

# \* Bivariate Random Variable

Let  $S$  be a Sample Space associated with random experiment. Let  $X(t), Y(t)$  be 2 random variable function each assigning a real number to each outcome  $t \in S$ . Then  $(X, Y)$  is called a Bivariate or 2-dimensional random variable.

→ PMF of  $(x, y)$  :-

If  $(x, y)$  is a 2-D discrete random variable, where  $x = x_i$  &  $y = y_j$ ,  $i, j = 1, 2, \dots$

Then  $P(X=x_i, Y=y_j) = P_{ij}$  is called Probability Mass Function of  $(X, Y)$  Provided.

$$\textcircled{1} \quad P_{ij} \geq 0$$

$$\textcircled{2} \quad \sum_{i=1}^n \sum_{j=1}^m P_{ij} = 1$$

The set of triple  $(x_i, y_j, P_{ij})$  is called Joint Probability distribution of  $(X, Y)$ .

→ Marginal PMF of  $x$

$$P(X=x_i) = \sum_{j=1}^n P_{ij} = P_{i1} + P_{i2} + \dots$$

Similarly

The collection of pair  $(x_i, P_i)$  is called marginal PMF of  $x$ ,

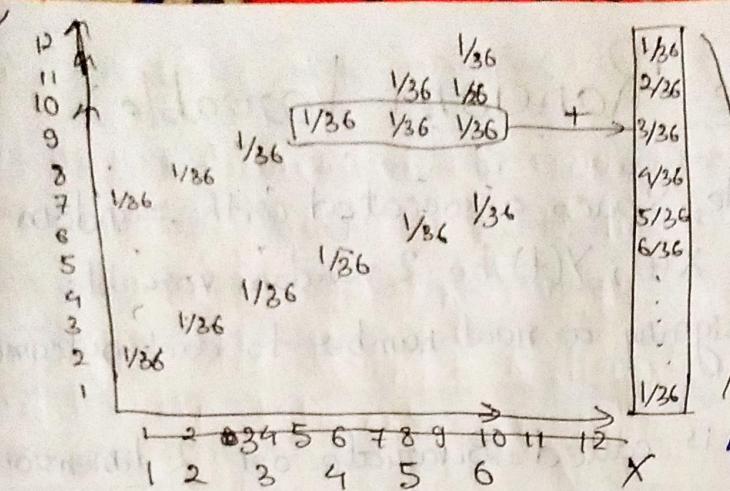
Marginal PMF of  $y$  :-

$$P(Y=y_j) = \sum_{i=1}^n P_{ij} = P_{j1} + P_{j2} + \dots$$

The pair  $(y_j, P_j)$  is called marginal PMF of  $y$ .

⇒ If  $X$  &  $Y$  are independent :-

$$P_{xy}(x, y) = P_x(x) P_y(y)$$



Similarly,

Marginal distribution

$X$ : The no. rolled on the 1st dice

$Y$ : Sum of the two dice

$$P_X(x_i) = \sum_j P_{XY}(x_i, y_j) \quad | \quad P_Y(y_j) = \sum_i P_{XY}(x_i, y_j)$$

→ Conditional distribution:

It is the probability distribution of a random variable, calculated according to the rules of conditional probability after observing the realization of another random variable.

$$P_{Y|X=x}(y) = P(Y=y | X=x) = \frac{P(X=x, Y=y)}{P(X=x)}$$

Taking,  $x=9$   
in  $y=49$

## \* Covariance:-

Covariance measures the direction of the relationship between two variables.

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

Eg:-

The first plot shows a positive linear relationship where both x and y increase as each other increases. The second plot shows no clear relationship where x and y are roughly constant or random. The third plot shows a negative linear relationship where x increases as y decreases.

$$\begin{aligned}\text{Cov}(x, y) &> 0 \\ \text{if } x \uparrow \text{ then } y \uparrow \\ \text{or } x \downarrow \text{ then } y \downarrow\end{aligned}$$

$$\text{Cov}(x, y) \approx 0$$

$$\begin{aligned}\text{Cov}(x, y) &< 0 \\ \text{if } x \uparrow \text{ then } y \downarrow \\ \text{or } x \downarrow \text{ then } y \uparrow\end{aligned}$$

## → Covariance Matrix:-

Covariance Matrix is a measure of how much 2 R.V. gets change together. It is actually used for computing the covariance in between every column of data matrix.

(Also known as Dispersion matrix  
or Variance-covariance matrix)

$$\begin{bmatrix} \text{Var}(x_1) & \dots & \text{Cov}(x_n, x_1) \\ \vdots & \ddots & \vdots \\ \text{Cov}(x_n, x_1) & \dots & \text{Var}(x_n) \end{bmatrix}$$

Eg:-

$$\begin{bmatrix} \text{Var}(x) & \text{Cov}(x, y) \\ \text{Cov}(x, y) & \text{Var}(y) \end{bmatrix}$$

for only (x, y)

$$\begin{bmatrix} \text{Var}(x) & \text{Cov}(x, y) & \text{Cov}(x, z) \\ \text{Cov}(x, y) & \text{Var}(y) & \text{Cov}(y, z) \\ \text{Cov}(x, z) & \text{Cov}(y, z) & \text{Var}(z) \end{bmatrix}$$

for only (x, y, z)

\* Pearson Correlation Coefficient :-

The Pearson Correlation Coefficient ( $r$ ) is the most common way of measuring a linear correlation.

It ranges between -1 and 1 that measures the strength & direction of the relationship between 2 variables.

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

(Best in linear)

$$-1 \leq f(x,y) \leq 1$$

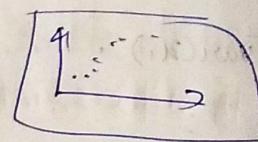
## \* Spearman Rank Correlation:-

While Pearson's correlation b/w the rank values of 2 R.V.s,

Spearman's correlation assesses monotonic relationships (whether linear or not).

It also ranges b/w -1 to 1.

$$\sigma_s = \text{Pr}(x), R(y) = \frac{\text{cov}(R(x), R(y))}{\sigma_{R(x)} \sigma_{R(y)}}$$



B91

$$g_3 = \frac{1 - 6 \sum d_i^2}{n(n^2 - 1)}$$

$$d_i = \text{rank } x_{1i} - \text{rank } y_i$$

## \* Multivariate Gaussian Distribution

Multivariate Gaussian distribution is a generalization of the one-dimensional (univariate) to higher dimension.

$$x = [x_1, x_2, \dots, x_n]^T$$

$$\mu = [\mu_1, \mu_2, \dots, \mu_n]^T$$

$$M = \begin{pmatrix} \mu_1 & \mu_2 & \dots & \mu_n \end{pmatrix}$$

$$\text{Mean}$$

$$f(x_1, x_2, \dots, x_n) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

Determinant of covariance matrix

## \* Moments

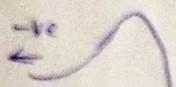
Moments are popularly used to describe the characteristic of a distribution.

Let's say the random variable of our interest is  $X$  then, moments are defined as the  $X$ 's expected values.

- ①  $\frac{\sum x_i}{n} \Rightarrow \text{Mean}$
- ②  $\frac{\sum x_i^2}{n} \xrightarrow{\text{(centralized)}} \frac{\sum (x_i - \bar{x})^2}{n} \Rightarrow \text{Variance}$
- ③  $\frac{\sum x_i^3}{n} \xrightarrow{\text{(standardized)}} \frac{\sum (x_i - \bar{x})^3}{n \cdot s^3} \xrightarrow{\text{Skew}} \frac{n}{(n-1)(n-2)} \frac{\sum (x_i - \bar{x})^3}{s^3}$
- ④  $\frac{\sum (x_i)^4}{n} \xrightarrow{\text{(population)}} \frac{1}{n} \left( \frac{\sum (x_i - \bar{x})^4}{s^4} \right) \xrightarrow{\text{(sample)}} \frac{n(n+1)}{(n-2)(n-3)} \left( \frac{\sum (x_i - \bar{x})^4}{s^4} \right) - \frac{3(n-1)^2}{(n-2)(n-3)}$

## → Skewness :-

Skewness is a degree of asymmetry observed in a probability distribution that deviates from the symmetrical normal distribution (bell curve) in a given set of data.



-ve skew



0 skew



+ve skew

### • Pearson's Method:-

#### ① Mode skewness:-

$$\text{skew} = \frac{\text{mean} - \text{mode}}{\sigma}$$

### • With moments :-

$$\frac{n}{(n-1)(n-2)} \times \frac{\sum (n-\bar{x})^3}{s^3}$$

#### ② Median skewness:-

$$\text{skew} = \frac{3(\text{mean} - \text{median})}{\sigma}$$

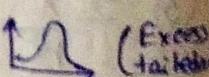
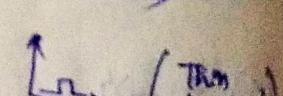
$$\therefore \text{mode} = 3\text{median} - 2\text{mean}$$

(In skewed data)

## → Kurtosis :-

Kurtosis is a measure of the tailedness of a distribution.

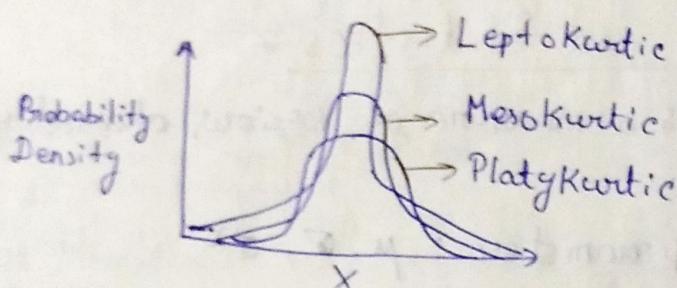
$$\text{Kurtosis (K)} = \frac{n(n+1)}{(n-1)(n-2)(n-3)} \left( \frac{\sum (n-\bar{x})^4}{s^4} \right) - \frac{3(n-1)^2}{(n-2)(n-3)}$$

- A normal distribution has  $K=3$ , (mesokurtic)
- Distribution with  $K > 3$ , (leptokurtic)  (Excess tailed)
- Distribution with  $K < 3$ , (platykurtic)  (Thin tailed)

→ Types of Kurtosis:-

	Mesokurtic	Platykurtic	Leptokurtic
Tailedness	Medium-tailed	Thin-tailed	Fat-tailed
Outlier frequency	Medium	Low	High
Kurtosis	Moderate ( $3$ )	Low ( $< 3$ )	High ( $> 3$ )
Excess Kurtosis	0	Negative (-ve)	Positive (+ve)

Example



# Sample & Sampling

\* Population :- A population is a distinct group of individuals, whether that group comprises a nation or group of people with a common characteristic. Instead of examining the entire group, called population.

\* Sample :- We may examine only a small part of this population which is called sample.

→ Parameter and Statistics :-

① Parameter is a measure of various characteristics of a population.

Population parameter :  $\mu, \sigma, \sigma^2$

② Statistics is a measure estimated through various sample of population.

Sample statistics :  $\bar{X}, S, S^2$

→ Sampling & Non-Sampling Error:-

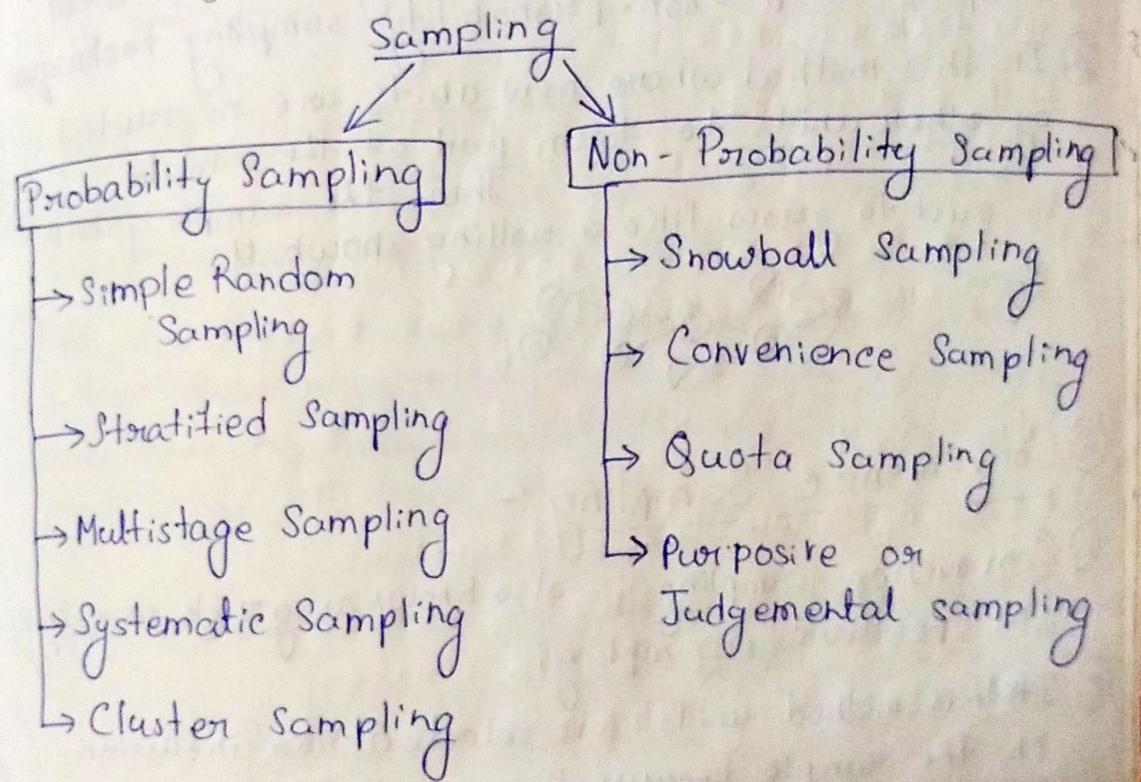
① The difference b/w parameters and statistics is known as the sampling error.

② These can't be controlled by humans.

③ These errors come into existence when various sample we taken.

- Non-Sampling Error :-
- ① There are various systematic errors and random errors that are not due to sampling.
  - ② These can be controlled by humans.  
Eg:- Errors in processing.
  - ③ These can control
  - ④ Non-Sampling errors are systematic error.

→ Sampling Techniques :-



\* When we use non-probability Sampling?

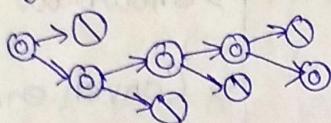
- ① Sample size should be small (less than 30)
- ② Biased (depend on research)
- ③ Qualitative research.
- ④ Also Known as non-random sampling.

→ Non-Probability Sampling Types:

① Snowball Sampling:-

- ① Snowball sampling / chain sampling / chain referral sampling is a non-probability sampling technique
- ② In this method where new units are recruited by other units to form part of the sample.  
Thus, the sample group is said to grow like a rolling snowball.

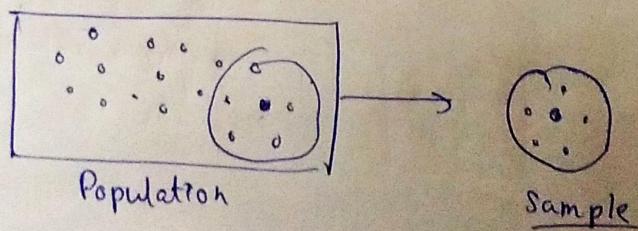
Eg:-



② Convenience Sampling:-

- ① Convenience Sampling is also known as grab sampling or opportunity sampling.
- ② In this method units are selected for inclusion in the sample because they are the easiest for the researcher to access.
- ③ This type of sampling is most useful for pilot testing.

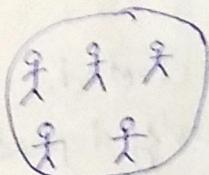
Eg:-



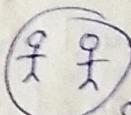
### ③ Quota Sampling :-

- ① Quota Sampling is the non-probability version of stratified sampling.
- ② In this assembled sample has the same proportions of individuals as the entire population with respect to known characteristics, traits or focused phenomena.

Eg:-



→  
Men above  
40



Sample

Population

### ④ Purposive Sampling :-

- ① Purposive sampling, also known as judgement, selective or subjective sampling.
- ② It is a form of non-probability sampling in which researchers rely on their own judgement when choosing members of the population to participate in their study.

Eg:-

↓ ↓ ↓



↓  
Sample of people  
living with diabetes

\* When we use probability sampling? (Big dataset)

- ① Almost no Biased  $\rightarrow$  Bias  $\approx 0$
- ② Quantitative Data
- ③ Sample size is large.

$\rightarrow$  Types of probability Sampling:-

① Simple Random Sampling (SRS):-

- ① Each item has equal chance to be selected.
- ② Here the selection of items completely depends on chance or by probability.

$\Rightarrow$  Two Types of SRS:-

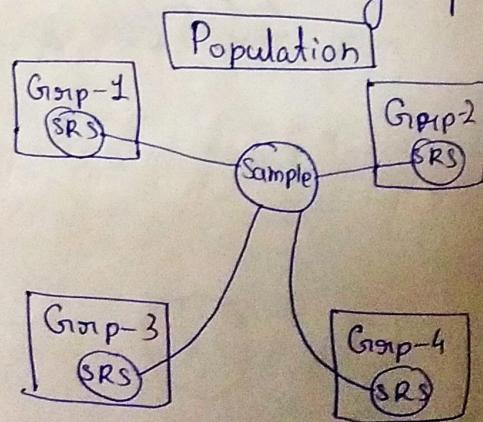
- ① With Replacement
- ② Without Replacement

Eg:- Choosing 25 employee names from a hat containing 250 employee names. (Every name has equal chance of being chosen)

② Stratified Sampling:-

- ① The samples are drawn at random.
- ② With stratified sampling researchers divides the population into separate groups called strata.
- ③ Then a probability sample is drawn from each group.

Eg:- Dividing the entire state into 3-4 geographical regions and then selecting required sample of districts from each of these regions.

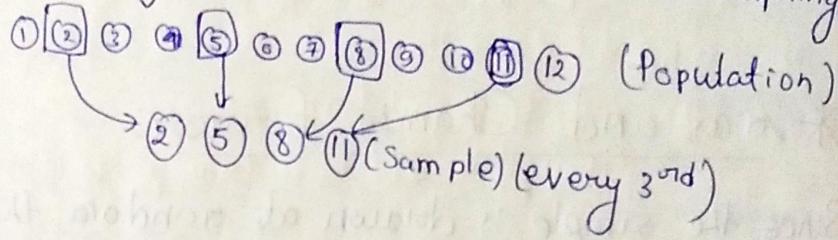


### ③ Systematic Sampling:-

- In this sample members from a larger population are selected according to random starting point but with a fixed, periodic interval.
- This interval is called the sample interval ( $K$ ).
- Sampling interval is calculated by dividing the population size ( $N$ ) by the desired sample size, ( $n$ )
$$K = \frac{N}{n}$$

④ It is functionally similar to simple random sampling,

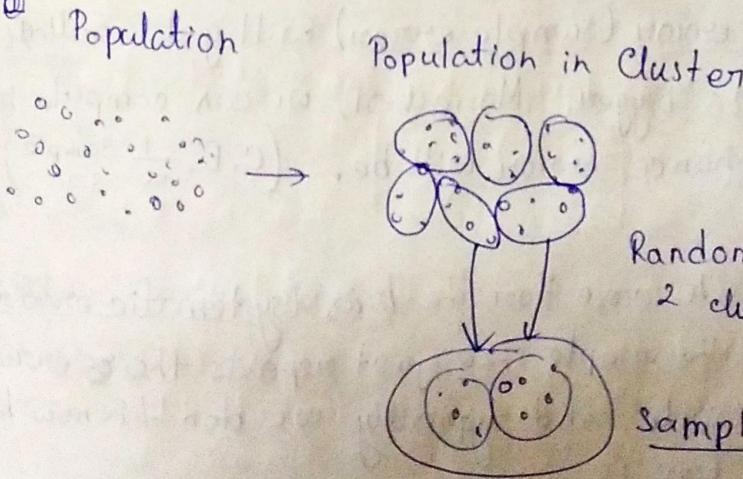
Eg:-



### ⑤ Cluster Sampling:-

A Probability sampling method in which we divide a population into cluster such as districts or schools and then randomly select some of these clusters as our sample.

Eg:-



## ② Multistage Sampling:

It is the taking of samples in stages using smaller and smaller sampling units in each stage.

Eg: Selection of households for the entire country.

1<sup>st</sup> stage unit  $\rightarrow$  states

2<sup>nd</sup> stage Units  $\rightarrow$  districts

3<sup>rd</sup> stage Units  $\rightarrow$  Villages

etc.

## \* Bias and Chance Error:

Since the sample is drawn at random, the estimate will be different from the parameter due to chance error.

Drawing another sample will result in a different chance error.

$$\text{Estimate} = \text{parameter} + \text{bias} + \text{chance error}$$

- The chance error (sample error) will get smaller as the size gets bigger. Moreover, we can compute how large the chance error will be. (C.E.  $\propto \frac{1}{\sqrt{\text{size}}}$ )
- This is not the case for the bias (systematic error): Increasing the sample size just repeats the error on a larger scale, and typically we don't know how larger the bias is.

## \* Observational Studies:

This is an observational study: I'

- It measures outcomes of interest and this can be used to establish association.
- But association is not causation, because there may be compounding factors such to other cause the association indirectly.

## \* Randomized Control Experiments / Trials (RCT):

RCT are prospective studies that measure the effectiveness of a new ~~intervention~~ treatment or treatment.

Although no study is likely on its own ~~to~~ prove causality, randomization reduces bias and provides a rigorous tool to examine ~~cause~~ cause-effect relations b/w an intervention and outcome.

→ The Logic of RCT :-

- ① It makes the treatment group similar to the control group. Therefore influences other than the treatment operate equally on both groups, apart from differences due to chance.
- ② It allows to assess how relevant the treatment effect is, by calculating the size of chance effects when comparing the outcomes in the 2 groups.

# \* Standard Error

The standard error is a statistical term that measures the accuracy with which a sample distribution represents a population by using standard deviation.

# The Standard Error of the mean is  
(the standard deviation of all the sample means  
from the actual population mean)  
(How far off from  $\mu$  will  $\bar{x}_n$  be?)

## The square root Law:

The accuracy of  $\bar{x}_n$  as an estimator of  $\mu$  (S.E.) is inversely proportional to the square root of the sample size  $n$ .

$$SE(\bar{x}_n) \propto \frac{1}{\sqrt{n}}$$

• The square root law is key for statistical inference

$$SE(\bar{x}_n) = \frac{\sigma}{\sqrt{n}}$$

→ The importance of the square root law is twofold:-

- ① It shows that the SE becomes smaller if we use a larger sample size  $n$ . We can use the formula to determine what sample size is required for a desired accuracy.
- ② The formula for the SE does not depend on the size of the population but only depends on the size of the sample

Note:-

There are 3 common error bars:-

- ① Standard deviation
- ② Standard error
- ③ Confidence Interval

## → Relation with Summation:-

The sum and the average are related by,  $S_n = n \bar{x}_n$

Similarly,

$$E(S_n) = n\mu ; \quad SE(S_n) = \sqrt{n}\sigma \quad \therefore SE = \frac{\sigma}{\sqrt{n}}$$

(Expected Value)

So, the variability of the sum of  $n$  draws increases at the rate  $\propto \sqrt{n}$ .

## → Relation with Percentages:-

Then, ~~SE~~  $\propto \frac{1}{\sigma}$

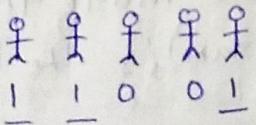
$$S.E. \propto \frac{1}{\sigma} \sqrt{n}$$

$\therefore S.E.(S_n) \propto \sqrt{n}$

Percentages:-  $E(\text{percentages of } 1s) = \mu \times 100\%$

$$SE(\text{percentages of } 1s) = \frac{\sigma}{\sqrt{n}} \times 100\%,$$

where we label people as 1 and 0.

Eg:-  then, sum = 3 =  $S_n$

## \* The Law of Large Numbers:-

The Law of Large numbers states that an observed sample avg. from a large sample will be close to the true population avg. and that it will get closer the larger the sample.

The Law of Large Numbers applies:-

Can be justified as according to the square root law

$$SE = \frac{1}{\sqrt{n}}$$

- ① for avg. and therefore percentages, but not for sums as their  $SE(S_n)$  increases.  $[SE(S_n) \propto \sqrt{n}]$
- ② for sampling with replacement from a population, or for simulating data from a probability histogram.

## \* The Central Limit Theorem:-

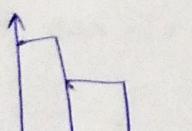
The Central Limit Theorem (CLT) states that the distribution of sample means approximates a normal distribution as the sample size gets larger, regardless of the population's distribution.

The sample size  $n$ . The more the skewed the population distribution (histogram) is, the larger the required sample size  $n$ .

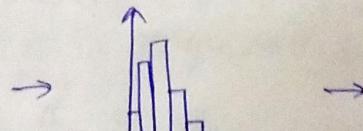
If there is no strong skewness then  $n \geq 15$  is sufficient otherwise  $n \geq 30$  or  $n \geq 40$  is sufficient in most cases

Eg:-

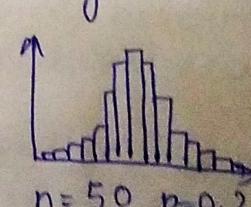
Binomial distribution (with winning probability,  $p=0.2$ )



$$n=1, p=0.2$$



$$n=10, p=0.2$$



$$n=50, p=0.2$$

## \* Normal Approximation:-

Finding percentages areas under the ~~normal curve~~ normal curve is called normal approximation.

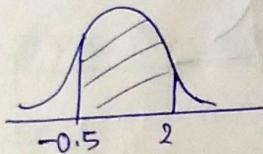
Eg:- Percentage of people with height b/w 67.4 in & 71.9 in?  
 $\mu_n = 68.3$  &  $\sigma_n = 1.8 \text{ in}$

→ Steps of Normal Approximation:-

① Standardize:

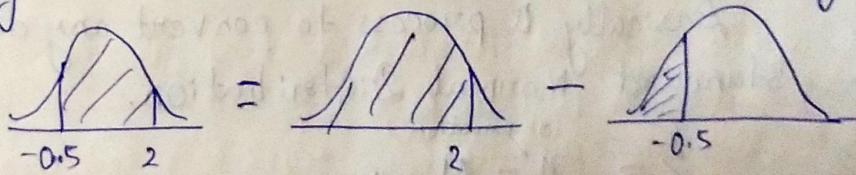
$$\frac{67.4 \text{ in} - 68.3 \text{ in}}{1.8 \text{ in}} = -0.5 ; \frac{71.9 \text{ in} - 68.3 \text{ in}}{1.8 \text{ in}} = 2$$

② Mark the area under the curve:-



③ Write the desired area in a form that can be computed by software:-

Typically we can look up the area to the left of a given value.



④ Use software or a table to find these values

$$97.7\% - 30.9\% = 68.8\%$$

→ Normal approximation to the binomial :-

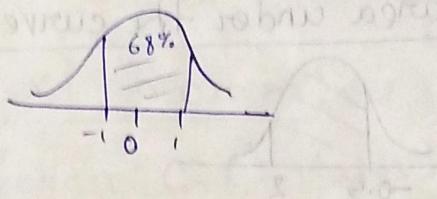
We can approximate binomial probabilities using normal approximation.

To standardize, subtract off  $np$  and then divide by  $\sqrt{np(1-p)}$

Eg: To find  $P(\text{at most } 12 \text{ small prizes})$   
We first standardize:  $\frac{12 - np}{\sqrt{np(1-p)}}$   $\rightarrow P(\text{win a small prize}) = 0.1$

### \* Standard Normal Distribution:

It is the Normal Distribution with  $\mu=0$  &  $\sigma=1$



### \* Standardization:-

Standardization is a scaling method where the values are centered around mean with a unit S.D ( $\sigma$ ).

Basically a process to convert any distribution to a standard Normal Distribution.

$$z = \frac{n - \mu}{\sigma} \quad \begin{matrix} (\text{observation}) \\ (\text{Mean}) \\ (\text{S.D.}) \end{matrix}$$

$(z\text{-score}) \rightarrow$

## \* Z-Score:

Z-Score is a statistical measurement that describes a value's relationship to the mean of a group of values.

Z-score is measured in terms of standard deviations from the mean.

Eg:-  $Z = 2.5$ , means  $2.5\sigma$  above  $\mu$

$Z = -1$ , means  $1\sigma$  below  $\mu$

## → Z-score Applications:-

① Standardization.

② Compare scores between different distributions.  
~~(After Standardization)~~

③ Eg:- Indian Cricket Team

2020

$$\mu = 181$$

$$\sigma = 12$$

2021

$$\mu = 182$$

$$\sigma = 5$$

Final Score = 187 runs

Final Score = 185 runs

$$Z_{2020} = \frac{187 - 181}{12} = \underline{0.5}$$

$$Z_{2021} = \frac{185 - 182}{5} = \underline{0.6}$$

∴  $Z_{2021} > Z_{2020}$ , Cricket Team did better in

2021.

## \* Confidence Interval:-

Confidence Interval, in statistics, refers to the probability that a population parameter will fall between a set of values for a certain proportion of times.

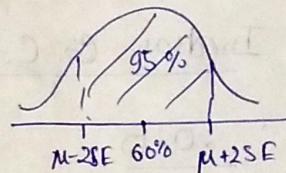
Eg:- If a point estimation is generated from a statistical model of 10 samples, with 95% confidence interval of 9.5 - 10.5, it can be inferred that there is a 95% probability that true value falls within that range.

• Confidence interval measures the degree of uncertainty or certainty in a sampling method.

Eg:- If we take a normal curve with  $\mu = 60\%$  &

$$SE = \frac{\sigma}{\sqrt{1000}} = 1.6\%$$

Then, 95% chance that the sample percentage is no more than 2SEs way from  $\mu$ , [54.8%, 61.2%]



## → Confidence Interval via CLT.

A confidence interval is centered at of plausible values for a population parameter  $\mu$ .

Usually confidence interval is centered at an estimated average.

Since the central limit theorem applies for averages, the confidence interval has a simple form:

$$\boxed{\text{estimate} \pm 2 \times S.E.}$$

→ Estimating the SE with the bootstrap:

SE is the standard error of the estimate.

$$SE = \frac{\sigma}{\sqrt{n}}$$

(But we don't know  $\sigma$ )

The bootstrap principle states that we can estimate  $\sigma$  by its sample version's and still get an approximately correct confidence interval.

$$SE = \frac{\sigma}{\sqrt{n}} \approx \frac{s}{\sqrt{n}}$$

⇒ Margin of error:

The width of the confidence interval is determined by  ~~$z\sigma$~~ , which is called margin of error.

It represents the degree of uncertainty that a survey result might have.

→ We can decrease the margin of error by:

- Increasing Sample size ( $n$ ). [We need  $4n$  to cut the MOE to  $\frac{MOE}{2}$ ]
- Decrease Confidence Interval.

(Trade off b/w Precision (MOE) and the Confidence Interval)

# \* Monte Carlo Method or (Simulation):

The Monte Carlo Method is a probabilistic method that can include an element of uncertainty or randomness in its prediction.

Monte Carlo simulation provides multiple possible outcomes and the probability of each from the large pool of random data samples.

→ In most cases the general problem ~~where~~ of interest is to find a unknown parameter  $\theta$  of a population.

→ We estimate  $\theta$  with a statistic (estimator)  $\hat{\theta}$  which is based on a sample of  $n$ . observations  $X_1, \dots, X_n$  drawn at random from the population:

$$\hat{\theta} = \text{avg of sample} = \frac{1}{n} \sum_{i=1}^n X_i$$

And  $\hat{\theta}$  tends to be close to the uncomputable population mean  $\theta$ , even for ~~so~~  $n=100$ .

OR,

## • Monte Carlo Method:

→ We approximate a fixed quantity  $\theta$  by the avg of independent random variables that have expected value  $\theta$ .

→ By the Law of large numbers, the approximation error can be made arbitrarily small by using large sample size ( $n$ ).

## Monte Carlo Method for S.E.

The Monte Carlo Method can also be used for more involved quantities, like S.E.

The process →

- Get many (say 1000) samples of 100 observations each.
- Compute  $\hat{\theta}$  for each sample, resulting in 1,000 estimates  $\hat{\theta}_1, \dots, \hat{\theta}_{1000}$
- Compute the standard deviation of these 1000 estimates.  
 $s(\hat{\theta}_1, \dots, \hat{\theta}_{1000}) = \sqrt{\frac{1}{999} \sum_{i=1}^{1000} (\hat{\theta}_i - \text{avg}(\hat{\theta}_i))^2}$

The Law of Large numbers is also applied here, making.  $s(\hat{\theta}_1, \dots, \hat{\theta}_{1000}) \approx SE(\hat{\theta})$

\*Note:-

Monte Carlo Method or Simulation only ~~exists~~ works when we can sample as much as we want.

→ Plug-in Principle:-

The plug-in Principle is a technique used to approximately compute or to estimate a feature of probability distribution (E.g. S.E., E.V.,  $\sigma$ ) that cannot be computed exactly. (Based on a sample distribution)

## \* The Bootstrap principle

- The basic idea of Bootstrapping is that inference about a population from sample data (sample  $\rightarrow$  population) can be modeled by resampling the sample data and performing inference about a sample data (resampled  $\rightarrow$  sample)
- The bootstrap uses the plug-in principle and the Monte Carlo Method to approximate quantities such as  $SE(\hat{\theta})$ .
- Simulating a bootstrap sample  $x_1^*, \dots, x_n^*$  means that we draw  $n$  times with replacement from  $x_1, \dots, x_n$ .

Eg:- from  $\underline{1 \ 2 \ 3 \ 4}$ ,

$\rightarrow 1 \ 2 \ 3 \ 3$   
 $\rightarrow 1 \ 2 \ 1 \ 3$   
 $\rightarrow 4 \ 4 \ 3 \ 1$   
etc.

\* The Bootstrap consists of 2 steps :-

- (1) Draw  $N$  bootstrap samples and compute  $\hat{\theta}^*$  for each bootstrap sample:

$$x_1^{*1}, \dots, x_n^{*1} \rightarrow \hat{\theta}_1^*$$

$$x_1^{*N}, \dots, x_n^{*N} \rightarrow \hat{\theta}_N^*$$

- (2) Use  $\hat{\theta}_1^*, \dots, \hat{\theta}_N^*$  to approximate the quantity of interest.

Eg:- We approximate  $SE(\hat{\theta})$  by the standard deviation of  $\hat{\theta}_1^*, \dots, \hat{\theta}_N^*$

→ Parametric Bootstrap:-

In this method we assume that the observed data is drawn from a specific parametric distribution (e.g:- Normal Distribution, exponential distribution, etc).

Then, we estimate the parameters of this distribution from our sample data, and generate new data samples from the assumed distributions.

→ Non Parametric Bootstrap:-

Non Parametric Bootstrap makes no specific assumptions about the underlying distribution of the data. It directly resamples observed data.

## \* Statistical Inferences :-

Statistical inference is a fundamental concept in statistics that involves drawing conclusions or making predictions about a population based on a sample of data.

It is a way of using sample data to make inferences or generalizations about a larger population from which the sample is drawn.

→ Statistical Inference can be approached from a Bayesian or Frequentist perspective.

⇒ Bayesian Inference uses prior knowledge and probability distributions to update beliefs, while, Frequentist Inference relies solely on sample data.

The choice b/w these approaches depends on the context and available information.

→ There are mainly 2 types Statistical Inferences:-

① ~~Point~~ Estimation

    └ (i) Point Estimation

    └ (ii) Interval Estimation

② Hypothesis Testing

### \* Estimation:-

A process in which we obtain the value of unknown population parameters with the help of sample data.

#### → Point Estimation:-

Point estimation involves providing a single value, often denoted as a point estimate, ~~as~~ as the best guess for the population parameter.

(Eg:  $\bar{x}$ ,  $s^2$ ,  $\hat{\sigma}$ )

#### → Interval Estimation:- (Confidence Interval)

Interval estimation provides a range of values within which the population parameter (Eg:  $\mu$ ) is likely to fall.

It is expressed as a Confidence Interval.

↙ [Watch Statquest for more clarity]

### \* Maximum Likelihood Estimator (MLE) :-

(Part of Point Estimation)

The MLE is a statistical method used to estimate the parameters of a statistical model. It is a powerful and widely used ~~the~~ technique in statistics and Machine Learning.

The key idea of MLE is to find the set of parameter values that maximizes the likelihood function, measures how well the observed data fits the model, or to find the parameter values that make the observed data the most probable under the assumed statistical model.

$$\text{MLE}(\theta) = \arg \max_{\theta} L(\theta / \text{data}), \text{ where,}$$

- $\hat{\theta}$  is Maximum Likelihood estimate of the parameter  $\theta$ .
- $L(\theta / \text{data})$  is the Likelihood function.

→ General steps to find MLE:-

- ① Define the Likelihood function.
- ② Take the natural logarithm for the equation.  
(Makes it easier to differentiate)
- ③ Differentiating the log likelihood.
- ④ Set the Score to zero and solve for parameter.  
Additionally,
- ⑤ We can check for 2<sup>nd</sup> derivative of log-likelihood  
to determine the nature of the critical point.  
(Maxima, Minima & Saddle Point)

⇒ MLE for Bernoulli Distribution:-

$$① l(p, x) = p^{\sum_{i=1}^n x_i} \times (1-p)^{n - \sum_{i=1}^n x_i}$$

$$② \log(l(p, x)) = \log\left(p^{\sum_{i=1}^n x_i} \times (1-p)^{n - \sum_{i=1}^n x_i}\right)$$
$$= \cancel{\log(l)} \left( \sum_{i=1}^n x_i \log(p) + (n - \sum_{i=1}^n x_i) \log(1-p) \right)$$

$$③ \frac{d}{dp} (\log(l(p, x))) = \frac{d}{dp} \left[ \left( \sum_{i=1}^n x_i \right) \log(p) + \left( n - \sum_{i=1}^n x_i \right) \log(1-p) \right]$$

$$④ \text{And } ⑤ \quad = \frac{\sum_{i=1}^n x_i}{p} + \frac{n - \sum_{i=1}^n x_i}{1-p} (1) = 0$$

$$\Rightarrow \sum_{i=1}^n x_i - p \times \cancel{\sum_{i=1}^n x_i} - np + p \times \cancel{\sum_{i=1}^n x_i} = 0$$

$$\Rightarrow np = \sum_{i=1}^n x_i$$

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

General form of finding MLE:

Let  $\theta$  be the population Parameter:

PD

$\downarrow$

$$\textcircled{1} \quad L(\theta) = f(x_1) \cdot f(x_2) \cdots f(x_n)$$

$$= \prod_{i=1}^n f(x_i)$$

$$\textcircled{2} \quad \log(L(\theta)) = \sum_{i=1}^n \log(f(x_i))$$

$$\begin{aligned} & \because \log(f(x_1) \cdot f(x_2) \cdots f(x_n)) \\ & = \log(f(x_1)) + \log(f(x_2)) \\ & + \cdots + \log(f(x_n)) \end{aligned}$$

$$\textcircled{3} \quad \frac{d}{d\theta} \log(L(\theta)) = 0$$

and solve for  $\theta$ .

→ MLE for Normal Distribution:

$$x_i \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$$

Note:-  
 $\sigma \neq 0$

$$\textcircled{1} \quad L(\mu, \sigma, x) = \prod_{i=1}^n f(x_i | \mu, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}}$$
$$= \frac{1}{(\sqrt{2\pi})^n \sigma^n} e^{-\frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2}}$$

$$\textcircled{2} \quad \log(L(\mu, \sigma, x)) = -\frac{n}{2} \log(2\pi) - n \log \sigma - \frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2}$$

To find  $\frac{\partial}{\partial \mu}$

$$\textcircled{3} \quad \frac{\partial \lambda(\mu, \sigma)}{\partial \mu} = -\frac{1}{2} \sum_{i=1}^n \frac{2(x_i - \mu)}{\sigma^2} + 1$$

To find  $\frac{\partial}{\partial \sigma}$

$$\textcircled{3} \quad \frac{\partial \lambda(\mu, \sigma)}{\partial \sigma} = -\frac{n}{\sigma} - \frac{1}{2} \left( \sum_{i=1}^n (x_i - \mu)^2 \right) \frac{1}{\sigma^3}$$

$$\textcircled{4} \quad \textcircled{5} \quad \frac{1}{\sigma^2} \left( \sum_{i=1}^n x_i - n\mu \right) = 0 \quad (\because \sigma \neq 0)$$

$$\textcircled{4} \quad \textcircled{5} \quad -\frac{n}{\sigma} + \left( \sum_{i=1}^n (x_i - \mu)^2 \right) \frac{1}{\sigma^3} = 0$$

$$\Rightarrow \sum_{i=1}^n x_i = n\mu$$

$$\Rightarrow -n + \left( \sum_{i=1}^n (x_i - \mu)^2 \right) \frac{1}{\sigma^2} = 0$$

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

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

Sample mean

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

$$\therefore \sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{n}} \quad \because \text{By estimate}$$

$$\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$$

$\bar{x}$

# \* Types of approaches for Statistical Inference

## Frequentist

① Probabilities represent long term frequencies

② Parameters are fixed (but unknown) constants, so we can make probability statement about them.

③ Find the model that better explains the observed data.

④ Statistical procedures have well-defined long run frequency properties

## Bayesian

① Probabilities represent a degree of belief, or certainty.

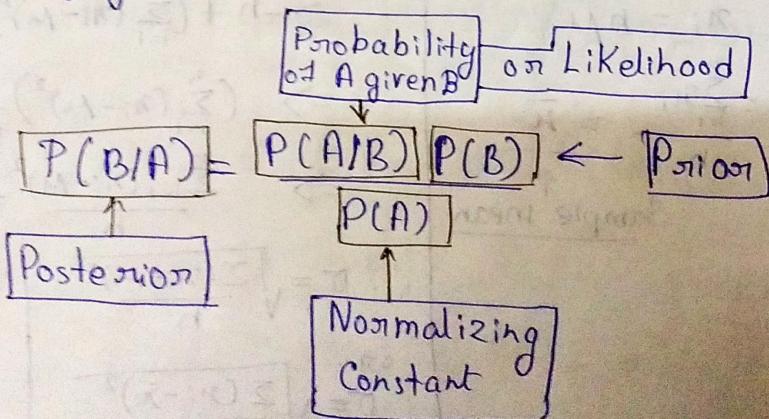
② We make probability statement about parameters, even if they are constant (beliefs)

③ Update beliefs on the model based on the observed data.

④ We make inferences about a parameter  $\theta$  by producing a probability distribution of  $\theta$ .

Basically, Frequentists imply that probabilities are objective properties of appearance of a certain event in the real world and that the parameters of the distributions are fixed constants. So, Frequentist interpret observed data as the sample of an unknown distribution.

While, Bayesians interpret probabilities as a degree of belief. The goal of Bayesians is to update this belief as we gather new data. Our result will be an updated probability distribution for the parameter we are trying to infer.



## \* Maximum a Posteriori (MAP):

Maximum a Posteriori (MAP) is an inference technique used to estimate the most likely value of an unknown variable on set of variables on the observed data and prior probability distribution.

It is used ~~as~~ in the context of Bayesian Inference and estimates the value that maximizes the posterior probability distribution.

$$MAP(\theta) = \arg \max_{\theta} P(\theta | D), \text{ where,}$$

$\theta$  = Parameter of interest

$P(\theta | D)$  = Posterior probability of parameter given  $D$ .

$P(\theta)$  = Prior probability of the parameter

$P(D|\theta)$  = Likelihood of observing data given the parameter

OR

$$\text{MAP: } \hat{\theta} = \arg \max_{\theta} f_{\theta | X=x}$$

Since, MAP considers both the data and prior information to provide more data, while MLE only considers the data.

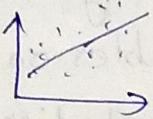
When there is no informative prior,  
MAP is same as MLE

# \* Regularization and MAP:

Using  $L_2$  regularization  
 $L_2 = a^2 + b^2$   
 $y = ax + b + \epsilon$

Scenario:-

Underfit



Loss

10

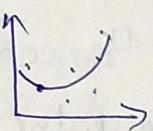
Equation

$$y = 4x + 3$$

Penalty

$$L_2 = 4^2 = 16$$

Good Fit

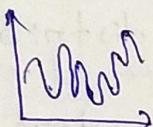


2

$$y = 2x^2 - 4x + 5$$

$$L_2 = 2^2 + (-4)^2 = 20$$

Overfit



0.1

$$y = 4x^{10} - 9x^8 - 2x^6 + 3x^5 - 6x^4 - 10x^3 + 4$$

$$L_2 = 4^2 + (-9)^2 + (-2)^2 + 3^2 + (-6)^2 + (-10)^2 = 246$$

New Loss  
(Loss + Penalty)

26

22

246.1

← Best-fit

→ Regularization Term:

Model:  $y = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$

Log-loss: ~~ll~~

$L_2$  Regularization Error:  $a_n^2 + a_{n-1}^2 + \dots + a_1^2$

Regularization parameter:  $\lambda$

Regularization Error:  $ll + \lambda(a_n^2 + a_{n-1}^2 + \dots + a_1^2)$

- Take MAP into account, we can use the prior information like simplicity of the model assuming simpler models are more likely to get selected.

Then, Likelihood of the Model,

$$\rightarrow P(\text{Model 1}) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}a_1^2}$$

$a_1 n + b$

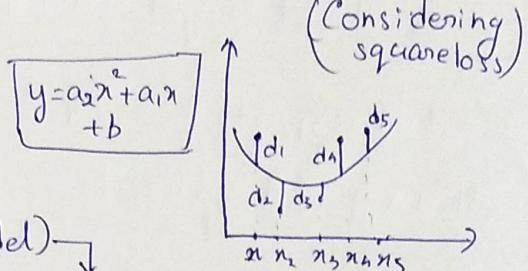
$$\rightarrow P(\text{Model 2}) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}a_1^2} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}a_2^2}$$

$a_1 n^2 + a_2 n + b$

$$\rightarrow P(\text{Model 3}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}a_i^2}$$

$a_1 n^0 + \dots + a_{10} n + b$

Then,



$$P = P(\text{Data} | \text{Model}) \times P(\text{Model})$$

$$= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2)} \times \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(a_1^2 + a_2^2)}$$

Then, to maximize  $P$  (Now finding the MAP)

$$\Rightarrow \log P = -\frac{1}{\sqrt{2\pi}} (d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2) - \frac{1}{\sqrt{2\pi}} (a_1^2 + a_2^2) = 0$$

$$\Rightarrow \underbrace{d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + a_1^2 + a_2^2}_{\downarrow \text{Minimize}} = 0$$

So, to maximize  $P$ , we need to minimize

$$\underbrace{d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2}_{\text{Square loss}} + \underbrace{a_1^2 + a_2^2}_{\text{Complexity of model (Lesser coefficients of n)}}$$