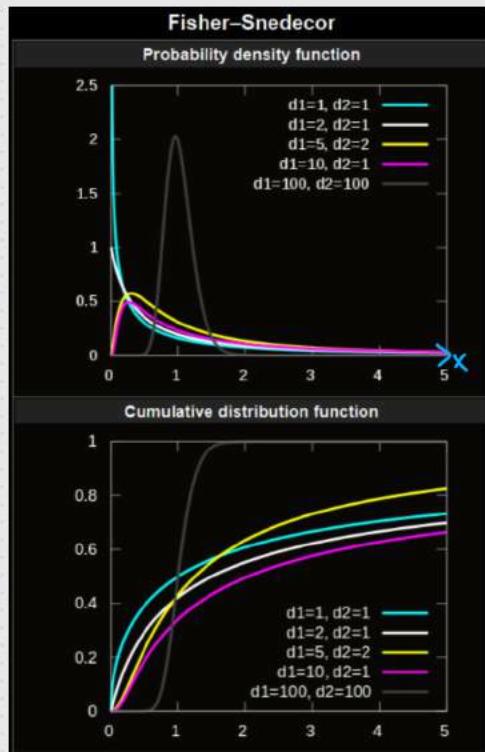
$$26.66 > 5.99, \text{ Reject } H_0$$

### Answer

The weights of 2020 population are different  
than those expected in the 2010 population

# F distribution

In probability theory and statistics, the F-distribution or F-ratio, also known as Snedecor's F distribution or the Fisher-Snedecor distribution (after Ronald Fisher and George W. Snedecor) is a continuous probability distribution that arises frequently as the null distribution of a test statistic, most notably in the analysis of variance (ANOVA) and other F-tests.



Parameters :  $d_1, d_2 > 0$  degree of freedom

Support  $x \in (0, +\infty)$

F table

$\alpha = 0.05$   $d_1, d_2$

$$pdf = f(x; d_1, d_2) = \frac{\sqrt{\frac{(d_1 x)^{d_1} d_2^{d_2}}{(d_1 x + d_2)^{d_1 + d_2}}}}{x^{\underline{B\left(\frac{d_1}{2}, \frac{d_2}{2}\right)}}}$$

Beta function

$$B(m, n) = \frac{(m-1)! (n-1)!}{(m+n-1)!} = \frac{m+n}{mn} \binom{m+n}{m}$$

The F distribution with  $d_1$  and  $d_2$  degree of freedom is the distribution of

$$X = \frac{S_1/d_1}{S_2/d_2}$$

$S_1$  = Independent Random Variables  $\left\{ \begin{array}{l} \text{Chi square} \\ \text{distribution} \end{array} \right\}$

$S_2$  = " " "

(S1)  $d_1$  = Degree of freedom

(S2)  $d_2$  = " " "

F Test : Variance Ratio Test

## F-Test [Variance Ratio Test]

① The following data shows the no. of bulbs produced daily for some days by 2 workers A and B

| A  | B  |
|----|----|
| 40 | 39 |
| 30 | 38 |
| 38 | 41 |
| 41 | 33 |
| 38 | 32 |
| 35 | 39 |
| 40 |    |
| 34 |    |

Can we consider based on the data

Worker B is more stable and efficient

$$\alpha = 0.05$$

Pw) ① Null Hypothesis  $H_0: \sigma_1^2 = \sigma_2^2$

Alternate Hypothesis  $H_1: \sigma_1^2 \neq \sigma_2^2$

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

② Calculation of Variance

| A     |           |                     | B     |             |                     |
|-------|-----------|---------------------|-------|-------------|---------------------|
| $x_i$ | $\bar{x}$ | $(x_i - \bar{x})^2$ | $x_2$ | $\bar{x}_2$ | $(x_2 - \bar{x})^2$ |
| 40    | 37        | 9                   | 39    | 37          | 4                   |
| 30    | 37        | 49                  | 38    | 37          | 1                   |
| 38    | 37        | 1                   | 41    | 37          | 36                  |
| 41    | 37        | 16                  | 33    | 37          | 16                  |
| 38    | 37        | 1                   | 32    | 37          | 25                  |
|       |           |                     | 39    | 37          | 4                   |
|       |           |                     | 40    | 37          | 9                   |

$$\frac{35}{\bar{x}_1} = 37$$

$$\sum (x_i - \bar{x})^2 = 80$$

$$\frac{34}{\bar{x}_2} = 37$$

$$\sum (x_i - \bar{x})^2 = 84$$

$$S_1^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$S_2^2 = \frac{84}{8-1}$$

$$= \frac{84}{7} = 12$$

$$S_1^2 = \frac{80}{6-1} = \frac{80}{5} = 16$$

→ Calculation of Variance Ratio F-test

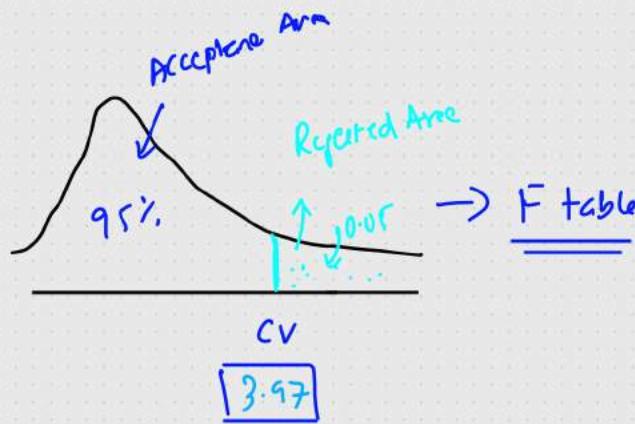
$$F = \frac{S_1^2}{S_2^2} = \frac{16}{12} = 1.33$$

### ③ Decision Rule

$$df_1 = 6-1 = 5$$

$$df_2 = 8-1 = 7$$

$$\alpha = 0.05$$



If F test is greater than 3.97, Reject the Null

Hypothesis

$1.33 < 3.97$ , we fail to Reject the Null

Conclusion

Hypothesis

Worker B is not efficient when worked to worker A.

## Analysis Of Variance (ANOVA)

Dfn : ANOVA is a statistical method used to compare the means of 2 or more group

### Anova

① Factors (variables)

② Levels

Eg : Factor = Medicines

levels = 5mg      10mg      20mg [Dosage]

Mode of Payment → Factor

GPay      PhonePE      IMPS      NEFT [Levels]

## Analysis Of Variance (ANOVA)

### Assumptions in ANOVA

#### ① Normality of Sampling Distribution of Means

The distribution of sample mean is normally distributed.

#### ② Absence of Outliers

Outlying score need to be removed from dataset.

#### ③ Homogeneity of Variance

Each one of the population has same variance

$$[\sigma_1^2 = \sigma_2^2 = \sigma_3^2]$$

Population variance in different levels of each independent variable are equal

#### ④ Samples are independent and random.

## Analysis Of Variance (ANOVA)

### Types of ANOVA

- ① One Way ANOVA : One factor with at least 2 levels, these levels are independent

Eg: Doctor want to test a new medication to decrease headache.

They split the participant in 3 conditions [10mg, 20mg, 30mg].  
Doctor ask the participants to rate the headache [1-10]

| Medication → Factor |      |      |
|---------------------|------|------|
| 10mg                | 20mg | 30mg |
| 5                   | 7    | 2    |
| 9                   | 8    | 7    |
| -                   | -    | -    |
| -                   | -    | -    |

- ② Repeated Measures Anova : One factor with at least 2 levels, levels are dependents

Running → Factor

| Day 1 | Day 2 | Day 3 |
|-------|-------|-------|
| 8     | 5     | 6     |
| 7     | 4     | 3     |

③ Factorial Anova : Two or more factors (each of which with at least 2 levels), levels can be either independent and dependent

|        |       | Running → Factor |       |       |
|--------|-------|------------------|-------|-------|
|        |       | Day 1            | Day 2 | Day 3 |
| Gender | Men   | 8                | 5     | 6     |
|        | Woman | 7                | 4     | 3     |
|        |       | 6                | 5     | 4     |
|        |       | 3                | 2     | 1     |

## Analysis Of Variance (ANOVA)

### Hypothesis Testing In ANOVA

Null hypothesis  $H_0$  :  $\mu_1 = \mu_2 = \mu_3 = \dots = \mu_K$

Alternate hypothesis  $H_1$  : At least one of the mean is not equal  
 $\times \quad \boxed{\mu_1 \neq \mu_2 \neq \mu_3 \dots \neq \mu_K}$

### Test Statistics

$$F = \frac{\text{Variation between Sample}}{\text{Variation within Sample}}$$

$$\underline{\underline{F \rightarrow dd}}$$

|                               |   | Variance between Sample |       |       |   |
|-------------------------------|---|-------------------------|-------|-------|---|
|                               |   | $x_1$                   | $x_2$ | $x_3$ |   |
| Variation<br>within<br>Sample | 1 | 1                       | 6     | 5     | $H_0: \bar{x}_1 = \bar{x}_2 = \bar{x}_3$            |
|                               | 2 | 2                       | 7     | 6     | $H_1: \text{At least one sample mean is not equal}$ |
|                               | 4 | 4                       | 3     | 3     |   |
|                               | 5 | 5                       | 2     | 2     |   |
|                               | 3 | 3                       | 1     | 4     |   |
|                               |   | <hr/>                   | <hr/> | <hr/> |   |

$$\sum x_1 = 15 \quad \sum x_2 = 19 \quad \sum x_3 = 20$$

$$\bar{x}_1 = 3 \quad \bar{x}_2 = 19/5 \quad \bar{x}_3 = 4$$

## One Way ANOVA

One factor with at least 2 levels, levels are independent

- ① Doctors want to test a new medication which reduces headache. They split the participant into 3 condition [15mg, 30mg, 45mg]. Later on the doctor ask the patient to rate the headache between [1-10]. Are there any differences between the 3 conditions using  $\alpha = 0.05$ ?

Ans)

| 15 mg | 30mg | 45mg |
|-------|------|------|
| 9     | 7    | 4    |
| 8     | 6    | 3    |
| 7     | 6    | 2    |
| 8     | 7    | 3    |
| 8     | 8    | 4    |
| 9     | 7    | 3    |
| 8     | 6    | 2    |

- ① Define Null and Alternative hypotheses?

$$H_0: \mu_{15} = \mu_{30} = \mu_{45}$$

$H_1$ : not all  $\mu$ 's are equal

- ② State Significance value

$$\alpha = 0.05 \Rightarrow C.I = 0.95$$

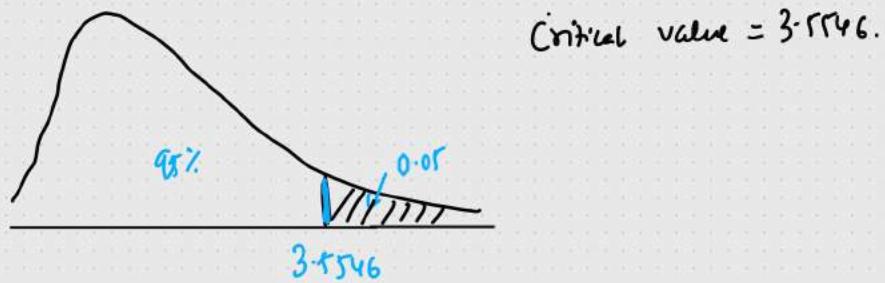
### ③ Calculate Degree of freedom

$$N = 21 \quad a = 3 \quad n = 7$$

$$df_{\text{between}} = a - 1 = 3 - 1 = 2 \quad (2, 18)$$

$$df_{\text{within}} = N - a = 21 - 3 = 18 \quad \Downarrow$$

$$df_{\text{total}} = N - 1 = 21 - 1 = 20 \quad F_{\text{table}}$$



### ④ State Decision Rule

If  $F$  is greater than 3.5546, reject the null hypothesis.

### ⑤ Calculate Test Statistics

| SS | df | MS | F |
|----|----|----|---|
|----|----|----|---|

Between

Within

Total

$SS_{\text{between}}$

$SS_{\text{within}}$

$SS_{\text{total}}$

$$\textcircled{1} \quad SS_{\text{between}} = \frac{\sum (\sum a_i)^2 - T^2}{n}$$

$$15\text{mg} = 9+8+7+8+8+9+8 = 57$$

$$30\text{mg} = 7+6+6+7+8+7+6 = 47$$

$$45\text{mg} = 4+3+2+3+4+3+2 = 21$$

$$= \frac{57^2 + 47^2 + 21^2}{7} - \frac{[57+47+21]}{21}$$

$$= \boxed{98.67}$$

$$\textcircled{2} \quad SS_{\text{within}} = \sum y^2 - \frac{\sum (\sum a_i)^2}{n}$$

$$= \sum y^2 - \left[ \frac{57^2 + 47^2 + 21^2}{7} \right]$$

$$\sum y^2 = 9^2 + 8^2 + 7^2 + 8^2 + 8^2 + 9^2 + 8^2 + 7^2 + 6^2$$

$$+ \quad - \quad - \quad - \quad - \quad -$$

$$= 853.$$

$$= 853 - \left[ \frac{57^2 + 47^2 + 21^2}{7} \right]$$

$$= \boxed{10.29.}$$

$$\textcircled{3} \quad SS_{\text{Total}} = \sum y^2 - \frac{T^2}{N}$$

$$= 853 - \frac{121^2}{21} = \boxed{108.95}$$

| 15mg | 30mg | 45mg |
|------|------|------|
| 9    | 7    | 4    |
| 8    | 6    | 3    |
| 7    | 6    | 2    |
| 8    | 7    | 3    |
| 8    | 8    | 4    |
| 9    | 7    | 3    |
| 8    | 6    | 2    |

|         | SS     | df | MS    | F |
|---------|--------|----|-------|---|
| Between | 98.67  | 2  | 49.34 |   |
| Within  | 102.9  | 18 | 0.54  |   |
| Total   | 108.95 | 20 |       |   |

$$F = \frac{\text{Variation between samples}}{\text{Variation within samples}}$$

$$F = \frac{MS_{\text{between}}}{MS_{\text{within}}}$$

$$F = \frac{49.34}{0.54} = \underline{\underline{86.56}}$$

If F is greater than 3.5546, reject the Null hypothesis.

$86.56 > 3.5546$ , Reject the Null hypothesis