

INTRODUCTION TO DATA ANALYTICS

Class #8

Probability Distributions

Dr. Sreeja S R

Assistant Professor

**Indian Institute of Information Technology
IIIT Sri City**

TODAY'S DISCUSSION...

- Probability vs. Statistics
- Concept of random variable
- Probability distribution concept
- Discrete probability distribution
 - Discrete uniform probability distribution
 - Binomial distribution
 - Multinomial distribution
 - Hypergeometric distribution
 - Poisson distribution

TODAY'S DISCUSSION

- Continuous probability distribution
 - Continuous uniform probability distribution
 - Normal distribution
 - Standard normal distribution
 - Chi-squared distribution
 - Gamma distribution
 - Exponential distribution
 - Lognormal distribution
 - Weibull distribution

Probability and Statistics

Probability is the chance of an **outcome** in an **experiment** (also called **event**).

Event: Tossing a fair coin

Outcome: Head, Tail

Probability deals with **predicting** the likelihood of **future** events.

Statistics involves the **analysis** of the **frequency** of **past** events

Example 8.1: Consider there is a drawer containing 100 socks: 30 red, 20 blue and 50 black socks.

We can use probability to answer questions about the selection of a random sample of these socks.

- **PQ1.** What is the probability that we draw two blue socks or two red socks from the drawer?
- **PQ2.** What is the probability that we pull out three socks or have matching pair?
- **PQ3.** What is the probability that we draw five socks and they are all black?

STATISTICS

Instead, if we have no knowledge about the type of socks in the drawers, then we enter into the realm of statistics. Statistics helps us to infer properties about the population on the basis of the random sample.

Questions that would be statistical in nature are:

- **SQ1:** A random sample of 10 socks from the drawer produced one blue, four red, five black socks. **What is the total population of black, blue or red socks in the drawer?**
- **SQ2:** We randomly sample 10 socks, and write down the number of black socks and then return the socks to the drawer. The process is done for five times. The mean number of socks for each of these trial is 7. **What is the true number of black socks in the drawer?**
- etc.

PROBABILITY VS. STATISTICS

In other words:

- In probability, we are **given a model** and asked **what kind of data** we are likely to see.
- In statistics, we are **given data** and asked **what kind of model** is likely to have generated it.

Example 8.2: Measles Study

- A study on health is concerned with the **incidence of childhood measles in parents of childbearing age** in a city. For each couple, we would like to know how likely, it is that either the mother or father or both have had childhood measles.
- The current census data indicates that 20% adults between the ages 17 and 35 (regardless of sex) have had childhood measles.
 - This give us the probability that an individual in the city has had childhood measles.

DEFINING RANDOM VARIABLE

Definition: Random Variable

A random variable is a rule that assigns a numerical value to an outcome of interest.

Example: In “measles Study”, we define a random variable X as the number of parents in a married couple who have had childhood measles.

This random variable can take values of 0, 1 *and* 2.

Note:

- Random variable is not exactly the same as the variable defining a data.
- The probability that the random variable takes a given value can be computed using the rules governing probability.
- For example, the probability that $X = 1$ means either mother or father but not both has had measles is 0.32. Symbolically, it is denoted as $P(X=1) = 0.32$

PROBABILITY DISTRIBUTION

Definition: Probability distribution

A probability distribution is a definition of probabilities of the values of random variable.

•
Example 8.3: Given that 0.2 is the probability that a person (in the ages between 17 and 35) has had childhood measles. Then the probability distribution is given by

X	Probability
0	0.64
1	0.32
2	0.04

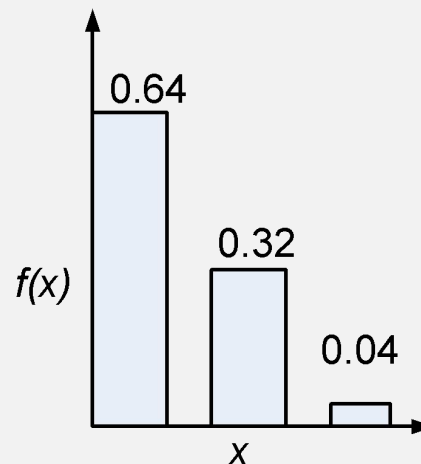


PROBABILITY DISTRIBUTION

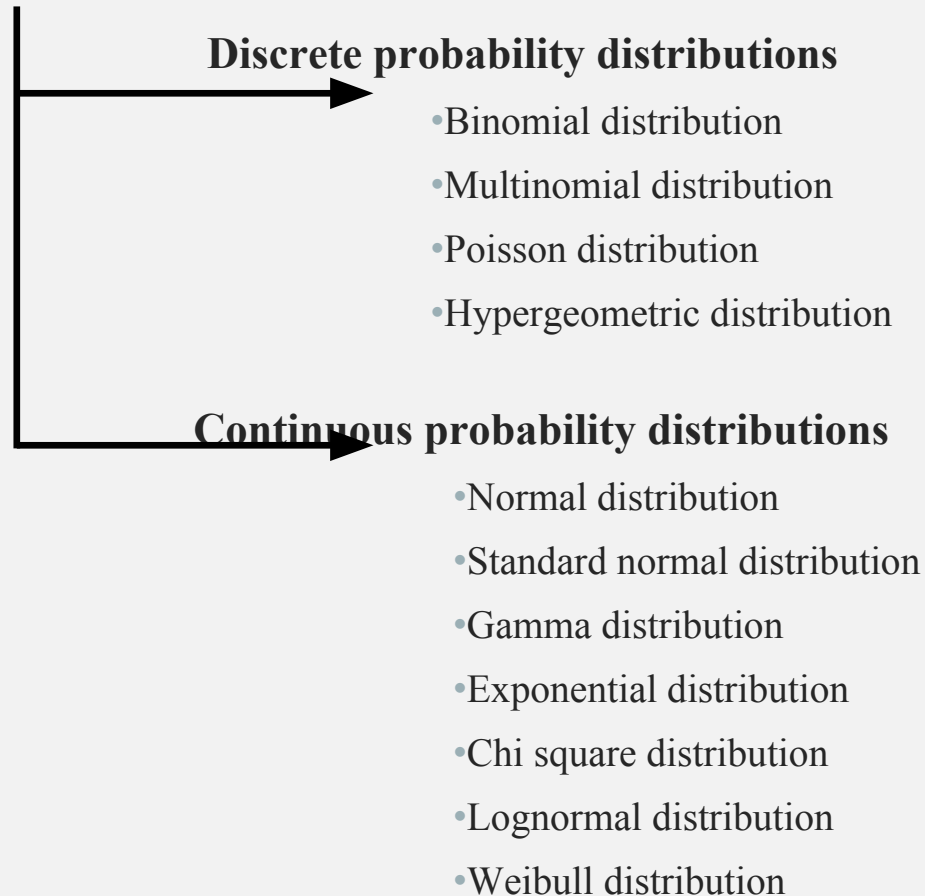
- In data analytics, the probability distribution is important with which many statistics making inferences about population can be derived .
- In general, a probability distribution function takes the following form

Example: Measles Study

	0	1	2
	0.64	0.32	0.04



TAXONOMY OF PROBABILITY DISTRIBUTIONS



USAGE OF PROBABILITY DISTRIBUTION

- Distribution (discrete/continuous) function is widely used in simulation studies.
 - A simulation study uses a computer to simulate a real phenomenon or process as closely as possible.
 - The use of simulation studies can often eliminate the need of costly experiments and is also often used to study problems where actual experimentation is impossible.

Examples 8.4:

- 1) A study involving testing the effectiveness of a new drug, the number of cured patients among all the patients who use such a drug approximately follows a binomial distribution.
- 2) Operation of ticketing system in a busy public establishment (e.g., airport), the arrival of passengers can be simulated using Poisson distribution.

DISCRETE PROBABILITY DISTRIBUTIONS

BINOMIAL DISTRIBUTION

- In many situations, an outcome has only two outcomes: **success** and **failure**.
 - Such outcome is called dichotomous outcome.
- An experiment when consists of repeated trials, each with dichotomous outcome is called **Bernoulli process**. Each trial in it is called a **Bernoulli trial**.

Example 8.5: Firing bullets to hit a target.

- Suppose, in a Bernoulli process, we define a random variable $X \equiv$ the number of successes in trials.
- Such a random variable obeys the binomial probability distribution, if the experiment satisfies the following conditions:
 - 1) The experiment consists of n trials.
 - 2) Each trial results in one of two mutually exclusive outcomes, one labelled a “*success*” and the other a “*failure*”.
 - 3) The probability of a success on a single trial is equal to p . The value of p remains constant throughout the experiment.
 - 4) The trials are independent.

DEFINING BINOMIAL DISTRIBUTION

Definition: **Binomial distribution**

The function for computing the probability for the binomial probability distribution is given by

$$f(x) = \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x}$$

for $x = 0, 1, 2, \dots, n$

Here, $f(x) = P(X = x)$, where X denotes “the number of success” and $X = x$ denotes the number of success in x trials.

BINOMIAL DISTRIBUTION

Example 8.6: Measles study

X = having had childhood measles a success

$p = 0.2$, the probability that a parent had childhood measles

$n = 2$, here a couple is an experiment and an individual a trial, and the number of trials is two.

Thus,

$$P(x = 0) = \frac{2!}{0!(2-0)!} (0.2)^0 (0.8)^{2-0} = \mathbf{0.64}$$

$$P(x = 1) = \frac{2!}{1!(2-1)!} (0.2)^1 (0.8)^{2-1} = \mathbf{0.32}$$

$$P(x = 2) = \frac{2!}{2!(2-2)!} (0.2)^2 (0.8)^{2-2} = \mathbf{0.04}$$

BINOMIAL DISTRIBUTION

Example 8.7: Verify with real-life experiment

Suppose, 10 pairs of random numbers are generated by a computer (Monte-Carlo method)

15 38 68 39 49 54 19 79 38 14

If the value of the digit is 0 or 1, the outcome is “had childhood measles”, otherwise, (digits 2 to 9), the outcome is “did not”.

For example, in the first pair (i.e., 15), representing a couple and for this couple, $x = 1$. The frequency distribution, for this sample is

x	0	1	2
$f(x)=P(X=x)$	0.7	0.3	0.0

Note: This has close similarity with binomial probability distribution!

THE MULTINOMIAL DISTRIBUTION

The binomial experiment becomes a multinomial experiment, if we let each trial has more than two possible outcome.

Definition: **Multinomial distribution**

If a given trial can result in the k outcomes E_1, E_2, \dots, E_k with probabilities p_1, p_2, \dots, p_k , then the probability distribution of the random variables X_1, X_2, \dots, X_k representing the number of occurrences for E_1, E_2, \dots, E_k in n independent trials is

$$f(x_1, x_2, \dots, x_k) = \binom{n}{x_1, x_2, \dots, x_k} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}$$

$$\text{where } \binom{n}{x_1, x_2, \dots, x_k} = \frac{n!}{x_1! x_2! \dots x_k!}$$

$$\sum_{i=1}^k x_i = n \text{ and } \sum_{i=1}^k p_i = 1$$

THE HYPERGEOMETRIC DISTRIBUTION

- Collection of samples with two strategies
 - With replacement
 - Without replacement
- A necessary condition of the binomial distribution is that all trials are **independent to each other**.
- When sample is collected “with replacement”, then each trial in sample collection is independent.

Example 8.8:

Probability of observing three red cards in 5 draws from an ordinary deck of 52 playing cards.

- You draw one card, note the result and then returned to the deck of cards
- Reshuffled the deck well before the next drawing is made
- The hypergeometric distribution *does not require independence* and is based on the sampling done **without replacement**.

THE HYPERGEOMETRIC DISTRIBUTION

- In general, the hypergeometric probability distribution enables us to find the probability of selecting x successes in n trials from N items.

Properties of Hypergeometric Distribution

- A random sample of size n is selected without replacement from N items.
- k of the N items may be classified as success and $N - k$ items are classified as failure.

Let X denotes a hypergeometric random variable defining the number of successes.

Definition: Hypergeometric Probability Distribution

The probability distribution of the hypergeometric random variable X , the number of successes in a random sample of size n selected from N items of which k are labelled success and $N - k$ labelled as failure is given by

$$f(x) = P(X = x) = \frac{\binom{k}{x} \binom{N-k}{n-x}}{\binom{N}{n}}$$
$$\max(0, n - (N - k)) \leq x \leq \min(n, k)$$

MULTIVARIATE HYPERGEOMETRIC DISTRIBUTION

The hypergeometric distribution can be extended to treat the case where the N items can be divided into k classes A_1, A_2, \dots, A_k with a_1 elements in the first class A_1 , ... and a_k elements in the k^{th} class. We are now interested in the probability that a random sample of size n yields x_1 elements from A_1 , x_2 elements from A_2 ,, x_k elements from A_k .

Definition: Multivariate Hypergeometric Distribution

If N items are partitioned into k classes a_1, a_2, \dots, a_k respectively, then the probability distribution of the random variables X_1, X_2, \dots, X_k , representing the number of elements selected from A_1, A_2, \dots, A_k in a random sample of size n , is

$$f(x_1, x_2, \dots, x_k) = P(X_1 = x_1, X_2 = x_2, \dots, X_k = x_k) = \frac{\binom{a_1}{x_1} \binom{a_2}{x_2} \dots \binom{a_k}{x_k}}{\binom{N}{n}}$$

with $\sum_{i=1}^k x_i = n$ and $\sum_{i=1}^k a_i = N$

THE POISSON DISTRIBUTION

There are some experiments, which involve the occurring of the number of outcomes during a given time interval (or in a region of space).

Such a process is called **Poisson process**.

Example 8.9:

Number of clients visiting a ticket selling counter in a metro station.



THE POISSON DISTRIBUTION

Properties of Poisson process

- The number of outcomes in one time interval is independent of the number that occurs in any other disjoint interval [[Poisson process has no memory](#)]
- The probability that a single outcome will occur during a very short interval is proportional to the length of the time interval and does not depend on the number of outcomes occurring outside this time interval.
- The probability that more than one outcome will occur in such a short time interval is negligible.

Definition: **Poisson distribution**

The probability distribution of the Poisson random variable X , representing the number of outcomes occurring in a given time interval t , is

$$f(x, \lambda t) = P(X = x) = \frac{e^{-\lambda t} \cdot (\lambda t)^x}{x!}, x = 0, 1, \dots$$

where λ is the average number of outcomes per unit time and $e = 2.71828 \dots$

DESCRIPTIVE MEASURES

Given a random variable X in an experiment, we have denoted $f(x) = P(X = x)$, the probability that $X = x$. For discrete events $f(x) = 0$ for all values of x except $x = 0, 1, 2, \dots$.

Properties of discrete probability distribution

1. $0 \leq f(x) \leq 1$
2. $\sum f(x) = 1$
3. $\mu = \sum x \cdot f(x)$ [is the **mean**]
4. $\sigma^2 = \sum (x - \mu)^2 \cdot f(x)$ [is the **variance**]

In 2, 3 and 4, summation is extended for all possible discrete values of x .

Note: For discrete **uniform** distribution, $f(x) = \frac{1}{n}$ with $x = 1, 2, \dots, n$

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

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

DESCRIPTIVE MEASURES

1. Binomial distribution

The binomial probability distribution is characterized with p (the probability of success) and n (is the number of trials). Then

$$\mu = n.p$$

$$\sigma^2 = np(1 - p)$$

2. Hypergeometric distribution

The hypergeometric distribution function is characterized with the size of a sample (n), the number of items (N) and k labelled success. Then

$$\mu = \frac{nk}{N}$$

$$\sigma^2 = \frac{N - n}{N - 1} \cdot n \cdot \frac{k}{N} \left(1 - \frac{k}{N}\right)$$

DESCRIPTIVE MEASURES

3. Poisson Distribution

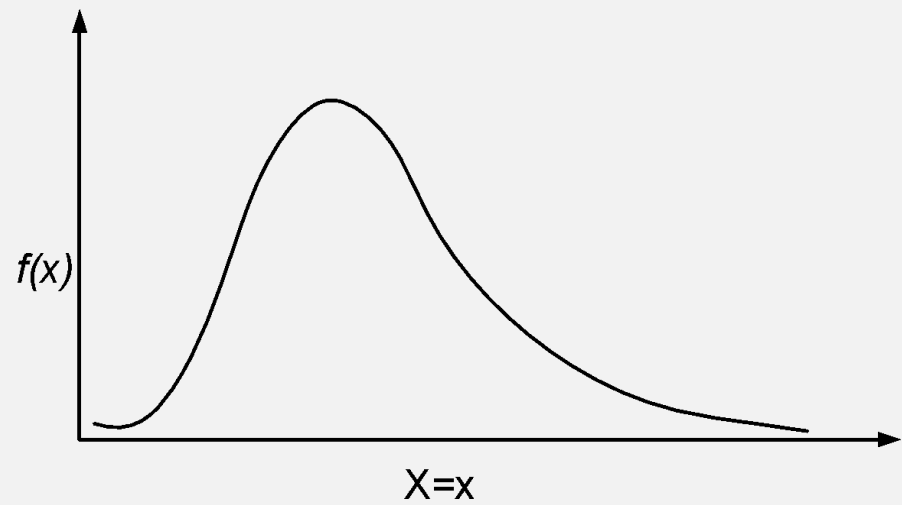
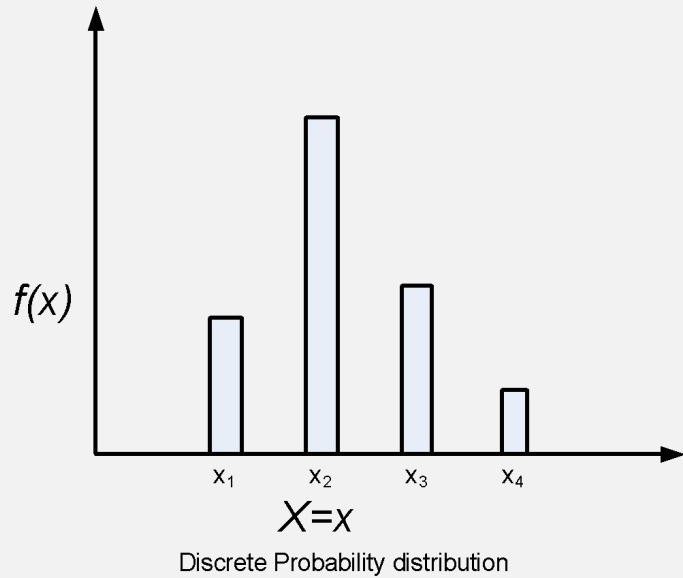
The Poisson distribution is characterized with λt where $\lambda = \text{the mean of outcomes}$ and $t = \text{time interval}$.

$$\mu = \lambda t$$

$$\sigma^2 = \lambda t$$

CONTINUOUS PROBABILITY DISTRIBUTIONS

CONTINUOUS PROBABILITY DISTRIBUTIONS



Continuous Probability Distribution

CONTINUOUS PROBABILITY DISTRIBUTIONS

- When the random variable of interest can take **any value in an interval**, it is called continuous random variable.
 - Every continuous random variable has **an infinite, uncountable number of possible values** (i.e., any value in an interval)
- Consequently, continuous random variable differs from discrete random variable.

PROPERTIES OF PROBABILITY DENSITY FUNCTION

The function $f(x)$ is a probability density function for the continuous random variable X , defined over the set of real numbers R , if

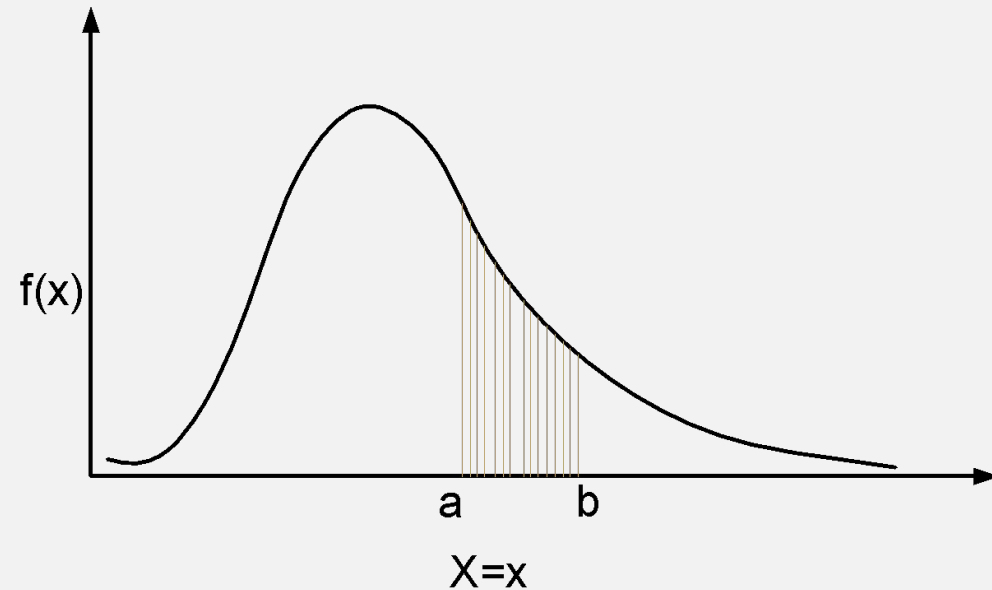
1. $f(x) \geq 0$, for all $x \in R$

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

3. $P(a \leq X \leq b) = \int_a^b f(x) dx$

4. $\mu = \int_{-\infty}^{\infty} xf(x) dx$

5. $\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$



CONTINUOUS UNIFORM DISTRIBUTION

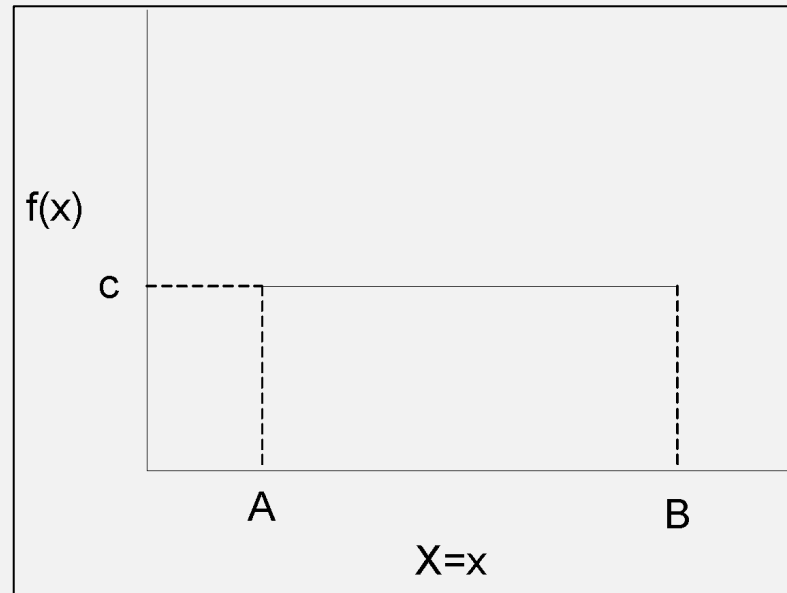
- One of the simplest continuous distribution in all of statistics is the continuous **uniform** distribution.

Definition: Continuous Uniform Distribution

The density function of the continuous uniform random variable X on the interval $[A, B]$ is:

$$f(x; A, B) = \begin{cases} \frac{1}{B - A} & A \leq x \leq B \\ 0 & \text{Otherwise} \end{cases}$$

CONTINUOUS UNIFORM DISTRIBUTION



Note:

a) $\int_{-\infty}^{\infty} f(x)dx = \frac{1}{B-A} \times (B - A) = 1$

b) $P(c < x < d) = \frac{d-c}{B-A}$ where both c and d are in the interval (A, B)

c) $\mu = \frac{A+B}{2}$

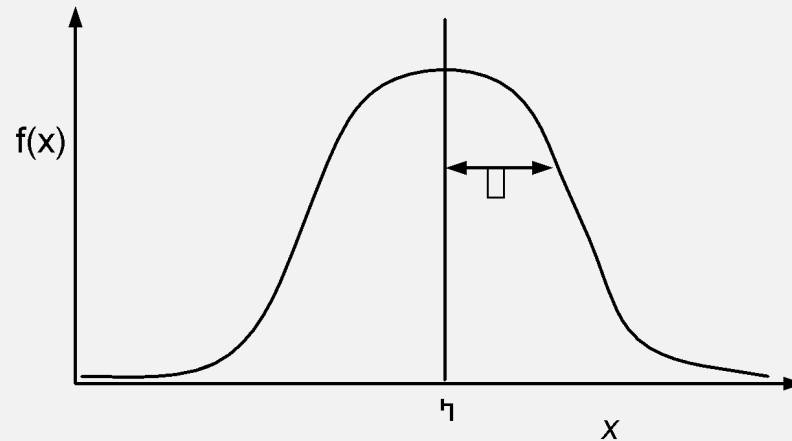
d) $\sigma^2 = \frac{(B-A)^2}{12}$

NORMAL DISTRIBUTION

- The most often used continuous probability distribution is the normal distribution; it is also known as **Gaussian distribution**.
- Its graph called the normal curve is the bell-shaped curve.
- Such a curve approximately describes many phenomenon occur in nature, industry and research.
 - Physical measurement in areas such as meteorological experiments, rainfall studies and measurement of manufacturing parts are often more than adequately explained with normal distribution.
- A continuous random variable X having the bell-shaped distribution is called a normal random variable.

Normal Distribution

- The mathematical equation for the probability distribution of the normal variable depends upon the two parameters μ and σ , its mean and standard deviation.



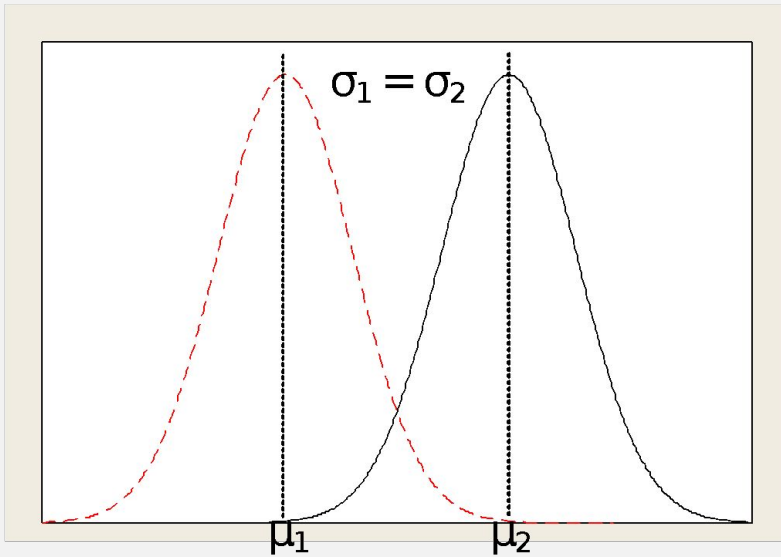
Definition: Normal distribution

The density of the normal variable x with mean μ and variance σ^2 is

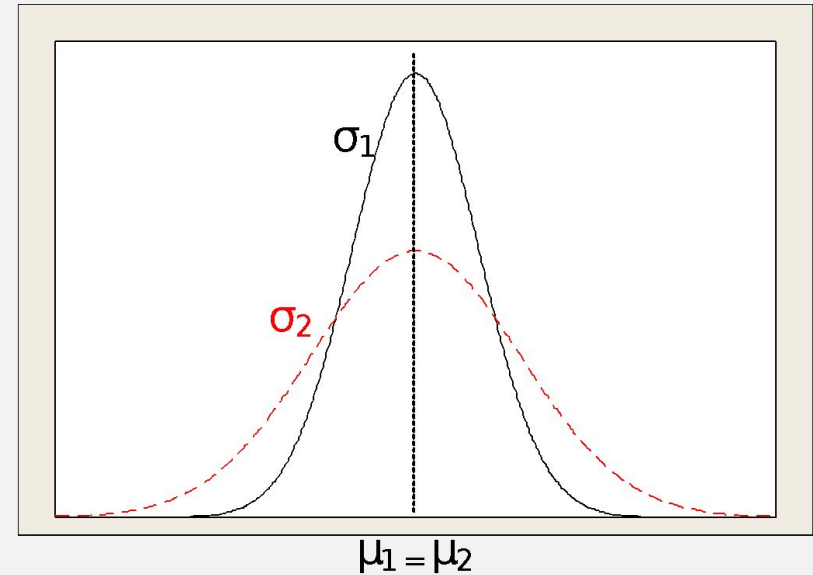
$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad -\infty < x < \infty$$

where $\pi = 3.14159 \dots$ and $e = 2.71828 \dots$, the Naperian constant

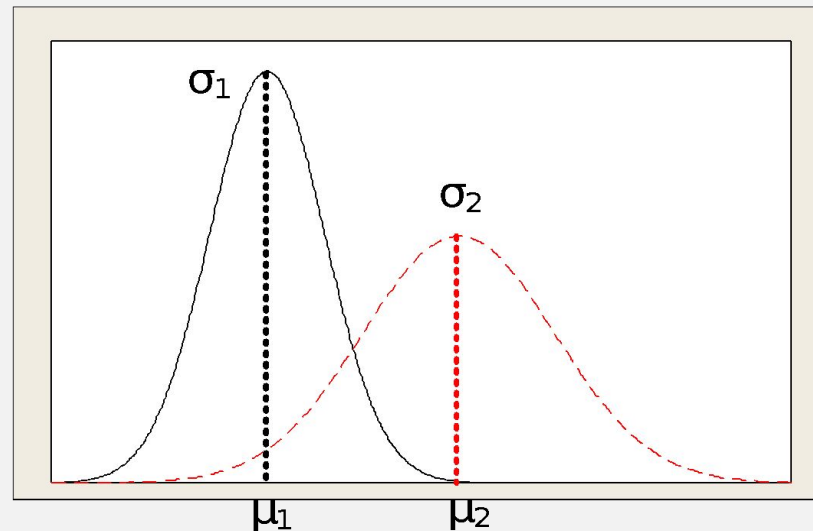
Normal Distribution



Normal curves with $\mu_1 < \mu_2$ and $\sigma_1 = \sigma_2$



Normal curves with $\mu_1 = \mu_2$ and $\sigma_1 < \sigma_2$



Normal curves with $\mu_1 < \mu_2$ and $\sigma_1 < \sigma_2$

PROPERTIES OF NORMAL DISTRIBUTION

- The curve is symmetric about a vertical axis through the mean μ .
- The random variable x can take any value from $-\infty$ to ∞ .
- The most frequently used descriptive parameters define the curve itself.
- The mode, which is the point on the horizontal axis where the curve is a maximum occurs at $x = \mu$.
- The total area under the curve and above the horizontal axis is equal to 1.

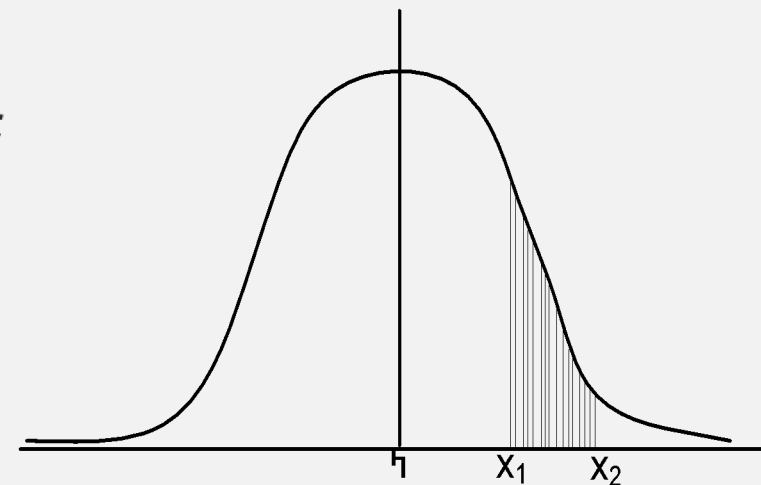
$$\int_{-\infty}^{\infty} f(x) dx = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx = 1$$

- $\mu = \int_{-\infty}^{\infty} x \cdot f(x) dx = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} x \cdot e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx$

- $\sigma^2 = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} (x - \mu)^2 \cdot e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx$

- $P(x_1 < x < x_2) = \frac{1}{\sigma\sqrt{2\pi}} \int_{x_1}^{x_2} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx$

denotes the probability of x in the interval (x_1, x_2) .



STANDARD NORMAL DISTRIBUTION

- The normal distribution has computational complexity to calculate $P(x_1 < x < x_2)$ for any two (x_1, x_2) and given μ and σ
- To avoid this difficulty, the concept of z-transformation is followed.



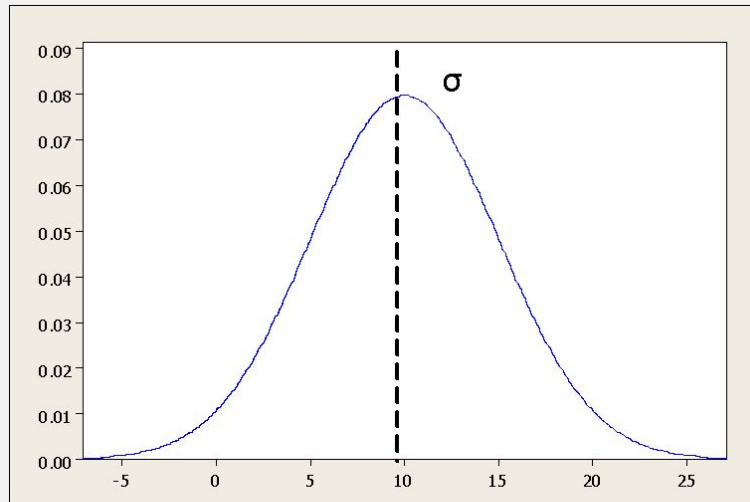
- X: Normal distribution with mean μ and variance σ^2 .
- Z: Standard normal distribution with mean $\mu = 0$ and variance $\sigma^2 = 1$.
- Therefore, if $f(x)$ assumes a value, then the corresponding value of $f(z)$ is given by

$$\begin{aligned} f(x; \mu, \sigma) : P(x_1 < x < x_2) &= \frac{1}{\sigma \sqrt{2\pi}} \int_{x_1}^{x_2} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx \\ &= \frac{1}{\sigma \sqrt{2\pi}} \int_{z_1}^{z_2} e^{-\frac{1}{2}z^2} dz = f(z; 0, \sigma) \end{aligned}$$

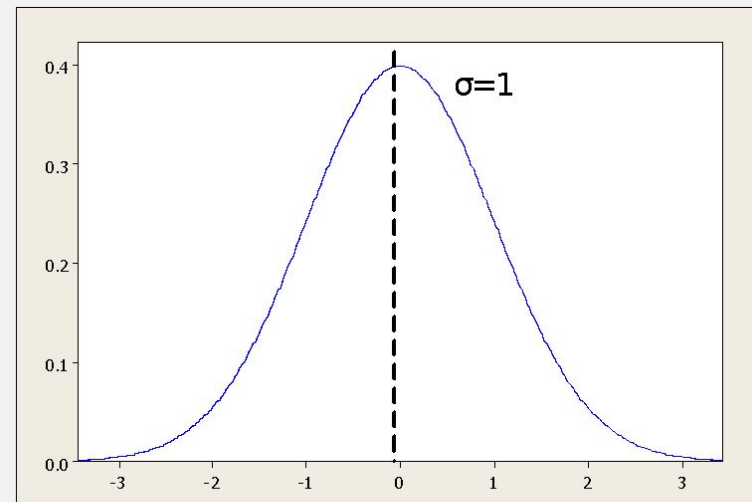
STANDARD NORMAL DISTRIBUTION

Definition: Standard normal distribution

The distribution of a normal random variable with mean 0 and variance 1 is called a standard normal distribution.



$$x = \mu$$
$$f(x; \mu, \sigma)$$



$$\mu = 0$$
$$f(z; 0, 1)$$

GAMMA DISTRIBUTION

The gamma distribution derives its name from the well known gamma function in mathematics.

Definition: Gamma Function

$$\Gamma(\alpha) = \int_0^{\alpha} x^{\alpha-1} e^{-x} dx \quad \text{for } \alpha > 0$$

Integrating by parts, we can write,

$$\begin{aligned} \Gamma(\alpha) &= (\alpha - 1) \int_0^{\alpha} x^{\alpha-2} e^{-x} dx \\ &= (\alpha - 1) \Gamma(\alpha - 1) \end{aligned}$$

Thus Γ function is defined as a recursive function.

GAMMA DISTRIBUTION

When $\alpha = n$, we can write,

$$\begin{aligned}\Gamma(n) &= (n-1)(n-2) \dots \dots \dots \Gamma(1) \\ &= (n-1)(n-2) \dots \dots \dots 3.2.1 \\ &= (n-1)!\end{aligned}$$

$$\text{Further, } \Gamma(1) = \int_0^{\infty} e^{-x} dx = 1$$

Note:

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$$

[An important property]

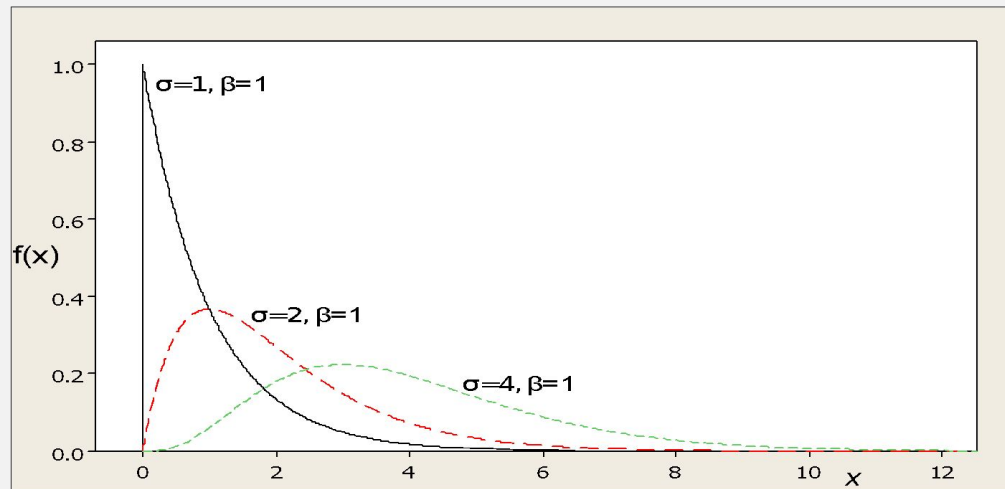
GAMMA DISTRIBUTION

Definition: Gamma Distribution

The continuous random variable x has a gamma distribution with parameters α and β such that:

$$f(x; \alpha, \beta) = \begin{cases} \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-\frac{x}{\beta}} & x > 0 \\ 0 & \text{Otherwise} \end{cases}$$

where $\alpha > 0$ and $\beta > 0$



EXPONENTIAL DISTRIBUTION

Definition: Exponential Distribution

The continuous random variable x has an exponential distribution with parameter β , where:

$$f(x; \beta) = \begin{cases} \frac{1}{\beta} e^{-\frac{x}{\beta}} & \text{where } \beta > 0 \\ 0 & \end{cases}$$

Note:

- 1) The mean and variance of gamma distribution are

$$\begin{aligned} \mu &= \alpha\beta \\ \sigma^2 &= \alpha\beta^2 \end{aligned}$$

- 2) The mean and variance of exponential distribution are

$$\begin{aligned} \mu &= \beta \\ \sigma^2 &= \beta^2 \end{aligned}$$

CHI-SQUARED DISTRIBUTION

Definition: Chi-squared distribution

The continuous random variable x has a Chi-squared distribution with ν degrees of freedom, is given by

$$f(x; \nu) = \begin{cases} \frac{1}{2^{\frac{\nu}{2}} \Gamma(\nu/2)} x^{\nu/2-1} e^{-\frac{x}{2}}, & x > 0 \\ 0 & \text{Otherwise} \end{cases}$$

where ν is a positive integer.

- The Chi-squared distribution plays an important role in statistical inference .
- The mean and variance of Chi-squared distribution are:

$$\mu = \nu \text{ and } \sigma^2 = 2\nu$$

LOGNORMAL DISTRIBUTION

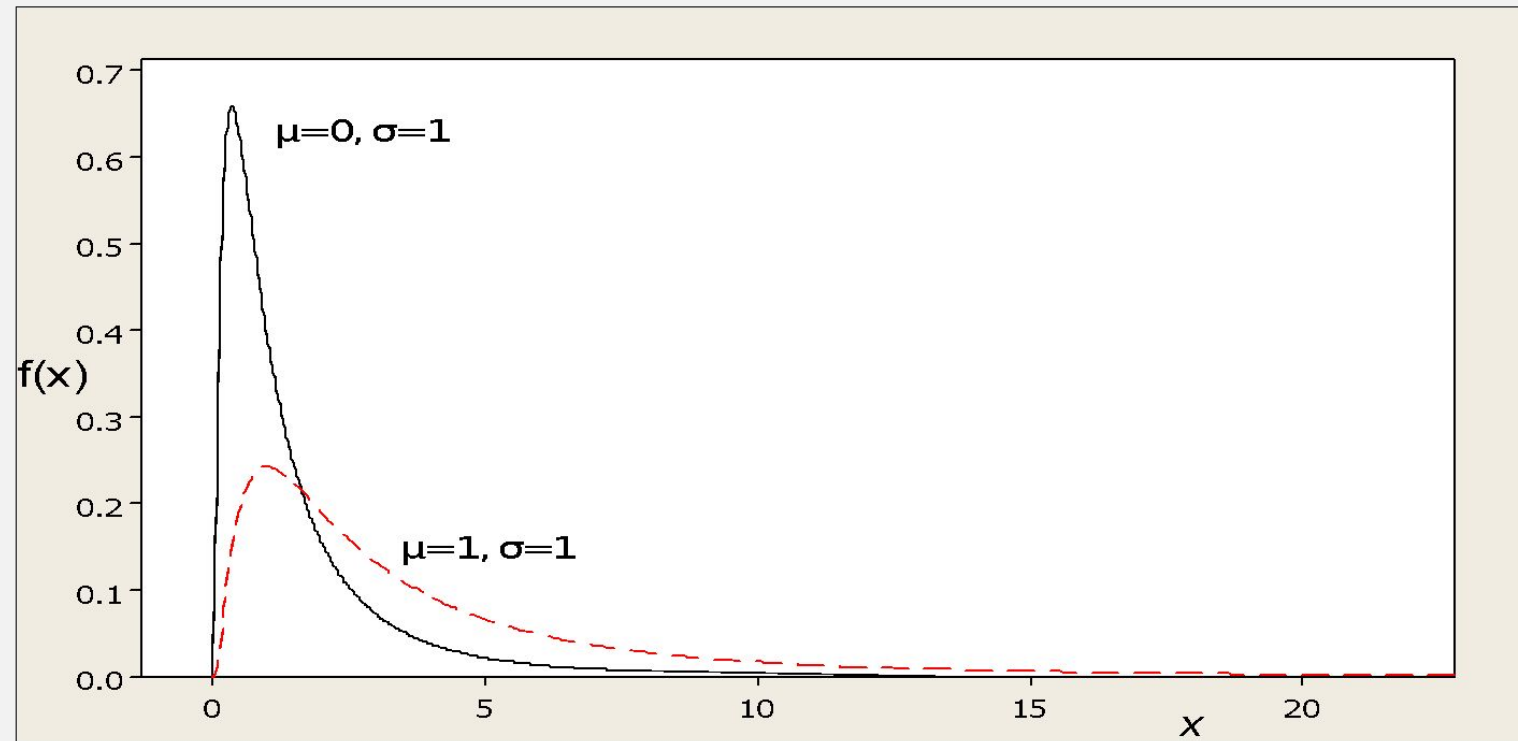
The lognormal distribution applies in cases where a natural log transformation results in a normal distribution.

Definition: Lognormal distribution

The continuous random variable x has a lognormal distribution if the random variable $y = \ln(x)$ has a normal distribution with mean μ and standard deviation σ . The resulting density function of x is:

$$f(x; \mu, \sigma) = \begin{cases} \frac{1}{\sigma x \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2} [\ln(x) - \mu]^2} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

LOGNORMAL DISTRIBUTION



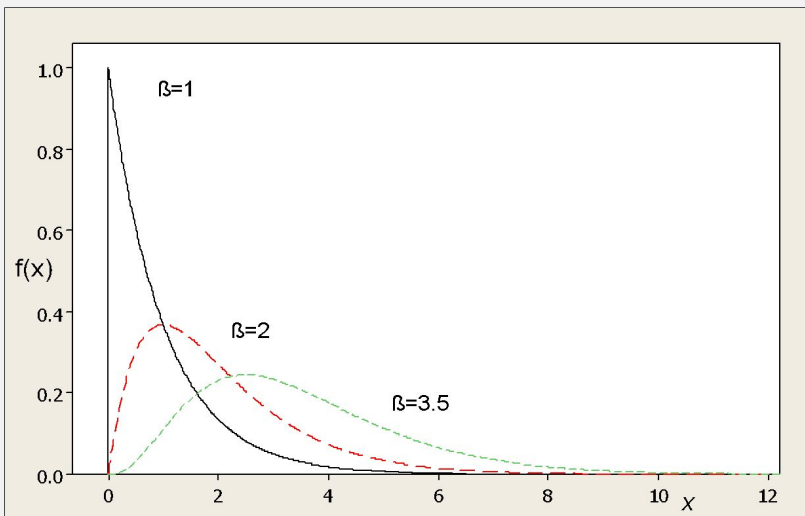
WEIBULL DISTRIBUTION

Definition: Weibull Distribution

The continuous random variable x has a Weibull distribution with parameter α and β such that.

$$f(x; \alpha, \beta) = \begin{cases} \alpha \beta x^{\beta-1} e^{-\alpha x^\beta} & x > 0 \\ 0 & \text{Otherwise} \end{cases}$$

where $\alpha > 0$ and $\beta > 0$



The mean and variance of Weibull distribution are:

$$\mu = \alpha^{-1/\beta} \Gamma\left(1 + \frac{1}{\beta}\right)$$

$$\sigma^2 = \alpha^{-2/\beta} \left\{ \Gamma\left(1 + \frac{2}{\beta}\right) - [\Gamma(1 + \frac{1}{\beta})]^2 \right\}$$

REFERENCE

- The detail material related to this lecture can be found in

Probability and Statistics for Engineers and Scientists (8th Ed.) by
Ronald E. Walpole, Sharon L. Myers, Keying Ye (Pearson), 2013.

Any question?

QUESTIONS OF THE DAY...

1. Give some examples of random variables? Also, tell the range of values and whether they are with continuous or discrete values.
2. In the following cases, what are the probability distributions are likely to be followed. In each case, you should mention the random variable and the parameter(s) influencing the probability distribution function.
 - a) In a retail source, how many counters should be opened at a given time period.
 - b) Number of people who are suffering from cancers in a town?

QUESTIONS OF THE DAY...

2. In the following cases, what are the probability distributions are likely to be followed. In each case, you should mention the random variable and the parameter(s) influencing the probability distribution function.
- c) A missile will hit the enemy's aircraft.
 - d) A student in the class will secure EX grade.
 - e) Salary of a person in an enterprise.
 - f) Accident made by cars in a city.
 - g) People quit education after i) primary ii) secondary and iii) higher secondary educations.

QUESTIONS OF THE DAY...

3. How you can calculate the mean and standard deviation of a population if the population follows the following probability distribution functions with respect to an event.
- a) Binomial distribution function.
 - b) Poisson's distribution function.
 - c) Hypergeometric distribution function.
 - d) Normal distribution function.
 - e) Standard normal distribution function.