

Population:

A group of individuals under study is called a "population".

Sample: (To generalise population).

A part of the population is called a sample.

Measure of Central tendency	Arithmetic Mean Median Mode Harmonic Mean Geometric Mean
" " average	
" " location	

HM:

$$HM = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}$$

reciprocal of mean of reciprocals.

H.M is calculated for a given set of observations x_1, x_2, \dots, x_n

$$x_n, HM = \frac{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}{n}$$

GM:

$$GM = \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n} = (x_1 \cdot x_2 \cdot x_3 \cdot \dots \cdot x_n)^{\frac{1}{n}}$$

$$\log G = \frac{1}{n} \log (x_1 \cdot x_2 \cdot x_3 \cdot \dots \cdot x_n)$$

$$= \frac{1}{n} (\log x_1 + \log x_2 + \dots + \log x_n)$$

$$G = \text{Anti log} \left[\frac{1}{n} (\log x_1 + \log x_2 + \dots + \log x_n) \right]$$

Average - Balancing the data

→ NSSO - National Sample Survey Organisation

(under Ministry of Finance Govt. of India)

Measure of skewness

Always tells you about whether the data is symmetric or not.

Measure of kurtosis:

The kurtosis always deals with the study of the frequency curves and its shapes.

Measure of central tendency

Measure of dispersion

Measure of skewness

Measure of kurtosis

Range

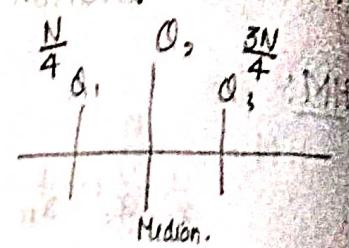
D.P. Problem related to range

M.D. $\frac{\sum |X - M|}{N}$ (M. mean, 100 median)

S.D. $\sqrt{\frac{\sum (X - \bar{X})^2}{N}}$ (S. standard deviation)

Quartile deviation:

$$Q.D. = \frac{Q_3 - Q_1}{2} \quad \text{third quartile - first quartile}$$



median - python prog.

$$Md = l + \frac{\frac{N}{2} - m}{f} \times c$$

$$\frac{l + l + f}{n} = MH$$

Standard deviation:

$$S.D. = \sqrt{\frac{1}{n} \sum (x_i - \bar{x})^2} \quad (\text{ungrouped data})$$

$$S.D. = \sqrt{\left(\frac{\sum f_i x_i^2}{\sum f_i} \right) - \left(\frac{\sum f_i x_i}{\sum f_i} \right)^2} \quad (\text{grouped data})$$

The square of the standard deviation is called variance.

$$\text{Var.} = f(S.D.)^2$$

Variation in the series of observations.

Ex: 15, 20, 15, 16, 17, 19, 25, 30, 45, 20 - Find Standard deviation.

x_i	$(x_i - \bar{x})$	$(x_i - \bar{x})^2$
15	-2.2	5.84
15	-2.2	5.84
16	-1.2	3.84
17	-0.2	0.04
19	0.8	0.64
20	1.8	3.24
20	1.8	3.24
25	6.8	46.24
30	11.8	139.24
45	22.8	519.84

$$\bar{x} = \frac{222}{10} = 22.2 \quad \sum (x_i - \bar{x})^2 = 777.6$$

$$S.D = \sqrt{\frac{1}{n} \sum (x_i - \bar{x})^2}$$
$$= \sqrt{\frac{1}{10} \times 777.6}$$
$$= \sqrt{77.76}$$
$$= 8.81816$$

Percentile:

Percentile is used to understand and interpret data.

The n^{th} percentile of a set of data is the value at which n percent of data is below it.

Ex: Calculate 45th percentile of the above data.

We use the following formula $L_p = (n+1) \frac{P}{100}$

15

15

16

17

19

20

20

25

30

45

$$L_{45} = (10+1) \frac{45}{100}$$

$$= 11 \times 0.45$$

$$= 4.95$$

$$P_{45} = 17 + (0.95)2$$

$$= 18.8$$

find 45th percentile

$$L_{45} = (10+1) \frac{45}{100}$$
$$= 11 \times 0.45$$
$$= 4.95$$

$$P_{45} = 0.5 + (0.25) 5$$
$$= Q_6 - Q_5$$

$$L_{25} = (10+1) \frac{25}{100}$$
$$= 11 \times 0.25$$

$$P_{25} = 15 + (0.75) 1$$
$$= 15.75$$

Range - 30
Q.O - 5.25
S.D - 8.81

n = int(input(" "))

x = len(n)

n.sort()

if x % 2 == 0:

median1 = n[x//2]

median2 = n[x//2 - 1]

median = (median1 + median2)/2

else:

median = n[x//2]

print("median is:", median)

Covariance :

Variation b/w 2 observations.

(x₁, y₁)

⋮

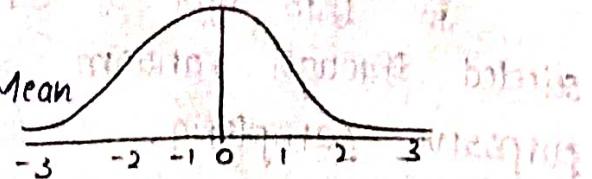
(x_n, y_n)

$$\text{Covariance} = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})$$

Skewness:-

$$S_k = \frac{\text{Mean} - \text{Mode}}{\sigma} \quad (-1 < S_k < 1)$$

$$\text{Mode} = 3\text{Median} - 2\text{Mean}$$

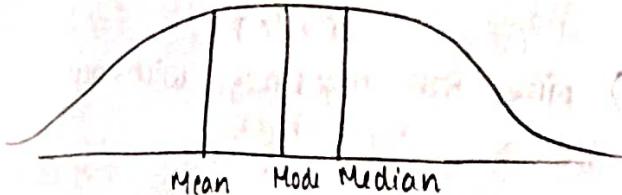


$$= \frac{\text{Mean} - (3\text{Median} - 2\text{Mean})}{\sigma}$$

$$S_k = \frac{3(\text{Mean} - \text{Median})}{\sigma}$$

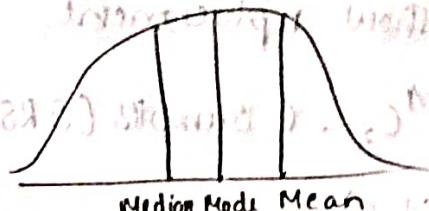
$$\text{Skewness} = 0 \\ (\text{Mean} = \text{Median} = \text{Mode})$$

↓
symmetrical distribution



(+vely skewed)

$$(0 < S_k < 1)$$



(-vely skewed)

$$(-1 < S_k < 0)$$

- To know the characteristics of data we find skewness, mean, median...

Sampling:-

Sampling can be done in 3 types

- Random Sampling
- Purposive Sampling
- Mixed Sampling

Random Sampling:-

In this type each unit in the population has an equal chance to include in the sample.

Purposive Sampling:-

Here the sampling units are selected from the population for the purpose.

Mixed Sampling:-

In this type of sampling some of the units are selected through random & remaining from purposive sampling.

→ In Random Sampling there are 2 types

- SRSWOR
- with replacement

Ex:- 2 numbers are to be drawn from 4 (4,5,6,7)
without replacement

${}^4C_2 = 6$ samples (SRSWOR) - simple Random Sample with out replacement

(4,5)

(4,6)

(4,7)

(5,6)

(5,7)

(6,7)

with replacement

${}^4^2 = 16$ samples (SRSWR) - simple Random sample with replacement.

(4,4) (5,4) (6,4) (7,4)

(4,5) (5,5) (6,5) (7,5)

(4,6) (5,6) (6,6) (7,6)

(4,7) (5,7) (6,7) (7,7)

problem:- (1M)
A population consist of 6 units draw the samples
with the two schemes simple Random sampling with
and without replacement of a sample 2 .

Population = 1, 2, 3, 4, 5, 6.

without replacement

$${}^3 C_2 = \frac{6 \times 5}{2} = 15 \text{ samples.}$$

(1, 2) (2, 3) (3, 4) (4, 5) (5, 6)

(1, 3) (2, 4) (3, 5) (4, 6)

(1, 4) (2, 5) (3, 6)

(1, 5) (2, 6)

(1, 6)

with replacement

$$6^2 = 36 \text{ samples}$$

(1, 1) (2, 1) (3, 1) (4, 1) (5, 1) (6, 1)

(1, 2) (2, 2) (3, 2) (4, 2) (5, 2) (6, 2)

(1, 3) (2, 3) (3, 3) (4, 3) (5, 3) (6, 3)

(1, 4) (2, 4) (3, 4) (4, 4) (5, 4) (6, 4)

(1, 5) (2, 5) (3, 5) (4, 5) (5, 5) (6, 5)

(1, 6) (2, 6) (3, 6) (4, 6) (5, 6) (6, 6)

- We have 3 techniques in Random Sampling
 - Simple Random Sampling
 - Stratified Random Sampling
 - Systematic Sampling
- Stratified R.S :- (no restrictions for population)

If the population is not homogeneous and if it is divided into k sections such that N_1, N_2, \dots, N_k are the respective population sizes in k sections of the pop. Then if I have to draw a sample of size n from the population then I can draw n_1, n_2, \dots, n_k samples from the respective sections in the population.

then

$$N_1 + N_2 + \dots + N_k = N$$

$$n_1 + n_2 + \dots + n_k = n$$

Then the sample is called stratified R.S.

- Systematic Sample :- (the population must be in systematic order)

In systematic sampling the population size must be known. Population must be arranged in an systematic order. So, that it is convenient to draw sample units from population.

Example:-

If I want to check 10 bills out of 100 bills generated by the shop keeper then divide the population size by sample size i.e., $\frac{100}{10} = 10$. ($\frac{N}{n} = k$)

Now the first unit in systematic sample is selected at random say 5.

For every next sample unit add k to the random number.

The next sample units from the bills are 5tk i.e., 15 tk
 The sample units are 5, 15, 25, 35, 45, 55, 65, 75, 85 & 95
 There is restriction on sample. The population size
 must be known. Should be in systematic order.

In any experiment if the result is unique, then it is a
 non-random experiment (example: Chemistry & Physics experiments)

If the result is coming from several outcomes then
 it is a Random Experiment.

Random Variable:

A variable associated with a random experiment is called a Random Variable.

Or - A function to map the sample space to the real line $(-\infty, \infty)$.

A function that maps the sample space of a random experiment to the real line $(-\infty, \infty)$.

Probability Distribution:-

A probability distribution is always defined for the values assumed by a RV. If x_1, x_2, \dots, x_n are the values assumed by 'X' for each value of the R.V. If we assign probabilities & sum of the Prob. is equal to 1. Then such a pair $(x_i, P(x_i))$ is called a prob. dist.

$$X: x_1, x_2, \dots, x_n$$

$$P(x_i): P(x_1), P(x_2), \dots, P(x_n)$$

$$\{x_i, P(x_i)\}$$

$$\sum_{i=1}^n P(x_i) = 1$$

Ex:- Write the prob. dist. if a single die is thrown, write for the faces of the die (X).

X	1	2	3	4	5	6
$P(X=x)$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

$$\sum P(x) = 1.$$

Binomial distribution:

A discrete R.V ' X ' is said to have a binomial distribution if its probability mass function is given by

$$P(x) = {}^n C_x p^x q^{n-x}, \quad x=0,1,2 \dots n$$

$$p+q=1.$$

Bernoulli R.V.:-

A R.V 'X' is said to be a Bernoulli R.V if $X=1$ with probability of success p .
 $X=0$ with probability of failure q ,
such that, $p+q=1$.

Problem:-

A python program is said to have 5% bugs. A sample of 10 lines were inspected what is the probability of getting 3 bugs in the sample.

$$p = 5\% = 0.05 \quad q = 1 - 0.05 = 0.95$$

$$n = 10 \quad x = 3 \\ P(x) = {}^{10} C_3 (0.05)^3 (0.95)^{10-3}$$

$$P(3) = {}^{10} C_3 (0.05)^3 (0.95)^7 \\ = 0.01047$$

2. A bolt manufacturing company claims that there will be 3% defectives in their total output. An inspection committee inspected a sample of 5 bolts. What is the prob that there will be 3 defectives in sample?

$$p = 0.03 \quad q = 0.97$$

$$n = 5 \quad x = 3$$

$$\begin{aligned} P(3) &= {}^5C_3 (0.03)^3 (0.97)^{5-3} \\ &= {}^5C_3 (0.03)^3 (0.97)^2 \\ &= 2.51013 \times 10^{-4} \end{aligned}$$

```
from scipy.stats import binom
```

```
import matplotlib.pyplot as plt
```

```
n = 5
```

```
p = 0.03
```

```
x = list(range(n+1))
```

```
dist = [binom.pmf(r, n, p) for r in x]
```

```
plt.plot(x, dist)      plt.bar(x, dist) → bar graph.
```

```
plt.show()
```

Poisson distribution:

A discrete R.V. 'X' is said to have poisson distribution if its pmf is given by

$$P(x) = \frac{e^{-\mu} \mu^x}{x!}, x=0,1,2,\dots,\infty, \mu > 0$$

The mean and variance of dist. are equal and equal to μ

Example:

In a bus stop junction the mean no. of accidents is '2'. What is the probability that there will be 4 accidents on a particular day. Since the no. of accidents is a poisson variate.

No. of accidents - poisson variate

$\therefore X$ is the no. of accidents.

$$P(x) = \frac{e^{-\mu} \mu^x}{x!}$$

$$P(4) = \frac{e^{-2} 2^4}{4!} = 0.09.$$

If newly published book have 500 pages and it is found to be 2 typographical errors per page on an average. Obtain the no. of pages with 0, 1, 2, 3, 4 errors.

The no. of errors in a page is a
poisson variate

and given as average 2 errors

$$\mu = 2$$

The poisson probability mass function

$$P(x) = \frac{e^{-\mu} \mu^x}{x!}$$

No. of pages with 0 errors = $P(0) \cdot 500 = 67.66 \approx 68$

$$\therefore P(1) \cdot 500 = 135.3 \approx 135$$

$$\therefore P(2) \cdot 500 = 135.5 \approx 135$$

$$\therefore P(3) \cdot 500 = 90.22 \approx 90$$

$$\therefore P(4) \cdot 500 = 45.1 \approx 45$$

```

from scipy.stats import poisson
import matplotlib.pyplot as plt
x = poisson.rvs(mu=2, size=10)
plt.hist(x, density=True, edgecolor='black')
plt.show()

```

Continuous Random Variable:

A R.V which assumes values bw certain limits

Ex: height and weight

*Normal distribution:

A continuous R.V is said to have normal distribution if its probability density function is given by

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

$-\infty < \mu < \infty$

$\sigma > 0$

where μ & σ are the arithmetic mean & standard deviation of the distribution.

$$\text{if } z = \frac{x-\mu}{\sigma}$$

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

where z is known as standard normal variate.

The normal distribution having the following characteristics:

1. The curve is bell shape and symmetrical about the line

$$x=\mu$$

2. Mean, Median & Mode of the distribution coincide

3. As x increases numerically, $f(x)$ decreases rapidly,

The max. probability occurring at the point $x=\mu$.

4. $\beta_1 = 0$ and $\beta_2 = 3$ (measure of kurtosis)
(skewness parameter)

5. Since $f(x)$ being the probability can never be -ve. No portion of the curve lies below the x -axis.

6. linear combination of independent normal variates is also

a normal variate.

7. x-axis is an asymptote to the curve (never touches x-axis)
8. The points of inflection of the curve are given by,
 $x = \mu \pm \sigma$ (third derivative)

In kurtosis we have 3 types of curves

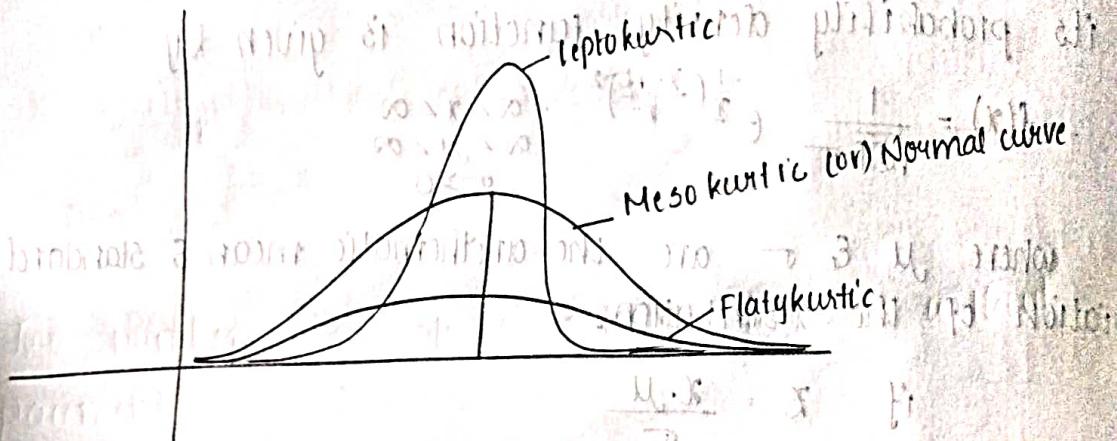
i. leptokurtic : The curve which is peaked.

ii. platykurtic : The curve which is almost flat

iii. Mesokurtic : curve which is neither peak nor flat.

or, $B_2 = 3$ (Always)

Normal curve



It seems that the x-axis is touching the curve but it will never touch the x-axis. Therefore, x-axis is an asymptote to the normal probability curve.

From the area property of the normal probability curve

$$P(|Z| \leq 3) = 0.9973.$$

Z always lies between -3 and 3.

Ex:- The mean and standard deviation of a normal distribution are 12 & 3 respectively. Find the prob. of

i. $x \geq 10$ $\sim N(\mu, \sigma^2)$

ii. $x \geq 15$

iii.

mean of $X = 12$

S.D. of $X = 3$:

$$N(\mu, \sigma^2) \quad X \sim N(12, 3^2)$$

Now if now $Z_1 = \frac{x - \mu}{\sigma}$ is a standard normal variable

standard normal distribution is the standard form

i.e., $P(X \geq 10) \approx$ probability of getting at least

$$\text{standard deviation} = Z = \frac{10 - 12}{3} = \frac{-2}{3} \approx -0.666$$

Normal variable

. or standard

Exponential distribution:-

A continuous R.V is said to have exponential distribution

its pdf $f(x) = \mu e^{-\mu x} \quad x > 0$

$$\text{mean} = \frac{1}{\mu}$$

$$\text{variance} = \frac{1}{\mu^2}$$

Testing of Hypothesis:

Since the population is very large to study the characteristics of population we take a sample from the population through this sample we try to estimate or test the population parameters.

Parameter:

A population const. is known as parameter.

Ex: Population arithmetic mean of a characteristic which can be measured.

Statistic:

A function of sample of the observations is known as statistics.

Ex: Sample mean, Sample Variance, Sample standard deviation

If there is a population of size 'N' and if we want to draw sample of 'n' we can have $N \times n = k$ samples, k samples can be drawn for each sample, if we calculated sample mean then it is called a sampling dist. of the statistic (\bar{x}) that is given in following table.

Sample no. Sample mean

1

$$\bar{x}_1$$

2

$$\bar{x}_2$$

⋮

$$\vdots$$

n

$$\bar{x}_n$$

k

$$\bar{x}_k$$

Ex: Two numbers are to be selected from 4 without replacement

$${}^4C_2 = \frac{4 \cdot 3}{1 \cdot 2} = 6$$

Sample no. Sample observations, Sample mean (\bar{x})

1 (4, 5) 4.5

2 (4, 6) 5.0

3 (4, 7) 5.5

4 (5, 6) 5.5

5 (5, 7) 6.0

6 (6, 7) 6.5

Standard error of the sample mean $\bar{x} \rightarrow S.E(\bar{x})$.

A population consists of

The above table represents the sampling dist. of statistic, Sample mean (\bar{x})

Standard error:

The S.D. of the sampling dist. of the statistic is known as standard error. It is denoted by S.E (statistic)

$$\bar{x} = \frac{4.5 + 5 + 5.5 + 5.5 + 6 + 6.5}{6} = 5.5$$

$$S.E = \sqrt{\frac{1}{n} \sum (x_i - \bar{x})^2}$$

x_i	$(x_i - \bar{x})$	$(x_i - \bar{x})^2$
1.5	-4	1
5	-0.5	0.25
5.5	0	0
5.5	0	0
6	0.5	0.25
6.5	1	1
		$\bar{x} = 5.5$
		$\sum = 2.5$

$$= 0.645$$

∴ Standard Error = 0.645