

# STATISTICS

## Outline # Probability

### Sample Space:

A set consisting of all possible outcomes that can result from a random experiments (real or conceptual), is defined to be a sample space for the experiment and is ~~defined~~ denoted by the letter S. Each possible outcome is a member of the sample space and is called a sample point in that space-

### Example:

Experiment of tossing a coin results in either of the two possible outcomes : a head (H) or a tail (T); Landing on its edge or rolling away is not considered. The sample space for this experiment is set notation as  $S = \{H, T\}$ .

# Events =



An event is an individual outcome or any number of outcomes (sample points) of a random experiment or a trial. In set terminology, any subset of a sample space  $S$  of the experiment, is called an event. An event that contains exactly one sample point, is defined a sample event. A compound event contains more than one sample point and is produced by the union of sample events.

## Counting Sample Events:

When the number of sample points in a sample space  $S$  is very large, it becomes very inconvenient and difficult to list them all and to count the number of points in the sample space  $S$  and in subsets of  $S$ . We then need some rules which helps us to count number of all sample points without actually listing them.

## ⇒ Rules / Methods:

### ij: Rule of Multiplication:

Compound Experiment consists of two experiments such that the first experiment has exactly  $n$  distinct outcomes and, if corresponding to each outcome of the first experiment there can be  $m$  distinct outcomes of the second experiment, then the compound experiment has exactly  $mn$  outcomes.

#### Example:

The compound experiment of tossing a coin and throwing a die together consist of two experiments - The total number of possible outcomes are  $2 \times 6 = 12$ .

### iij: Rule of Permutation:

A permutation is any ordered subset from a set of  $n$  distinct objects. The number of permutation of  $r$  objects, selected in a definite order from  $n$  distinct objects is denoted by the symbol  ${}^n P_r$ .

$${}^n P_r = \frac{n!}{(n-r)!}$$

### iii. Rules of Combination:

A combination is any subset of  $r$  objects, selected without regard to their order, from a set of  $n$  distinct objects. The total number of such combinations is denoted by the symbol  ${}^n C_r$  or  $\binom{n}{r}$ .

where  $r \leq n$ .

$$\binom{n}{r} = \frac{n!}{r!(n-r)!}$$

$${}^n C_r = \frac{n!}{r!(n-r)!}$$

### Probability of an Event:

If an event  $A$  is defined in a sample space  $S$ , then its probability will be equal to the sum of the probabilities of all sample points that are included in  $A$ .

i.e;  $P(A) = \sum P(E_i)$  - When all  $n$  possible outcomes of a random experiment are equally likely to occur, then

$$P(A) = \frac{m}{n} = \frac{\text{Number of sample points in } A}{\text{Number of sample points in } S} = \frac{m(A)}{n(S)}$$

$\Rightarrow \frac{1}{n}$  is the probability of each outcome

and its remains same-

## Additive Law:

If A and B are any two events defined in a sample space S, then

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

### Proof:

Two events  $A \cup B$  may be written as the union of two mutually exclusive events  $A$  and  $B \cap \bar{A}$ . That is.

$$A \cup B = A \cup (B \cap \bar{A})$$

Then

$$P(A \cup B) = P(A) + P(B \cap \bar{A}) \quad \text{--- (i)}$$

Again the event B may also be decomposed into two mutually exclusive events as-

$$B = (A \cap B) \cup (\bar{A} \cap B)$$

$$P(B) = P(A \cap B) + P(\bar{A} \cap B) \quad \text{--- (ii)}$$

By subtracting (ii) and (iii):

$$P(A \cup B) - P(B) = P(A) - P(A \cap B)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

Hence proved -

# Conditional Probability:

It is the probability of one event occurring with some relationship to one or more other events.

## Example:

- Event A is that it is raining outside, and it has a 0.3 [30%] chance of raining today.
- Event B is that you will need to go outside, and that has a probability of 0.5 [50%].

A conditional probability would look at these two events in relationship with one another, such as the probability that it is both raining and you will need to go outside.

## Formula:

$$P(A|B) = \frac{\text{number of sample points in } A \cap B}{\text{number of sample points in } B}$$

$$= \frac{n(A \cap B)}{n(B)}$$

we get

$$P(A|B) = \frac{n(A \cap B)}{n(S)} \cdot \frac{n(S)}{n(B)} = \frac{P(A \cap B)}{P(B)}$$

so, if  $P(B) = 0$ , the conditional probability  $P(A|B)$  remains undefined.

Similarly,

$$P(B|A) = \frac{P(A \cap B)}{P(A)} \quad \text{where } P(A) > 0.$$

## Baye's Theorem =

$$P(A_i|B) = \frac{P(A_i) P(B|A_i)}{\sum_{i=1}^k P(A_i) P(B|A_i)}$$

for  $i = 1, 2, 3, \dots, k$

Proof:

By multiplicative law of probabilities, we have-

$$P(B \cap A_i) = P(B) P(A_i|B)$$

$$P(B \cap A_i) = P(A_i) P(B|A_i)$$

Evaluating the equivalent relations of  $P(B \cap A_i)$  and

Solving for  $P(A_i|B)$  we get

$$P(B)P(A_i|B) = P(A_i)P(B|A_i)$$

$$P(A_i|B) = \frac{P(A_i)P(B|A_i)}{P(B)}$$

We may write  $B$  as  $B = S \cap B$ :

$$B = (A_1 \cup A_2 \cup A_3 \cup \dots \cup A_k) \cap B$$

$$B = (A_1 \cap B) \cup (A_2 \cap B) \cup (A_3 \cap B) \cup \dots \cup (A_k \cap B)$$

$$\text{Now } P(B) =$$

$$P(B) = P(A_1 \cap B) + P(A_2 \cap B) + \dots + P(A_k \cap B)$$

$$P(B) = \sum_{i=1}^k P(A_i \cap B)$$

Using the multiplicative law of probability, we may express each term  $P(A_i \cap B)$  as

$$P(A_i)P(B|A_i) \text{ Then:}$$

$$P(B) = P(A_1)P(B|A_1) + P(A_2)P(B|A_2) + \dots + P(A_k)P(B|A_k).$$

$$P(B) = \sum_{i=1}^k P(A_i)P(B|A_i) -$$

So

$$P(A_i|B) = \frac{P(A_i)P(B|A_i)}{\sum_{i=1}^k P(A_i)P(B|A_i)}$$

Hence proved

# Random Variables

## Outline #

### Concept Of Random Variables:-

Such a numerical quantity whose value is determined by the outcome of a random experiment is called a random variables.

Mathematically we assign a single value / real number to each outcome of the sample space, and hence we state that a random variable is a real valued function defined on a sample space -

A random variable is also called a chance variable, a stochastic variable or simply a variate and is abbreviated as r.v. The random variables are usually denoted by capital Latin letters such as  $X, Y, Z$ ; while the values taken by them are represented by the corresponding small letters such as  $x, y, z$ . There are two types of random variables:

i) Discrete-

ii) Continuous:

# Discrete Probability Distribution:

A discrete distribution describes the probability of occurrence of each value of a discrete random variable - A discrete random variable is a random variable that has a countable values, such as a list of non-negative integers -

With a discrete probability distribution each possible value of the discrete random variables can be associated with a non-zero probability - Thus a discrete probability distribution is often presented in tabular form -

# CONTINUOUS PROBABILITY DISTRIBUTION

A continuous distribution describes the probability of the possible values of a continuous random variable.

A continuous random variable is a random variable with a set of possible values that is infinite and uncountable.

Probabilities of a continuous random variable ( $x$ ) are defined as the area under the curve of its PDF - Thus only ranges of values can have a nonzero probability - The probability that a continuous random variable equals some values is always zero.

**Probability Density Function:** A probability density function (p.d.f) has the following properties -

$$i) f(x) \geq 0, \text{ for all } x$$

$$ii) \int_{-\infty}^{\infty} f(x) dx = 1$$

iii) The probability that  $x$  takes on a value in the interval  $[c, d]$ ,  $c < d$  is given by -

$$P(c < x \leq d) = F(d) - F(c)$$

$$= \int_{-\infty}^d f(u) du - \int_{-\infty}^c f(u) du$$

$$= \int_c^d f(u) du$$

which is the area under curve.

# Joint Distributions:

The distribution of two or more random variables which are observed simultaneously when an experiment is performed is called their joint distribution. It is necessary to call the distribution of a single r.v. as univariate. Likewise, a distribution involving two, three or many r.v.'s simultaneously is referred to be a bivariate, trivariate or multivariate.

## ⇒ Bivariate Distribution Function:

Let  $x$  and  $y$  be two r.v.'s defined on the same sample space  $S$ . Then, the function  $f(x,y)$  defined by  $F(x,y) = P(X < x \text{ and } Y < y)$ , where  $F(x,y)$  gives the probability that  $X$  will take on a value less than or equal to  $x$  and at the same time,  $Y$  will take on a value less

than or equal to  $y$  is called a bivariate or joint distribution function of  $x$  and  $y$ .

## ⇒ Bivariate Probability Function:

Let  $x$  and  $y$  be two discrete r.v's defined on the same sample space  $S$ ,  $x$  taking the values  $u_1, u_2, \dots, u_n$  and  $y$  takes on the value  $y_1, y_2, \dots, y_n$ .

Then the probability that  $x$  takes on the values  $u_i$  and  $y$  takes on the values  $y_j$ , denoted by  $f(u_i, y_j)$  of  $P_{ij}$  is defined to be

Bivariate Probability Function or joint probability function and its values at point  $(u_i, y_j)$  is given by

$$f(u_i, y_j) = P(x=u_i \text{ and } y=y_j).$$

## ⇒ Conditional Probability Functions:

Let  $x$  and  $y$  be two discrete r.v's with joint probability function  $f(u_i, y_j)$ . Then the condition probability for  $x$  given  $y=y_j$ , denoted as  $f(x|y_j)$ , is defined by -

$$f(u_i|y_j) = P(x=u_i | y=y_j)$$

$$= \frac{P(X = x_i \text{ and } Y = y_j)}{P(Y = y_j)}$$

$$= \frac{f(x_i, y_j)}{h(y_j)} \quad \text{for}$$

$$i = 1, 2, \dots$$

$$j = 1, 2, \dots$$

where  $h(y_j)$  is the marginal probability  
and  $h(y_j) > 0$ .

# Continuous Probability

## Distribution:

### Outline #6

#### Continuous Uniform Distribution:

The density function of a continuous r.v.  $X$  is called a uniform distribution when between the end points any two subintervals of the same length containing  $X$ , have the same probability.

Alternatively, a r.v.  $X$  is said to be uniformly distributed if its density function is defined as:

$$f(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{elsewhere} \end{cases}$$

The distribution derives its name from the fact that its density is constant or uniform over the interval  $[a, b]$  and is 0 elsewhere.

#### Properties:

i) Let  $X$  be the uniform distribution over  $[a, b]$ . Then its mean is  $a+b/2$  and variance is  $(b-a)^2/12$ .

Now  $E(u) = \int_{-\infty}^{\infty} f(u) u du = \frac{1}{b-a} \int_a^b u du$

$$= \frac{1}{b-a} \left[ \frac{u^2}{2} \right]_a^b$$

$$= \frac{b^2 - a^2}{2(b-a)}$$

$$= \frac{a+b}{2} \Rightarrow \text{mid-point}$$

And

$$\text{Var}(x) = E(x^2) - [E(x)]^2$$

$$E(x^2) = \int_a^b u^2 \cdot \frac{1}{b-a} du$$

$$= \frac{1}{b-a} \left[ \frac{u^3}{3} \right]_a^b$$

$$= \frac{b^3 - a^3}{3(b-a)}$$

$$= \frac{(a^2 + ab + b^2)(b-a)}{3(b-a)}$$

$$= \frac{a^2 + ab + b^2}{3}$$

$$\therefore \text{Var}(x) = E(x^2) - [E(x)]^2$$

$$= \frac{a^2 + ab + b^2}{3} - \frac{(a+b)^3}{4}$$

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

2j: The shape of distribution is rectangular.

## NORMAL DISTRIBUTION

The normal distribution is a continuous probability distribution that is symmetrical on both sides of the ~~center is~~ mean, so the right side of the center is a mirror image of the left side. The area under the curve of normal distribution represents probability and the total area under the curve sums to one.

$$f(u) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{u-\mu}{\sigma}\right)^2}$$

where,

$f(u)$  = probability density function-

$\sigma$  = standard deviation-

$\mu$  = mean-

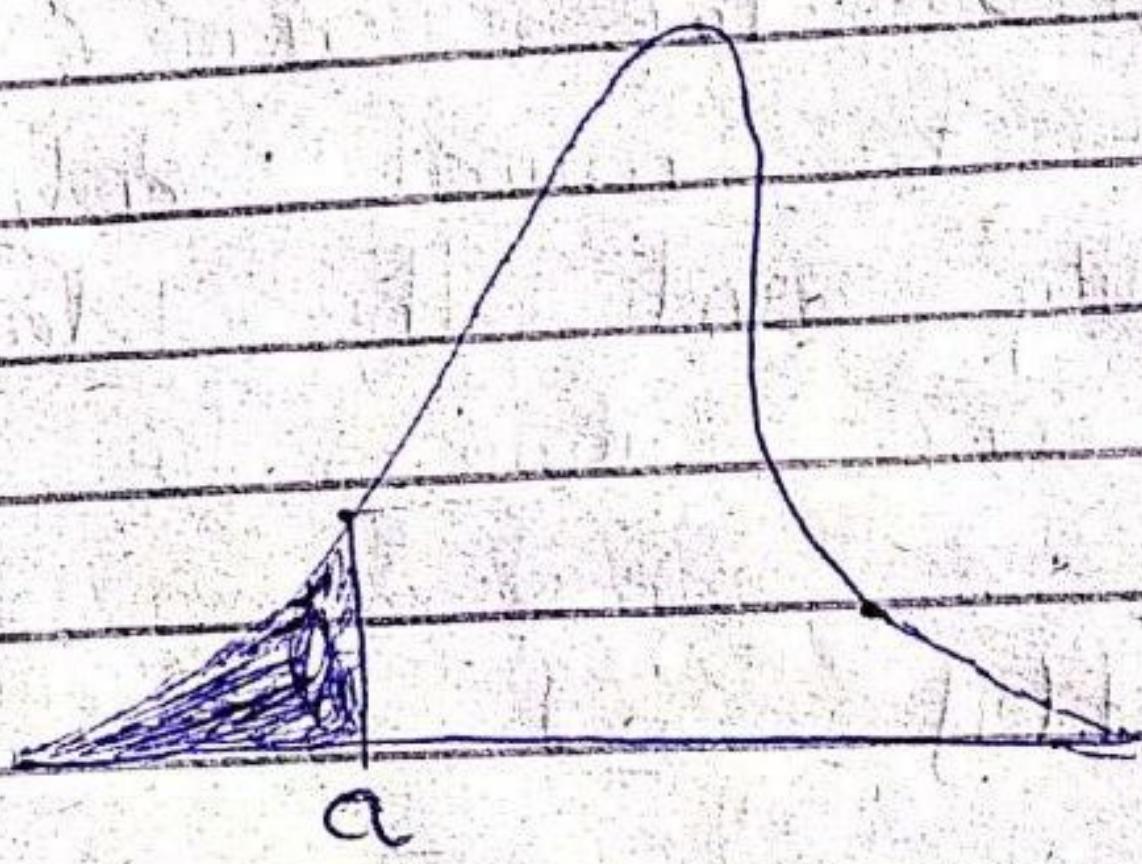
### Area Under the Normal Curve:

o) The total area under the normal curve is equal to 1.

•) The probability that a normal random variable  $x$  equals any particular value is 0.

e) The probability that  $X$  is greater than  $a$  equals the area under the normal curve bounded by  $a$  and infinity (indicated by non-shaded area)-

f) The probability that  $X$  is less than  $a$  equals the area under the normal curve bounded by  $a$  and minus infinity (indicated by shaded area)-



## Applications of Normal Distribution:

OR

## Properties =

The main properties of normal distribution are:-

i) The function  $f(x)$  defining the normal distribution is a proper pdf ie  $f(x) \geq 0$

and the total area under the normal curve is unity.

2) = The mean and variance of the normal distribution are  $\mu$  and  $\sigma^2$ .

3) = The median and mode of the normal distribution are each equal to  $\mu$ , the mean of distribution.

4) = The mean deviation of the normal distribution is approximately  $4/5$  of its standard deviation.

5) = The normal curve has points of inflection which are equidistant from the mean.

6) = For the normal distribution, the odd order moments about the mean are all zero and even order moments are given by =

$$\mu_{2n} = (2n-1)(2n-3) \dots 5, 3, 1, \sigma^{2n}$$

7) = If  $X$  is  $N(\mu, \sigma^2)$  and if  $Y = a + bx$  then  $Y = N(a + b\mu, b^2\sigma^2)$

8) = The sum of independent normal variables is a normal variable.

# Normal Approximation to the Binomial Distribution:

The binomial distribution  $b(u; n, p)$  can be closely approximated by the normal distribution when  $n$  is sufficiently large and neither  $p$  nor  $q$  is close to zero. As a rule of thumb, normal distribution provides a responsible approximation to the binomial distribution if both  $np$  and  $nq$  are equal to or greater than 5.

Now, the probability for a binomial random variable  $X$  to the value  $u$  is:

$$F(u) = \binom{n}{u} p^u q^{n-u}$$

For  $0 \leq u \leq n$  and  $q/p=2$

# Gamma Distribution:

A continuous r.v.  $X$  is said to have a gamma distribution with parameter  $m > 0$ , if its p.d.f is defined by

$$f(u) = \begin{cases} \frac{1}{\Gamma(m)} u^{m-1} e^{-u} & \text{for } 0 \leq u < \infty \\ 0, & \text{otherwise} \end{cases}$$

A gamma variable with parameter  $m$  is actually usually denoted by  $\gamma(m)$ . A straightforward integration shows that  $\int_0^\infty f(u) du = 1$ , and hence it represents a p.d.f.

The distribution function  $F(x)$  is

$$F(x) = \int_0^x \frac{1}{\Gamma(m)} u^{m-1} e^{-u} du, \quad u \geq 0$$

which is also called incomplete gamma distribution.

## ⇒ Properties:

1) The mean and variance of gamma distribution is equal to its parameters  $m$ .

as  $\mu = m$

and  $\sigma^2 = m$

2) The sum of two independent Gamma distributions with parameter  $m$  and  $n$  is Gamma distribution with parameter  $(m+n)$ .

## Exponential Distribution:

A random variable  $X$  is said to have an exponential distribution with parameter  $\lambda$ , if its p.d.f is defined by -

$$f(x) = \lambda e^{-\lambda x} \quad \text{for } x > 0 \\ = 0$$

where  $\lambda > 0$ . The p.d.f may also written as

$$f(x) = \frac{1}{\beta} e^{-x/\beta} \quad \text{for } x > 0$$

The function  $f(x)$  is a proper p.d.f so

$$\int_{-\infty}^{\infty} f(x) dx = \int_0^{\infty} \lambda e^{-\lambda x} dx = [-e^{-\lambda x}]_0^{\infty} = 1$$

The distribution function of the exponential r.v.  $X$  is given by -

$$F(x) = P(X \leq x)$$

$$= \int_0^x \lambda e^{-\lambda t} dt$$

$$[-e^{-\lambda t}]_0^n$$

$$= 1 - e^{-\lambda n}, \text{ for } n > 0.$$

$$= 0$$

$$\text{and hence } P(X > n) = e^{-\lambda n}.$$

## $\Rightarrow$ Properties:

1) The mean and standard deviation of the exponential distribution are equal-

2) The distribution is extremely skewed and thus there does not exist any mode-

## BETA DISTRIBUTION:

A continuous r.v.,  $x$  is said to be beta distribution with two parameters  $m$  and  $n$ , if its p.d.f is defined by:

$$f(u) = \begin{cases} \frac{1}{B(m,n)} u^{m-1} (1-u)^{n-1}, & 0 < u < 1; \\ 0; & m, n > 0 \end{cases}$$

This is called beta distribution of first kind and referred as  $B(m,n)$ .

## $\Rightarrow$ Properties of $B_1(m, n)$ :

• i: The mean and variance of this distribution are  $m$  and  $\frac{mn}{(m+n)^2(m+n+1)}$  respectively.

## Beta Distribution of Second kind:

A continuous r.v.  $X$  is said to have a beta continuous distribution of second kind with the parameters  $m$  and  $n$ ; if its pdf is defined by -

$$f(u) = \begin{cases} \frac{1}{B(m,n)} \cdot \frac{u^{m-1}}{(1-u)^{n-1}} & 0 \leq u \leq \infty \\ 0 & m, n > 0 \end{cases}$$

It is generally denoted by  $B_2(m, n)$ .

# Chi-Square Distribution:

In probability theory and statistics the chi-square distribution with  $k$ -degree of freedom is the distribution of a sum of the squares of  $k$  independent standard normal r.v. The chi-square is a special case of the gamma distribution.

g.f.s parameter is  $k$ .

The p.d.f of chi-square is -

$$f(u) = \begin{cases} \frac{1}{2^{k/2} \Gamma(k/2)} u^{k/2-1} e^{-u/2} & \text{for } 0 \leq u < \infty \\ 0 & \text{and } u < 0 \end{cases}$$

# INTRODUCTION TO STATISTICS AND Data Analysis.

## OUTLINE #1

### Statistical Interface:

obj: It is the branch of Statistics which is concerned with using probability concept to deal with uncertainty in decision making -

process of selecting and using a sample to draw inference about population from which sample is drawn-

### Samples:

A sample is a part or a subset of population. Generally it contains some of the observations but in certain situations, it may include the whole of the population. The number of observations included in a sample is called the size of the sample and is denoted by letter  $n$ .

# (1) Populations:-

A population or a statistical population is the collection or set of all possible observations whether finite or infinite, relevant to some characteristics of interest. A statistical population may be real such as the height of all college students or hypothetical such as the possible outcomes from the toss of a coin.

## Sampling Procedures:-

o:- Choosing part of a population to use to test hypotheses about the entire population.

o:- Used to choose the number of participants, interviews, or work samples to use in the assessment process.

## Collection of Data:-

The most important part of statistical work is perhaps the collection of data.

There are two types of collection of data.

o:- Primary Collections

o:- Secondary Collections

## ⇒ Collection of Primary Data:

One or more of the following methods are employed to collect primary data:-

i) Direct Personal Investigation

ii) Indirect Investigation or Personal Interviews.

iii) Collection through Questionnaires.

iv) Collection through Enumerators.

v) Collection through local sources.

vi) Computer Interviews.

## ⇒ Collection of Secondary Data:

The secondary data may be obtained from following sources:-

i) Internal secondary data-

ii) External secondary data-

• i) Official

• ii) Semi-official

## Measure of Location:

The three most common measure of location are the mean, the median and the mode.

## Mean:

The mean of a data set is found by adding all numbers in a data set and then dividing by the numbers of value in a set-

## Median:

The median is the middle value when a data set is ordered from least to greatest

## Mode:

The mode is a number that occurs most often in a data set-

# Measures of Variability:

The most common measures of variability are-

## ii: Range:

The range  $R$  is defined as the difference b/w the largest and the smallest observations in a set of data. Symbolically,

$$R = u_m - u_o$$

## iii: Interquartile Range:

The interquartile range is a measure of dispersion, defined by the difference between the third and first quartiles. Symbolically,

$$Q.D = \frac{Q_3 - Q_1}{2}$$

## iv: Mean Deviation:

The mean deviation of a set of data is defined as the arithmetic mean of the deviations measured either from the mean or from the median, all deviations being counted as positive.

$$M.D = \frac{\sum |u_i - \bar{u}|}{n} \quad \text{For sample data}$$

$$M.D = \frac{\sum |u_i - \mu|}{n} \quad \text{For population data}$$

## iv) Variance:-

The variance of the set of observations is defined as the mean of squares of deviations of all the observations from their mean. Symbolically:-

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N} \text{ for population data}$$

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n} \text{ for sample data}$$

## v) Standard Deviation:-

The positive square root of the variance is called the standard deviation. Symbolically:-

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}, \text{ for population data.}$$

$$s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}, \text{ for sample data.}$$

## DISCRETE AND CONTINUOUS DATA:-

"Discrete data" is information that can only take certain values.

"Continuous data" is data that can take any value.

i) Define Probability?

Ans- The word probability has two basic

meanings: (i) a quantitative measure of uncertainty and (ii) a measure of degree of belief in a particular statement or problem.

ii) Mutually and non-mutually exclusive events?

⇒ Mutually exclusive events-

Two events A and B of a single experiment are said to be mutually exclusive or disjoint if and only if they cannot both occur at the same time. That is they have no points in common.

⇒ Non-mutually exclusive events-

Two events A and B are said to be non-mutually exclusive if both events can occur at same time- eg; if we draw a card from an ordinary deck of 52 playing cards, it can be both a king and a diamond. Therefore,

kings and diamonds are not mutually exclusive.

## ⇒ Independent and Dependent Events:

### Independent Events:

Two events A and B in the same sample space S, are defined to be independent, if the probability that one event occurs, is not affected by whether the other event has or has not occurred that is;

$$P(A/B) = P(A)$$

and

$$P(B/A) = P(B)$$

It then follows that two events A and B are independent if and only if

$$P(A \cap B) = P(A) P(B)$$

### Dependent Events:

Two events A and B are defined to be dependent if  $P(A \cap B) \neq P(A) \times P(B)$ . This means that the occurrence of one of the events in some way affects the probability of the occurrence of the other events -

⇒ Statement of general law of Addition:  
If A and B are any two events defined in a sample space, then

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

This is the statement for general law of addition.

⇒ Definition of expectation -

## $\Rightarrow$ Properties of Normal Distribution:

The main properties of Normal distribution are given below:-

1) = The function  $f(u)$  defining the normal distribution is a proper p.d.f ie,  $f(u) \geq 0$  and the total area under the normal curve is unity.

2) = The mean and variance of the normal distribution are  $\mu$  and  $\sigma^2$  respectively.

3) = The median and the mode of the normal distribution are each equal to  $\mu$ , the mean of the distribution.

4) = The mean deviation of the normal distribution is approximately  $\frac{4}{5}$  of its standard deviation.

5) = The normal curve has points of inflection which are equidistant from the mean.

6) = For the normal distribution, the odd number movements about mean are all zero and the even order movements are  $\Rightarrow n_{2n} = (2n-1)(2n-3) \dots 5, 3, 1$