

II Descriptive & Inferential Statistics

* Descriptive statistics :-

1. Population & sample :-

- Population is the entire group that you want to draw conclusions about.
- A sample is the specific group that you will collect data from. size of the sample is always less than the total size of the population.

Population

sample.

I Advertisements of IT jobs
in Netherland.

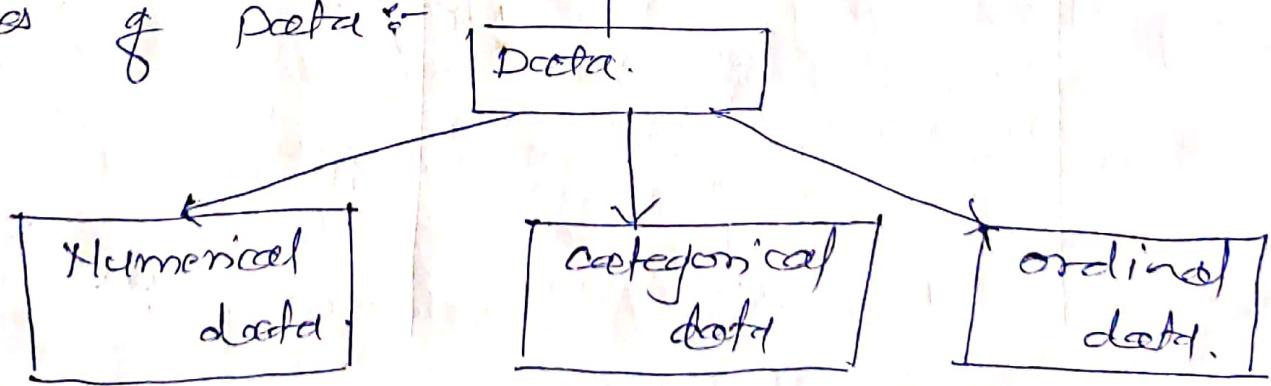
Top 20 search results for
advertisements for IT jobs
in Netherlands on June 1 2022
including songs from Eurovision
song contest that were
performed in English.

- songs from Eurovision song
contest.

200 eug students from
8 universities for research
study.

2. Undergraduate students in
Netherlands

* Types of Data :-



* Numerical data :-

- Data includes count or measurement of any object or person such as mass, volume, height, sugar level etc.

* Categorical data :-

- The characteristic or behavioural attribute of person or object is included under categorical data. It can be in any form such as sex/gender, nativity, caste, marital status or a favorite thing.
- It can be either true or false.

* Ordinal data :-

- Mixture of numerical & categorical data.
- Data is arranged under categories & numbers are placed in the categories which give a correct meaning.
e.g. Hotel rating can be given from one to five with the category from poor to excellent considering its ambience, taste, service, cost facilities etc. By merging data, the chart is prepared to analyze the performance of the hotel.

④ Measurement Levels :-

1. Nominal	2. ordinal	3. interval	4. ratio
- categorized.	categorized & ranked data.	categorized, ranked & evenly spaced.	has a natural zero.

1. Nominal Level :-

- You can categorize your data by labelling them in mutually exclusive groups, but there is no order betⁿ the categories.

2. ordinal Level :-

- You can categorize & rank your data in an order, but you cannot say anything about the intervals betⁿ the rankings.
 - e.g. You can rank top 5 Olympic medalists. This scale does not tell you how close or far apart they are in number of wins.

3. Interval Level :-

- You can categorize, rank & infer equal intervals betⁿ neighbouring data pts, but there is no true zero point.
- The diff. betⁿ any two adjacent temp. is the same: one degree. But zero degrees is defined differently depending on the scale.

(4)

* Representation of categorical variables :-

- Categorical data is the statistical data type consisting of variables or of data that has been converted into that form, for example as grouped data.
- More specifically, categorical data may derive from observations made of qualitative data that are summarised as counts or cross tabulations or from observations of quantitative data grouped within given intervals. Often, purely categorical data are summarised in the form of a contingency table. However, particularly when considering data analysis, it is common to use the term 'categorical data' to apply to data sets that, while containing some categorical variables, may also contain non-categorical variables.
- A categorised variable that can take on exactly two values is termed a binary variable or dichotomous variable.
- Categorical variables with more than possible values are called polytomous variables: categorical variables are often.

* Dummy coding

- The reference group is assigned a value of 0 for each code variable
- Data are analyzed through comparing one group to all others, groups.

"Effects coding"

* Contrast coding :-

- It allows researcher to directly ask questions.
- ~~This tailored hypothesis~~
- sum of contrast codes is restricted by three rules:
 - sum of contrast coefficients per each code variable must equal zero.
 - diff. betⁿ sum of positive coefficients & the sum of the negative coefficients should equal 1.
 - coded variables should be orthogonal.

* Nonsense coding :-

- It occurs when one uses arbitrary values in place of designated 0's 1's & -1' seen in previous coding systems.

* Embeddings :- Embeddings are codings of categorical values into high-dimensional real valued vector spaces usually in such a way that similar values are assigned similar vectors or with respect to some other kind of criterion making the vectors useful for resp. applⁿ.

⑨ * Measures of central tendency :-

* mean :- The most commonly used measure of central tendency is the mean. To calculate the mean of a dataset, you simply add up all of the individual values & divide by the total no. of values.

$$\text{Mean} = (\text{sum of all values}) / (\text{total no. of values})$$

* median :- The median is the middle values in a dataset. Arrange all values in a dataset from smallest to largest & finding the middle value will give you median. If there are odd no. of values middle one is median. If there are even no. of values the median is the average of the two middle values.

e.g. Player.	1	2	3	4	5	6
Home Runs.	8	9	10	12	13	14

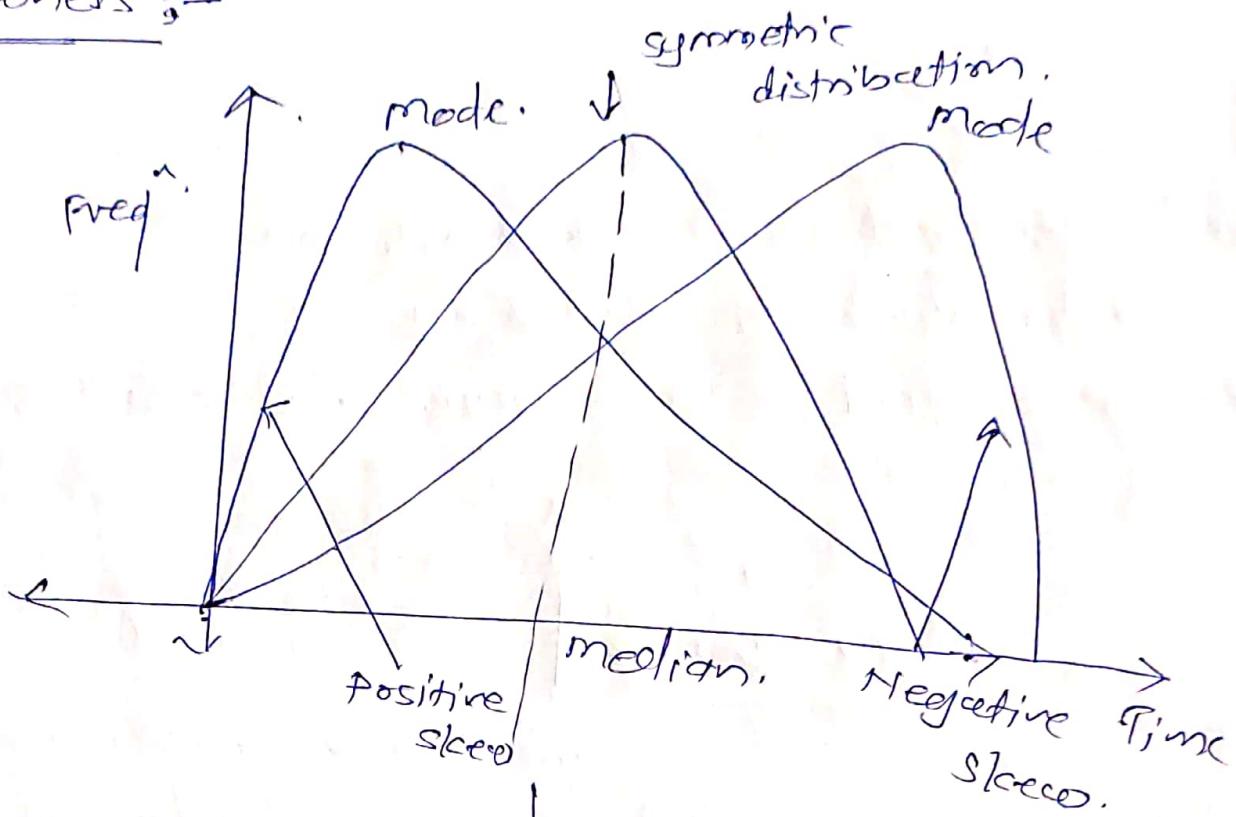
$$= \frac{10+12}{2} = \frac{22}{2} = \underline{\underline{11}}$$

* mode :- Mode is the value that occurs most often in a dataset. Dataset can have no mode, one mode or multiple modes.

e.g.	1	2	3	4	5	mode - 6 is
	8	9	10	10	11	

SKEWNESS :-

(7)



* Positive skewness

- If the given distribution is shifted to the left with tail on the right side, it is a positively skewed distribution.
- It is also called a right-skewed distribution.
- It assumes skewness value of more than zero.
- mean value is greater than the median & moves towards right. ^{mode} occurs at highest freqⁿ of distribution.

Negative skewness.

- If the given distribution is shifted to the right & with the tail on the left side, it is a negatively skewed distribution. It is also called a left-skewed distribution.

- The skewness value of any distribution showing a negative skew is always less than zero.

$$\text{skewness} = \frac{\bar{x} - m_o}{s}$$

\bar{x} = mean value

m_o = mode value

(8)

s = standard deviation of sample data.

* Variance :- is the expected value of the squared variation of a random variable from its mean value in probability & statistics. Informally, variance estimates how far a set of numbers are spread out from their mean value.

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

* Standard Deviation :- it measures dispersion of a dataset relative to its mean & is calculated as the square root of variance.

$$\text{standard deviation} = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}}$$

where, x_i = value of the i^{th} point in the dataset

\bar{x} = The mean value of the dataset

n = The number of data points in the data

(9)

* Variance vs std. deviation :-

Parameters	Variance	std. deviation
meaning	Numerical value that describes variability of observations from its arithmetic mean.	— It is measure of dispersion of observations within dataset.
What is it?	It's average of squared deviations.	It is root mean square deviation.
labelled as	σ^2 (sigma-squared)	σ (sigma).
Expressed in	squared units	same units as the values in the set of data.
Indicates	How far individuals in a group are spread out.	How much observations of a dataset differs from its mean.

* Coefficient of variation: 10

- It is a statistical measure of dispersion of data points around the mean.
- The metric is commonly used to compare the data dispersion between distinct series of data.
- By identifying determining coefficient of variation of different securities, an investor identifies the risk to reward ratio of each security & develops an investment decision.

Definition: Coefficient of variation = $\frac{\text{standard deviation}}{\text{mean}} \times 100\%$

S - standard deviation M - mean

Coefficient of variation = $\frac{\text{volatility}}{\text{Expected Returns}} \times 100\%$

ABC Corp. offered stock of variation 10% & financial performance. The volatility of stock is 10% & the expected return is 10%.

ANSWER

* covariance & correlation between two variables. (11)

Parameters

Meaning.

It indicates extent of variables being dependent on each other. Higher values denotes higher dependency.

Relationship:

→ Correlation can be gathered from covariance.

Values

Lie between -∞ to +∞

Scalability

Affects covariance

Units

Covariance will have a definite unit as it is concluded from multiplication of numbers & their units.

Covariance

Correlation

→ It signifies strength of association bet the variables when other things are constant.

Correlation gives value of covariance on a std. scale.

Correlation has limited values in range of -1 to +1. Correlation is not affected by a change in scale.

Correlation is a number without units but includes decimal values.

9096116898

* Inferential statistics :-

Distribution :- 3 types: sampling, central limit theorem, confidence interval.

- Sampling distribution is a type of probability distribution created by drawing many random samples of a given size from population.
- Distribution is simply collection of data or scores or variables ordered from smallest to largest & can be presented graphically.

Standard Error It measures accuracy with which a sample distribution represents a population by using std. deviation.

- Sample mean deviates from actual mean of a population. This deviation is std. error of mean.

$$SE = \frac{s}{\sqrt{n}} \quad s = \text{std. error}$$

- It indicates how different population mean is likely to be from sample mean.

Estimators & Estimates:

- Estimator is a statistic that estimates some fact about population. It is a rule that creates an estimate.

e.g. Sample mean (\bar{x}) is estimator of population mean.

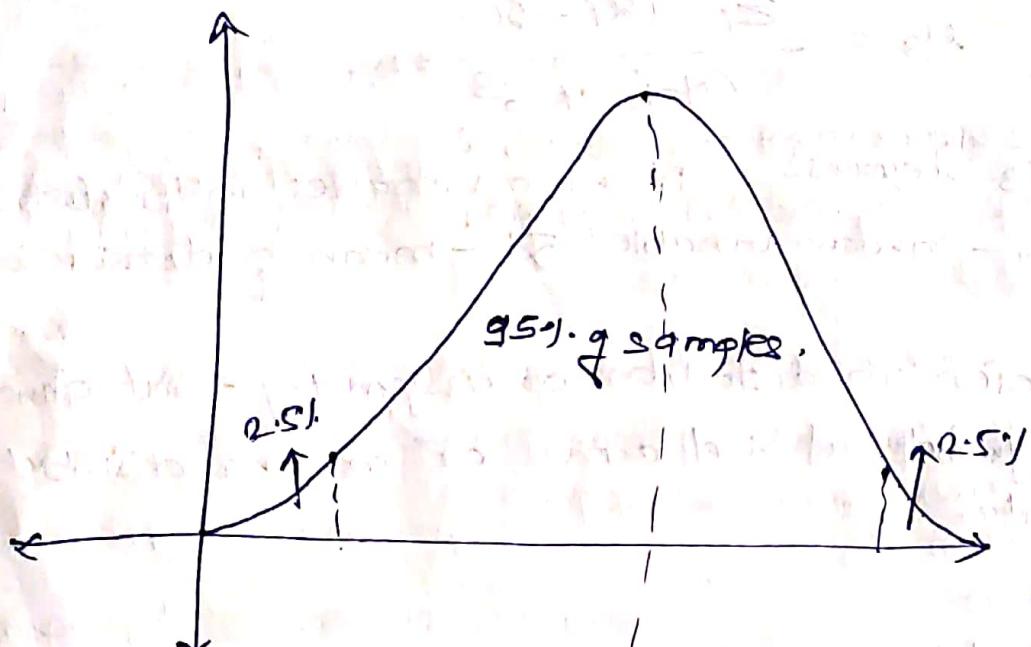
- Estimator is a function of sample & estimate is a value of an estimator calculated from a sample.

(B)

* central limit theorem :-

- Average of your sample means will be the population mean.
- Add up the means from all of your samples, find the average & that average will be your closest population mean.

Confidence interval :-



$$CI = \bar{x} + 2 \frac{s}{\sqrt{n}}$$

CI - confidence interval,

\bar{x} - sample mean

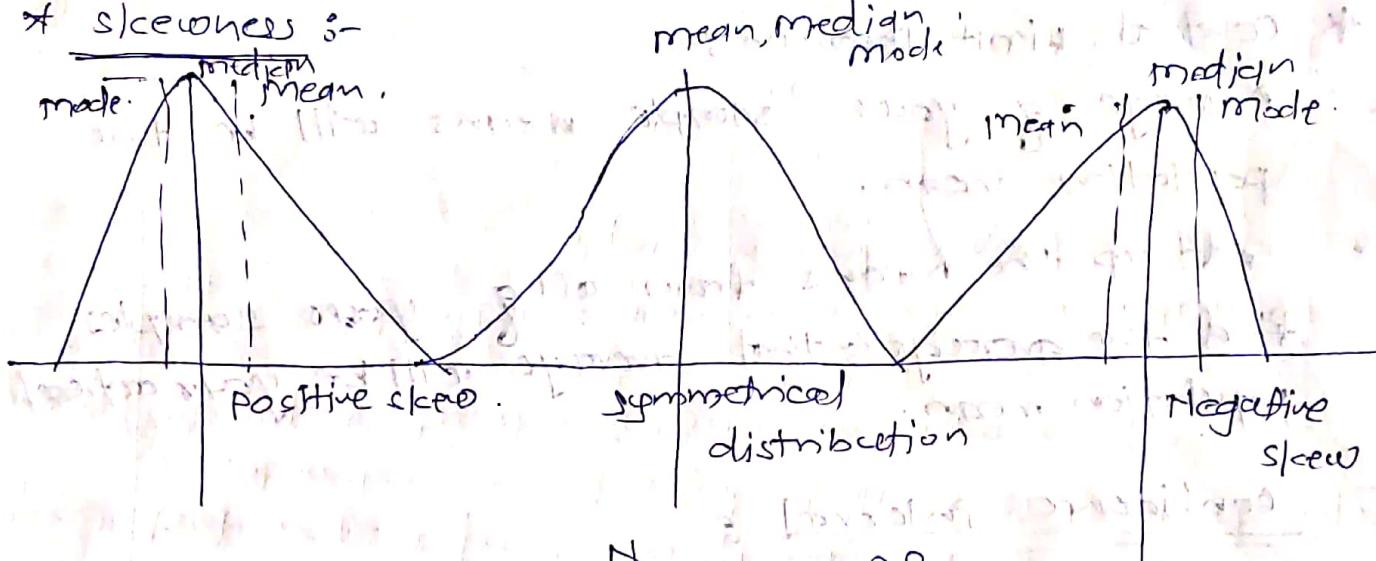
ϵ - confidence level value

s - sample std. deviat'

n - sample size

- CI is range of values that we observe in our sample & for which we expect to find the value that accurately reflects the population.

* skewness :-



$$Eg = \frac{\sum_i^N (x_i - \bar{x})^3}{(N-1) * 6^3}$$

Eg - skewness N - no. of variables in the distribution
 x_i - random variable \bar{x} - mean of distribution

- skewness refers to distortion or asymmetry that deviates from the symmetrical bell curve or normal distribution in a set of data.

* Variance :-

Measure of variability or dispersion

- It measures how far a set of numbers is spread out from their average value.

- measurement of spread between numbers in a data

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

- Measure how far data points from their mean.

x_i = value of one observation
 \bar{x} = mean value of all observations

n = no. of observations.

* Types of data:-

- Nominal.
- ordinal
- Discrete

- Continuous

Defn.

categorical or

quantitative data.

Nominal or

Quantitative data.

Normal Data.

Ordinal Discrete.

e.g. gender, hometown, colour etc.

Birthday, favorite sport, school postcode.

Nominal values but don't have numerical meaning.

- Nominal :- Label variables without providing numerical value. Also called nominal scale. It cannot be ordered or measured. e.g. Letters, symbols, words, gender etc.

- ordinal :- It follows natural order. Represented using bar chart found in surveys, finance, economics, questionnaires etc.

- Discrete data :- Takes only discrete values. contains only finite no. of possible values. Things can be counted in whole numbers & cannot be subdivided meaningfully.
e.g. no. of students in the class.

- continuous data :- It can be calculated & has an infinite no. of possible values that can be selected within a given specific range.

* coefficient of variation :-

$$CV = \frac{\sigma}{\bar{x}}$$

σ - population std dev
 \bar{x} - H-mean

- Ratio of std deviation to mean higher the coeff. of variation, greater the level of dispersion around the mean.
- Used to determine variability of data.

* Find the coeff. of variatn of following sample:

$$\{1, 5, 6, 8, 10, 40, 65, 88\}.$$

Ans:-

$$\text{Mean} = \frac{1+5+6+8+10+40+65+88}{8} = 223/8 = 27.875$$

$$\text{Variance} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} = \frac{7578.875}{8-1} = 1082.696$$

$$\text{std deviatn} = \sqrt{\text{variance}} = \sqrt{1082.696} = 32.904$$

$$CV = \frac{\sigma}{\bar{x}} = \frac{32.904}{27.875} = 1.180$$

* std. deviation :-

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

σ - std. dev.

N - size of population

x_i - each value of population

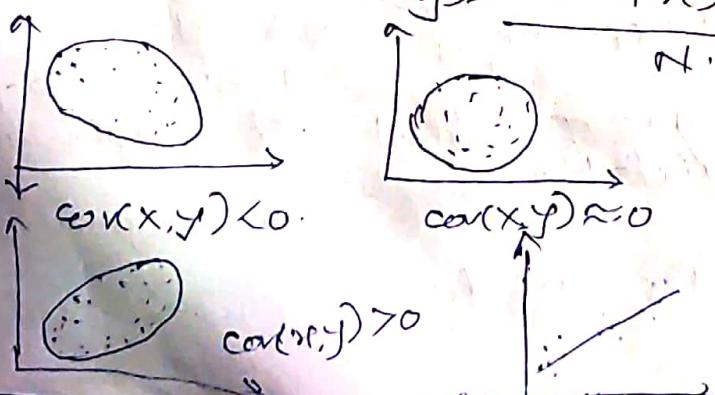
\bar{x} - population mean

- It is the measure of dispersion of dataset relative to its mean & calculated by square root of variance.

- It shows how much variation from the mean exists.

- key parameters
- meaning
- what pt is.
- values range
- change in scale

covariance.
Indicates the extent to which two random variables change in tandem.
Measure of correlation.
$-\infty$ to $+\infty$.
Affects covariance.
Affected by scale of variable.
Indicates direction of linear relationship.
not std. values



correlation.
statistical measure that indicates how strongly two variables are related.
scaled version of covariance.
-1 to $+1$.
Does not affect correlation.
not affected by scale of variable.
-1 - Direction & strength +1 -
standardized values.

covariance indicates how two variables are related to one another.
 covariance shows how two variables differ, & correlation shows you how two variables are related.