

Motivating examples:

- ① who will win the next general election in India
- ② Investigate people's belief:
 - (a) Do you believe in life after death?
 - (b) Would you be willing to pay higher prices to protect environment?
 - (c) How much TV do you watch per day?

Inferential statistics:- The word infer means to arrive at a decision or prediction by reasoning from known evidence. Statistical inference does this using data as evidence.

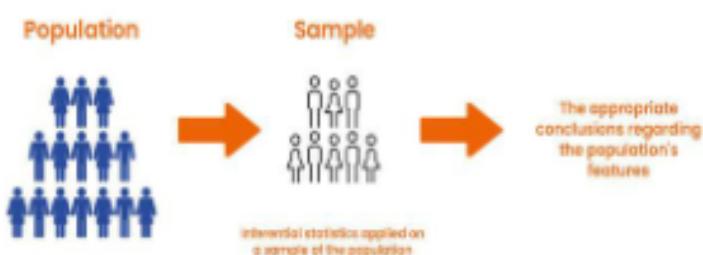
It is the process of generating conclusion about a population from sample(s) from the population.

Sample versus population:- population is a total set of similar items or events which is of interest for some statistical question or experiment.

Example: All voter in India

A sample is a subset of the population for whom we have (or plan to have) data, often randomly selected.

INFERENTIAL STATISTICS



How does inferential statistics answer a statistical question?

Let us consider a question

How likely is that 50% or more population think that Indian economy is getting worse.

Population: All voters in India, All graduates in economics in India etc.

To answer the above question, let us take a sample of 250 people. One can then take sample mean to arrive to a conclusion.

- ① Can we directly use this sample mean as the population mean?
- ② Will we be able to make 100% correct estimate or we have to go with approximate value?

We can work this out in two ways

- ① We draw whole lot of samples of size 250. Take the average of the sample mean of the each sample.

Challenge: Since the population involves people, taking lots of samples may be difficult and costly. There is an element of uncertainty as how well the sample represents the population. The way the sample is taken matters.

- ② Other way to work out this problem is to use probability theory.

The second approach is widely popular. For this we need to understand the fundamentals of probability and random variable.

Probability :- With any random phenomenon, the probability of a particular outcome is the proportion of times that outcome would occur in long run of observations.

Example: A weather forecaster might say that probability of rain today is 0.70. This means that in a large number of days with atmospheric condition like those today, the proportion of days in which rain occurs is 0.70.

Sometimes probabilities are expressed as percentage, such as the weather forecaster might say that the probability of rain is 70%.

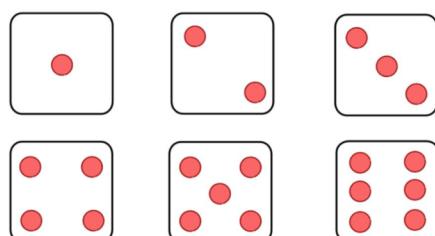
Probability Terminology :- Let us consider an example of rolling a die.



Sample space :- Sample space of an experiment consists of the set of all possible outcomes.

In the case of rolling a die, the sample space is $S = \{1, 2, 3, 4, 5, 6\}$ - denoted by S.

Sample Space for Rolling a Die:



6 outcomes

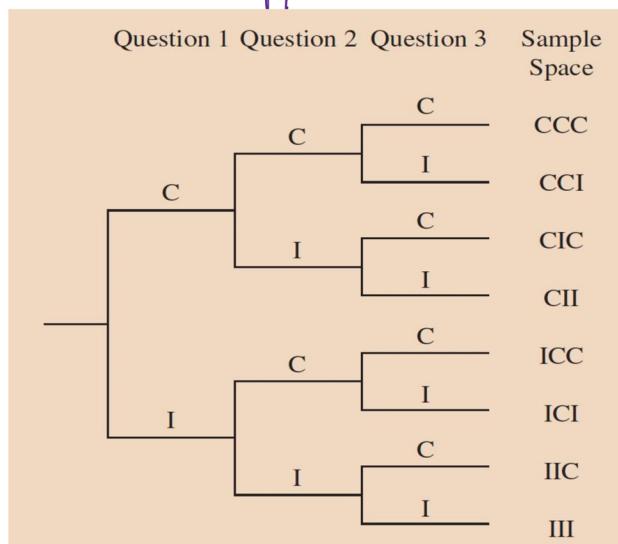
The integers 1 - - 6 represent the number of dots on the six faces of the die.

Real-life example for sample space:- Let a statistics instructor gives a pop quiz with three multiple choice questions.

The student's answer is either correct (C) or incorrect (I). For example, if a student answered the first two questions correctly & the last question incorrectly, the student's outcome on the quiz can be symbolized by CCI.

What is the sample space of the pop quiz?

Solution:- One technique for listing the possible outcomes for finding the sample space is to draw a tree-diagram, with branches showing what can happen on subsequent trials.



From the diagram, a student's performance has eight possible outcomes. Thus, the sample space is

$$\{ CCC, CCI, CIC, CII, ICC, ICI, IIC, III \}$$

Sample points:- The six possible outcomes of the rolling a die are the sample points of the experiment.

Events:- An event is a subset of 'S' and may consider any number of sample points.

For the pop quiz example, an event may be

$A = \text{student answer all three question correctly} = \{\text{ccc}\}$

$B = \text{student passes (at least two questions are correct)} = \{\text{CCI, CIC, ICC, CCC}\}$

Complement of an event:-

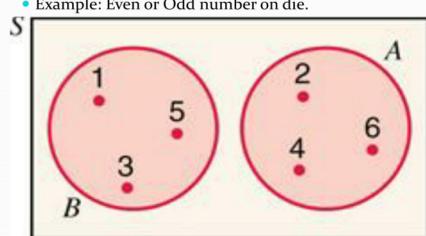
The complement of an event A, denoted by \bar{A} , consists of all the sample points in 'S' that are not in A. If $A = \{2, 4\} \Rightarrow \bar{A} = \{1, 3, 5, 6\}$

Mutually Exclusive events:- Two events are said

to be mutually exclusive if they have no sample point in common.

Mutually Exclusive Events for rolling a die example.

Example;



The above two events are mutually exclusive or disjoint events.

Union of events:- The union (sum) of two events is an event that consists all the sample points in the two events. Let $A = \{1, 2\}$, $B = \{5, 6\}$

$$A \cup B = \{1, 2, 5, 6\}$$

The above implies $A \cup \bar{A} = S$ (sample space)

Intersection of events:- Intersection two events is an event that consists of the points that are common to the two events.

If $A = \{2, 4\}$ and $B = \{1, 2, 3\} \Rightarrow A \cap B = \{2\}$

For mutually exclusive events (see fig)

$$A \cap B = \emptyset \text{ (Null event)}$$

Thus, $A \cap \bar{A} = \emptyset$

Probability rules:— Associated with each event A in S has probability of occurrence $P(A)$, which is defined as

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

① The probability of an event A satisfies the condition

$$P(A) \geq 0$$

② The probability of the sample space 'S'

$$P(S) = 1$$

① and ② imply that

$$0 \leq P(A) \leq 1$$

③ Let A_i for $i = 1, 2, \dots, n$ are mutually exclusive events in 'S', that is

$$A_i \cap A_j = \emptyset, i \neq j = 1, 2, \dots, n$$

Then, the probability of union of these mutually exclusive events satisfies the condition

$$P(\bigcup_i A_i) = \sum_i P(A_i)$$

Addition rule
of probability

If two events A and B are not disjoint or mutually exclusive

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

Example: Consider a family with two children. Let F and M denote child gender for female and male respectively

(a) find the sample-space

- (b) find the prob of the event A = first child is a girl
- (c) find the prob of the event B = second child is a girl
- (d) find the prob of the event $A \cup B$

Solution:- (a) Sample Space $S = \{FF, FM, MF, MM\}$

$$(b) A = \{FF, FM\} \Rightarrow P(A) = 2/4 = 0.5$$

$$(c) B = \{FF, MF\} \Rightarrow P(B) = 2/4 = 0.5$$

$$(d) A \cup B = \{FM, MF, FF\} \Rightarrow P(A \cup B) = 3/4 = 0.75$$

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

$$A \cap B = \{FF\} \Rightarrow P(A \cap B) = 1/4 = 0.25$$

$$P(A \cup B) = 0.5 + 0.5 - 0.25 = \underline{\underline{0.75}}$$

Random variable:- consider an experiment with a sample space 'S' and sample points $s \in S$. Let us define a function $x(s)$, which maps the outcome $s \in S$ of the experiment on the real line $(-\infty, \infty)$. The function $x(s)$ is called a random variable. In other words, A random variable is a numerical measurement of the outcome of a random phenomenon.

Example ① flipping a coin

the possible outcomes are Head (H) and Tail (T)

sample space $S = \{H, T\}$

let us define a function $x(s)$, $s \in S$, such that

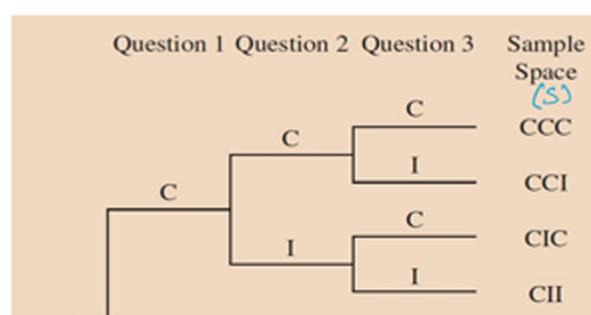
$$x(s) = \begin{cases} 1 & s = H \\ -1 & s = T \end{cases}$$

$x(s)$ is a random variable that can take value ± 1

Types of random variable 

Discrete random variable:- If a random variable $x(s)$ [Typically denoted by x] takes finite many values, it is known as discrete random variable.

Example :- Recall the top quiz example having sample space given below.



$$S = \{ccc, cci, cic, cII, Icc, ICI, IIc, III\}$$

$X = \{7, 6, 5, 4, 3, 2, 1, 0\} \rightarrow$ discrete random variable.

Example 2: X corresponding to tossing a coin and rolling a die.

Probability distribution of discrete random variable

Recall that a discrete RV takes finite many distinct values. Its probability distribution assigns a probability to each possible value, and the sum of probabilities of all possible values equals 1.

It is denoted by $p_X(x)$, which gives the probability that the random variable X takes a value x .

Example:- plot the probability distribution of the RV X that count the number of correct questions in the pop quiz exam. Assume each outcome to be equally likely.

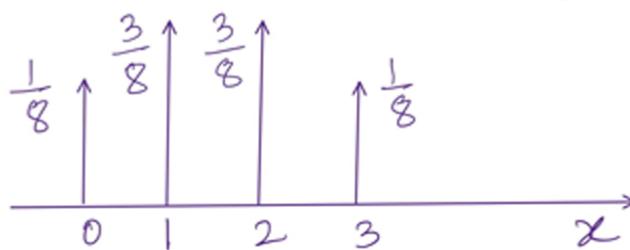
Sample Space (S) | $X = \# \text{ of correct questions}$

ccc	3	The values taken by the random variable $X = 0, 1, 2, 3$
cci	2	
cic	2	
cII	1	
Icc	2	
ICI	1	
IIC	1	
III	0	

$p_X(0) = P(X=0) = \frac{1}{8} = 0.125$
 $p_X(1) = P(X=1) = \frac{3}{8} = 0.375$
 $p_X(2) = P(X=2) = \frac{3}{8} = 0.375$
 $p_X(3) = P(X=3) = \frac{1}{8} = 0.125$

Note that $p_X(0) + p_X(1) + p_X(2) + p_X(3) = \frac{1}{8} + \frac{3}{8} + \frac{3}{8} + \frac{1}{8} = 1$

Plot of prob distribution



once you have prob distribution, you can answer the questions like

- ① Prob that a student passes (at least two questions are correct)

$$= P(X \geq 2) = P_X(2) + P_X(3) = \frac{3}{8} + \frac{1}{8} = \frac{4}{8}$$

Continuous Random Variable: Many experiments have continuous outcomes. For example,

- ① Time that people take to commute to work in Bangalore
② Salary of an engineer with one year of experience

The outcome in this can take any value in a range.

Thus, the sample space 's' is continuous so is the mapping $X(s) = s$.

The random variable X in this case can take any value in a certain range. If a random variable can take any value in a range, it is known as continuous random variable.

Probability distribution of continuous RV:-

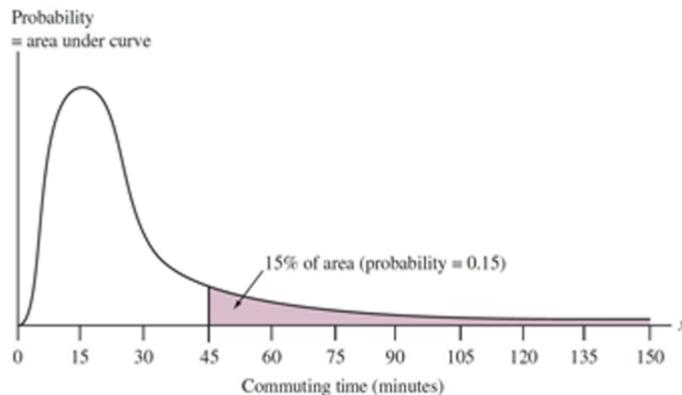
The probability distribution of a continuous random variable describes how probabilities are distributed over the range of possible values the variable can take.

It is represented by a probability density function (pdf), $f_X(x)$, which satisfies two key properties.

- ① $f_X(x) \geq 0$ for all x , ensuring probabilities are non-negative
- ② The total area under the curve of $f_X(x)$ over the variable's range equals 1.

Fig shows the curve for the probability distribution X = Commuting time to work in the united state.

Fig shows the curve for the probability distribution $x = \text{commuting time to work}$ in the united state.



The area of the shaded portion gives the prob that commuti time is more than 45 minutes. Thus

$$P(X > 45 \text{ minutes}) = \int_{45}^{150} f_X(x) dx = 0.15$$

statistical averages. of random variable:

Mean or expected value of a random variable X is

$$E[X] = \int_{-\infty}^{\infty} x f_X(x) dx \text{ For continuous random variable (CRV)}$$

$$E[X] = \sum_i x_i p_X(x_i) \text{ For discrete random variable (DRV)}$$

It provides the long-term average of the variable.

$E[X] = \mu$ also known as the first moment of X .

The variance of a rv X

$$\sigma_x^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f_X(x) dx \text{ For CRV}$$

$$\sigma_x^2 = \sum_i (x_i - \mu)^2 p_X(x_i) \text{ For DRV}$$

Variance σ_x^2 provides a measure of the dispersion of the random variable. In other words, it measures how much the random variable X varies from its mean.

Higher the value of σ_x^2 , larger is the difference between values taken by X and its mean value.

standard deviation $\sigma = \sqrt{\sigma_x^2} = \sigma_x$

with the help of σ , we can find dispersion of values of X relative to its mean.

Cumulative Distribution Function (CDF):

The CDF, denoted by $F_X(x)$, gives the probability that the random variable X takes value less than or equal to x .

$$F_X(x) = P(X \leq x) = \int_{-\infty}^x f_X(t) dt \text{ for CRV}$$

$$F_X(x) = P(X \leq x) = \sum_{x_i \leq x} P(X = x_i) \text{ for DRV}$$

x_i : i th value taken by the DRV X .

Example 1: Find the probability that the commute time in the USA is less than equal to 45 minutes.

$$\begin{aligned} \text{We need to find } P[X \leq 45] &= \int_{-\infty}^{45} f_X(x) dx \\ &= 1 - P[X > 45] = 1 - \int_{45}^{150} f_X(x) dx \end{aligned}$$

$$P[X \leq 45] = 1 - 0.15 = 0.85.$$

Example 2: In the pop quiz example, find the prob that a student fails in the exam.

Recall, condition of passing the exam = responses to at least two questions are correct.

Prob that a student fails in the exam = $P[X \leq 1]$

Example 2: In the pop quiz example, find the prob that a student fails in the exam.

Recall, condition of passing the exam = responses to at least two questions are correct.

Prob that a student fails in the exam = $P[X \leq 1]$

$$P[X \leq 1] = \sum_{x_i \leq 1} P[X=x_i] = P[X=0] + P[X=1]$$

$$P[X \leq 1] = \frac{1}{8} + \frac{3}{8} = \frac{4}{8}$$

$$\text{or } P[X \leq 1] = 1 - P[X \geq 2] = 1 - \frac{4}{8} = \frac{4}{8}.$$

Properties of a CDF:-

- ① $F_X(x)$ is monotonically non-decreasing function
 $F_X(x_1) \geq F_X(x_2)$ if $x_1 > x_2$
- ② $0 \leq F_X(x) \leq 1$ for all x .
- ③ $F_X(\infty) = P(X \leq \infty) = \int_{-\infty}^{\infty} f_X(t) dt = 1$
- ④ $F_X(-\infty) = \int_{-\infty}^{\infty} f_X(t) dt = 0$
- ⑤ If X is continuous, the pdf can be obtained $f_X(x) = \frac{d}{dx} F_X(x)$

Some useful Probability distributions:

- 1 Normal distribution:- It is the most important continuous distribution because in many applications random variables are normal random variables (they have a normal distribution) or they are approximately normal or can be transformed into normal random variable. Furthermore, it is very useful due to Central limit

A random variable x is said to be Normal distributed or Gaussian distributed if it obeys the PDF

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

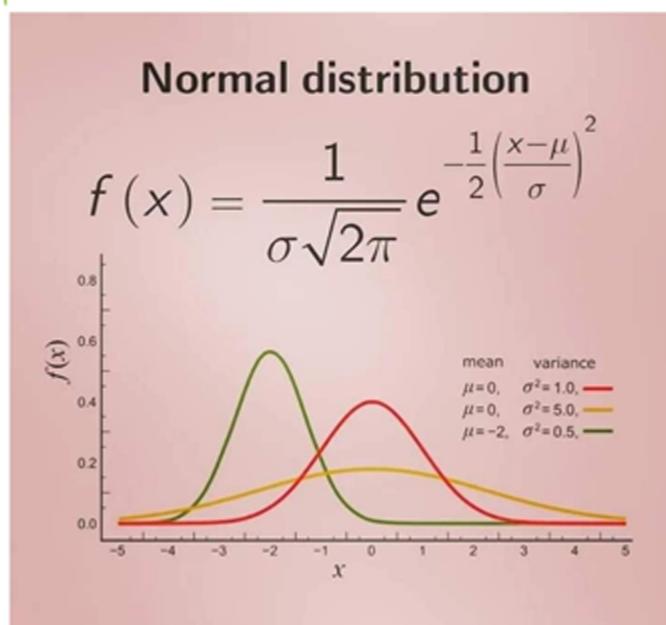
where μ is the mean and σ is the standard deviation

$$E[x] = \mu, E[(x-\mu)^2] = \sigma_x^2 = \sigma^2, \text{ and}$$

$\frac{1}{\sigma\sqrt{2\pi}}$ is a constant term that makes the area under the curve $f_x(x)$ from $-\infty$ to ∞ equal to 1.

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

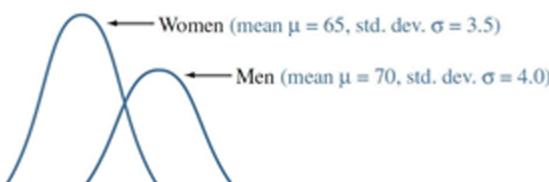
The PDF is symmetric with respect to $x=\mu$. It has a bell-shaped curve.



The PDF goes to zero very fast as σ or σ^2 decrease.

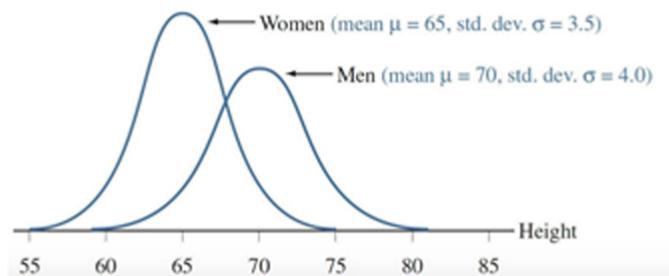
Some real-life example for normal distribution:-

- ① Heights of female and male students in the US follow a bell-shaped distribution. This can be approximated to a normal distribution as shown below.



Some real-life example for normal distribution:-

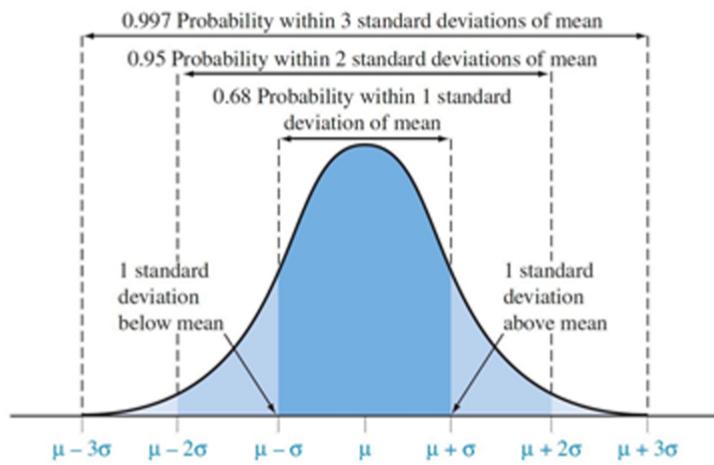
- ① Heights of female and male students in the US follow a bell-shaped distribution. This can be approximated to a normal distribution as shown below.



Important properties of normal distribution:-

- ① In practical work with the normal distribution, it is good to remember the probabilities within 1, 2, and 3 standard deviation away from the mean.

- ② $P(\mu - \sigma < X \leq \mu + \sigma) \approx 68\%$.
- ③ $P(\mu - 2\sigma < X \leq \mu + 2\sigma) \approx 95.5\%$.
- ④ $P(\mu - 3\sigma < X \leq \mu + 3\sigma) \approx 99.7\%$.



The multiple of standard deviation from the mean is denoted by a symbol z . For example, $z=2$ is 2 standard deviation away from the mean.

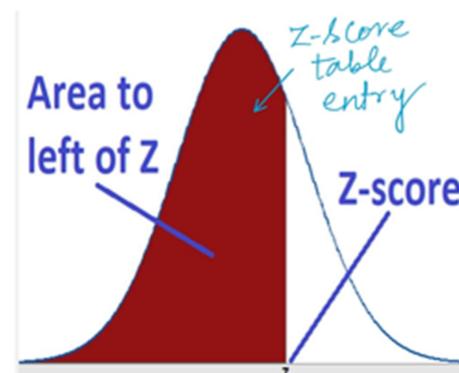
For a given z , the prob within z standard deviations

For a given Z , the prob within Z standard deviations away from the mean is the area under the normal curve between $\mu - Z\sigma$ to $\mu + Z\sigma$. This prob. is 0.68 for $Z=1$, as shown in Fig.

How to calculate Z for a certain probability?

It is calculated using z-score table. A part of this table is shown below. Table entry for Z is the area under the standard normal curve to the left of Z .

Second Decimal Place of z											
z	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09	
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359	
...											
1.3	0.9032	0.9049	0.9066	0.9082	0.9099	0.9115	0.9139	0.9147	0.9162	0.9177	
1.4	0.9192	0.9207	0.9222	0.9236	0.9251	0.9265	0.9278	0.9292	0.9306	0.9319	
1.5	0.9332	0.9345	0.9357	0.9370	0.9382	0.9394	0.9406	0.9418	0.9429	0.9441	



For example, let us find the value of Z for cumulative prob of 0.923. From the table, this corresponds to $Z = 1.4 + 0.03 = 1.43$. Thus, to get a cumulative prob of 0.923, we have to move 1.43 σ away from the mean.

$$\int_{\mu - 1.43\sigma}^{\mu + 1.43\sigma} f(x) dx = 0.923$$

where $f(x)$ is the normal distribution.

For a given cumulative probability, one can find Z-score from the table!

② Binomial distribution: (probabilities when each observation has two possible outcomes): In many applications, each observation is binary (has two possible outcomes).

For example, a person may

- Ⓐ accept or decline an offer from a bank for a credit card
- Ⓑ have or not have health insurance
- Ⓒ vote yes or no in a referendum, such as whether to provide additional funds to school.

② Binomial distribution: (Probabilities when each observation has two possible outcomes): In many applications, each observation is binary (has two possible outcomes).

For example, a person may

- (a) accept or decline an offer from a bank for a credit card
- (b) have or not have health insurance
- (c) vote yes or no in a referendum, such as whether to provide additional funds to school.

Binomial distribution is a prob distribution for a discrete random variable that takes two values.

Consider n cases, called trials, in which we observe a binary random variable. Let ' k ' denote the number outcomes of interest in n trial.

The prob of getting ' k ' success in ' n ' independent trials is given by the Binomial distribution as

$$P_X(x=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

where p is the prob of success and

$\binom{n}{k} = \frac{n!}{k!(n-k)!}$ is the binomial coefficient

$$E[X] = np, \sigma_X^2 = np(1-p)$$

The probability distribution can be interpreted as:

k success with prob p^k and $(n-k)$ failures occur with prob $(1-p)^{n-k}$. However, k success can occur anywhere among the n trials, hence there are $\binom{n}{k}$ different ways of distributing k -success in a sequence of n -trials.

Central limit theorem:-

Let $X_i, i=1, 2, \dots, n$, are independent and identically distributed (iid) random variables, each having

Let X_i , $i = 1, 2, \dots, n$, are independent and identically distributed (IID) random variables, each having a finite mean ' μ ' and a finite variance σ^2

Let Y be defined as the normalized sum, called sample mean

$$Y = \frac{1}{n} \sum_{i=1}^n X_i$$

The random variable Y is frequently encountered in estimating the mean of a random variable X from a number of observations X_i , for $i = 1, 2, \dots, n$.

$$E[Y] = \mu_Y = \frac{1}{n} \sum_{i=1}^n E[X_i] = \frac{1}{n} \sum_{i=1}^n \mu = \mu$$

$$\text{Var}[Y] = \sigma_Y^2 = E[Y^2] - \mu^2 = \frac{\sigma^2}{n}$$

$$\mu_Y = \mu, \quad \sigma_Y^2 = \frac{\sigma^2}{n}$$

When Y is viewed as point estimate for the mean μ , we note that its expected value is μ and its variance is decreases as n increases.

The central limit theorem states that the sum of statistically independent and identically distributed (IID) random variable approached to a Normal (Gaussian) distribution as $n \rightarrow \infty$.

$$Y \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

In more simple form, the central limit theorem states that if we take sufficiently large number of samples from a population, the sample mean is Normally distributed even if the population is not Normally distributed.

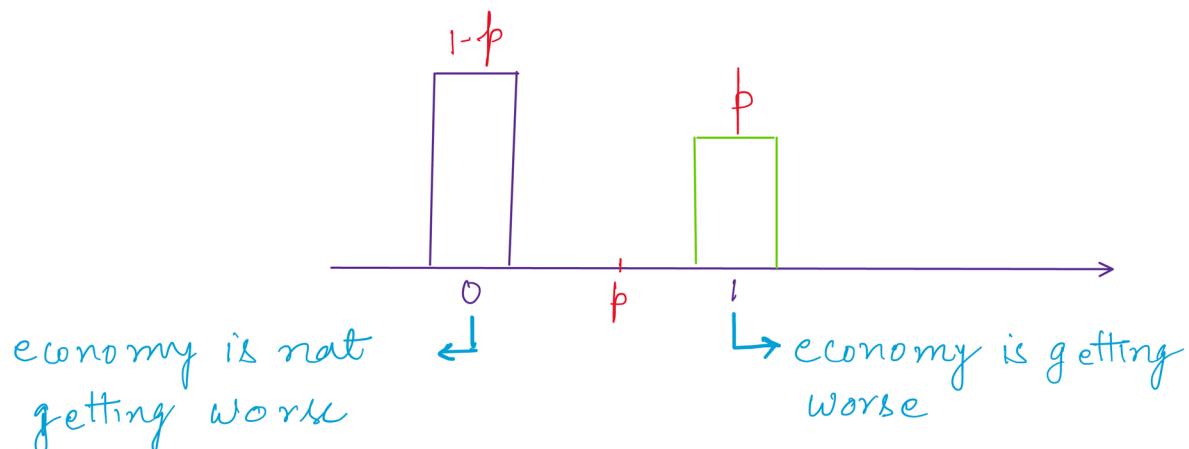
Theory SI Problem Solving

Friday, 27 June 2025 8:45 PM

Since we have understood probability, random variable, Central limit theorem, we can now solve the problem stated at the starting.

Problem: Assume that 40% population in India says that the economy is getting worse. If we take a sample of 250 people, how likely is that 50% or more of them will say they think the economy is getting worse. Calculate 99% confidence interval for proportion of the population who felt that the economy is getting worse.

Solution: We will not be able to survey the entire population to answer this question. But the entire population can be put in two buckets as follows



This is a binomial distribution.

Let we derive a sample of 250 people, in which 108 say the economy is getting worse.

We have 250 samples with 108 in a sample

$$\text{The sample mean } \bar{x} = \frac{1(108) + 0(142)}{250}$$

$$\bar{x} = \frac{108}{250} = 0.432 \quad (\text{point estimate of the population who think that the economy})$$

250 population who think that the economy is getting worse)

$$\text{Sample variance } \sigma_1^2 = \frac{108(1-0.432)^2 + 142(0-0.432)^2}{250-1}$$

$$\sigma_1^2 = \frac{34.8434 + 26.5}{249} = \frac{61.344}{249} = 0.246$$

$$\sigma_1 = \sqrt{0.246} = 0.4963$$

We want a 99.1% confidence interval.

Confidence interval :- It is a range of values

that likely contains the true population parameter, with a specified level of confidence (like 95% or 99%)

Formula for a confidence interval (for population mean)

$$CI = \bar{x} \pm z \left(\frac{\sigma}{\sqrt{N}} \right)$$

where \bar{x} = sample mean

z = z-score for the desired confidence level

σ = population standard deviation (or sample standard deviation if σ is unknown)

N = sample size

In our case, where $N=1$ and σ is the standard deviation of the population.

Using CLT, we can say that \bar{x} is coming from a Normal distribution with mean $\mu=\bar{x}$ and $\sigma_{\bar{x}} = \sigma/\sqrt{N}$

Recall that the variance of the sample mean

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} = \frac{\sigma^2}{250}$$

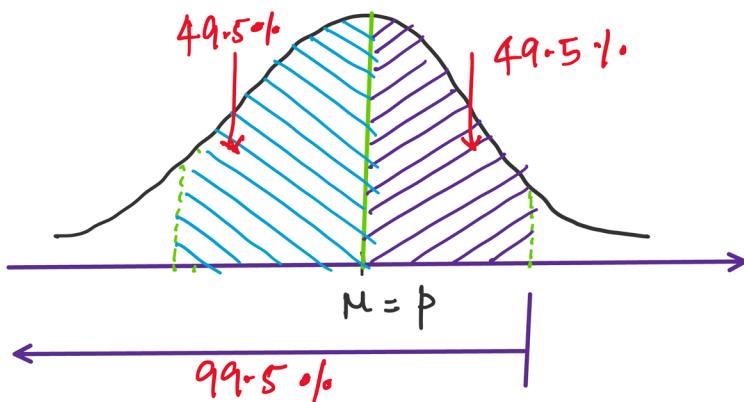
Recall that the variance of the sample mean

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} = \frac{\sigma^2}{250}$$

We do not know σ , so

$$\sigma_{\bar{x}} \approx \frac{\sigma_1}{\sqrt{250}} = \frac{0.4963}{\sqrt{250}} = 0.0314$$

For calculating the CI, we need to find Z .



from the z-table, $0.5 + 0.495 = 0.995$ corresponds to $Z = 2.58$. Thus, $2.58 \sigma_{\bar{x}}$ away from the mean μ would give us the 99% confidence interval.

$$\text{Thus the 99% CI} = \bar{x} \pm 2.58 \sigma_{\bar{x}}$$

$$2.58 \sigma_{\bar{x}} = 2.58 \times 0.0314 = 0.0810$$

$$\text{CI upper} \quad 0.432 + 0.08 = 0.513$$

$$\text{CI lower} \quad 0.432 - 0.08 = 0.352$$

We are 99% confident that true population proportion is within the range 35.2% to 51.3%
or

the true %age of the people who think that the economy is getting worse is in the range 35.2% to 51.3%. There is 99% chance of this.