

BIOSTATICS



Biostatistics index

STATISTICS

Measures of Central tendency and Variation Normal distribution curve. Tests of significance Concept of Probability value Correlation, Regression and Skew Sampling methods and calculation Probability and Tiles Graphs Biostats review and QA round

EPIDEMIOLOGY

Analytical Epidemiology Advanced analytical study designs Experimental Epidemiology Evidence based medicine

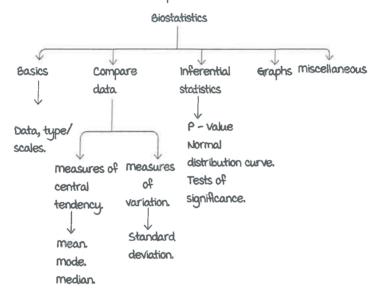
SCREENING

Advanced concepts in screening of disease

INTRODUCTION TO DATA IN **BIOSTATISTICS**

uses:

- · Define cut-offs.
- · Understand variation.
- To present data.
- · To make inference (provide evidence).



Data

00:07:54

Quantitative	Qualitative	
· Continuous.	• Discrete.	
 measurable. 	· Countable.	
• E.g. weight, height,	• E.g. No. of people	
AST, ALT levels.	who are sick/	
· mean of data can	healthy, alive/dead.	
be calculated.	Gender.	
	Proportions/	
	percentages can be	
	calculated	

Pulse rate is a data which is discrete and countable, however it is quantitative as we calculate its mean. SP is quantitative data.

Scales of data

00:14:25

Inherentorder.Has asequence.	•	Interval between two values is present.	 Ratio can be calculated There is
* Has a		two values	calculated
sequence.		is present.	• There is
			THE IS
• E.g. Stage,	•	No start	zero point/
grade, the		point/no	absolute
severity of		absolute	zero.
the disease.		zero.	• E.g. Na, 15,
		E.g. °C, d8.	FEV levels.
	grade, the severity of	grade, the severity of the disease.	grade, the point/no severity of absolute

interval type of data:

example: a0 °C is not half as hot as 40 °C, but colder compared to 40 °C. Here the intensity of data is measured. Also, the temperature can go below 0 °C (in minus °C), which means there is no absolute zero.

Ratios:

Example: A weak fragile child weighs 20 kg when the ideal weight should have been 40 kg in the same age group. The ideal weight is 2 x child's age, which means the values can be expressed in multiples (double, triple) of each other i.e calculation of ratios is possible.

Also, there is absolute zero/ no value below zero.

MEASURES OF CENTRAL TENDENCY AND VARIATION

Measures of Central tendency

00:01:59

mean:

- 1. Arithmetic mean:
 - Average = Σ (summation)

n

- a. Geometric mean:
 - · Calculated in case of: Exponential data.

Extreme values.

· Example: Human development index

(India = 0.647, ranked at 129 in 2019)

- 3. Harmonic mean:
 - · Calculated in case of: Inverse data.

Fractional values.

Advantages:

- · Best measure of central tendency.
- Easiest to calculate.

Disadvantages:

· most affected by extreme values.

median:

Central value after arranging in ascending or descending order.

Advantages:

· Least affected by extreme values.

mode:

The most frequently occurring value.

mode = 3 median - a mean.

Advantages:

- · The most robust measure of central tendency.
- . The last to be affected by extreme values.

Data with extreme values: Preferred measure is median.

Preferred mean is geometric mean.

1. Range:

Range = maximum to Minimum.

a. Standard deviation:

Gives the mean deviation of every value from the mean. Formula: The root of the mean of squared deviation.

$$SD = \sqrt{\frac{\sum (x - \bar{x})^a}{n}}$$

In case of a small sample,

SD =
$$\sqrt{\frac{\sum (x - \bar{x})^{a}}{n - 1}}$$
 n - 1 is the correction for the small sample (n < 30).

3. Variance:

Variance (V)= SDa

ril.dem=
$$\sum (x - \bar{x})^a$$

4. coefficient of variation (cv):

Absolute variation between a different populations.

$$CV = S.D. \times 100$$
mean

5. Standard error:

Gives the error in different studies in terms of standard deviation.

Alternatively, gives the variation between values when different researches are done.

a. Standard error for mean:

For quantitative data.

•
$$SE_m = \frac{SD}{\sqrt{n}}$$



Measures of Central Tendency and Variation

b. Standard error for proportions:

· For qualitative data.

• SE =
$$\sqrt{\frac{PQ}{N}}$$

P: Prevalence.

Q:100 - prevalence.

n: Sample size.

If p-value or Confidence interval is provided as input, Standard error has to be calculated and not the Standard deviation.

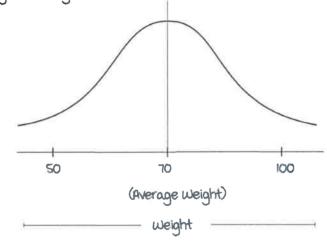
NORMAL DISTRIBUTION CURVE

Normal distribution curve

00:00:08

It represents the distribution of data in a bell-shaped čiúrve, in a large sample.

Eq: The weight of students in the class.



Features of Normal distribution curve:

It is also known as the Gaussian distribution curve.

It is a bilaterally symmetrical bell-shaped curve.

The ends never touch the baseline.

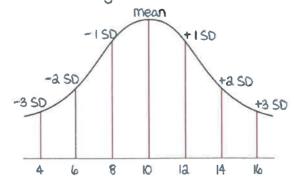
mean = median = mode \Rightarrow Coincide at 0 or the centrepoint.

SD = 1.

AUC = 1 (Area Under Curve), means the whole population is accounted for:

Eg: mean Hb (\bar{x} Hb) at a place = 10 gm% \pm 2 g%.

where ISO = a 9%

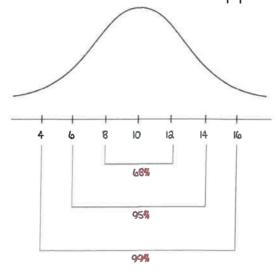


Active space

Assumptions in normal distribution curve: First assumption

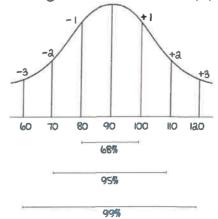
00:07:28

Between the -1 SD and +1 SD: 68% of the population lies. Between the -2 SD and +2 SD: 95% of the population lies. Between the -3 SD and +3 SD: 99% of the population lies.



Eq: Mean blood glucose = $90 \pm 10 \text{ SD}$. How much of the population will be expected to fall between:

- 80 to 100 mg/dl = 68% population.
- 70 to 110 mg/dl = 95% population.
- 70 to 100 mg/dl = 68% + 13.5 % population [(95-68)/a=13.5]
- more than 70 mg/dl = 100 2.5% = 97.5% population.
- Less than 100 mg/dl = 84% population (100-13.5+2+0.5).
- more than 100 mg/dl = 16 % population.
- Less than 60 mg/dl = 100 99 = 1/2 = 0.5% population.
- Less than 120 mg/dl = 100 0.5 % = 99.5% population.



MedPlus - 7850010383 Normal Distribution

Curve

Q. The mean blood glucose from 5929 ANC females in the state of maharashtra was found to be 130 ± 5 mg/dl. The cut off for diagnosing 610m was kept as higher than 140 mg/dl. How many pregnant females are expected to be 610m diagnosed?

To be 60M diagnosed, they must belong to above + 2 50 of population.

Above +a SD = 100 - 95% (between +a and -a SD) - 2.5% (less than -a SD) = 2.5%

a.5% of 59a9 $^{\sim}$ 150 females, which falls under range of 100-a00.

Second assumption: Zone of Normalcy 00:20:51

zone of normal cy/normal zone:

Between the - a SD and + a SD =
$$95\%$$
 of population.

2. score:

It is also called standard deviate.

It gives the location of the value in terms of the standard deviation (SD).

The cut off for Z score: ± a SD/± 1.96 SD.

If the 2 score is > a: Abnormal 2 score.

It is calculated by = Observed value - Expected value

50

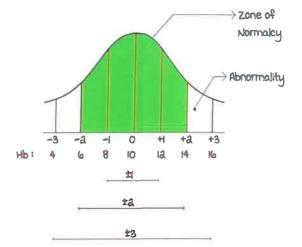
Eg: Observed value of Hb = 15 gm/dl. Expected value (always the mean value) = 10

$$SD = a$$

$$2. score = 15 - 10 = 2.5$$

а

2. score a.s : It lies a.s so away from the mean.



Active space

CONCEPT OF PROBABILITY VALUE

P value

00:01:05

P value:

Probability value (chance of events expressed in decimals).

Normal value ranges from 0 to 1.

0: Lowest probability.

1: Maximum probability.

Standard errors (SE):

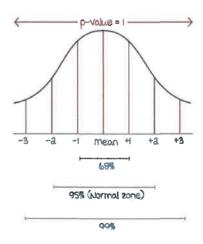
±1, ±2, ±3,...

confidence limit/interval:

+1 to -1 = 68% confidence interval

+a to - a = 95% confidence interval.

+3 to -3 = 99% confidence interval.



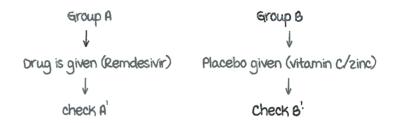
in the normal distribution curve:

The highest probability is towards the centre: I. The lowest probability lies on either side of the curve. At +a to -a standard deviation the P value is: 0.05 - 2 one of normalcy.

P value - abnormal zone

00:06:20

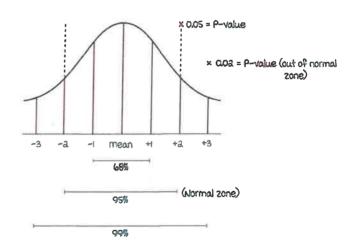
Example: Randomised clinical trial - two groups A and B



The collected data is incorporated in a machine : Gives P Value.

If the P value is 0.02: Abnormal/out of the normal zone.

P value > 0.05	P value < 0.05		
Normal variant	Abnormal variant		
Non-significant	significant		
No effect found	Effect is found		
Null hypothesis : Accepted	Null hypothesis : Rejected		



P value - normal zone and changes

00:16:42

The normal zone for P value - 95% confidence interval.

If the normal zone moved from 95% to 68%:

Previously non-significant becomes significant.

Chances of finding an effect increases.

The chances of reject of null hypothesis increases.

The chances of alpha error increases.

If the normal zone moves from 95% to 99%:

Previously significant becomes non-significant.

The chances of finding an effect decreases.

The chances of accepting of null hypothesis increases.

The chances of beta errors increase.

Alpha error, type I & II error

00:23:08

Definition:

It is the probability of finding an effect (just by chance) which in reality does not exist.

It corresponds to the P value/confidence interval/limit.

Example: P value of 0.02 corresponds to a value 2%.

It means there is 2% chance of error in the study.

It also means there is 98% of confidence in the study.

68% corresponds to 32% alpha. 95% corresponds to 5% alpha. 99% corresponds to 1% alpha.

FPER: The chance of finding disease in a healthy patient.

Tupe I error:

Rejecting a null hypothesis, which in reality is true.

Type II error:

Accepting a null hypothesis, which is false in reality.

TESTS OF SIGNIFICANCE

Statistical mathematical formula to derive a p-value. Determines if P-value is significant or non significant.

Types of tests of significance

00:02:59

Types

1. Parametric 2. Non-parametric

Quantitative Qualitative

Normal distribution data. Non-normal distribution data

Parametric test	Situation	Non-parametric test
Paired 't' test.	Single group	mc nermar's test.
unpaired 4 test A/K/A independent sample 4 test.	Two groups	Chi square test (y ^a).
Analysis of variance (ANOVA)	Three or more groups	Mruskal-wallis test. Chi square for trend.

Advance tests of significance

00:08:59

- Large sample (n > 30) = '2' test.
- · Ordinal data: wilcoxan rank test (W/R)

W/R sign test

W/R sum test

For grouped data

For ungrouped data

- · Normalcy of data: Holmogorov smirnov test.
- · Outliers : Dixon's Q test.
- Internal consistency of questionnaire : Cronbach's lpha score
- Compare a new test with a gold standard test: Bland altman analysis.

· Level of agreement : KAPPA test.

Formula = Observed level of agreement - expected level of agreement

1 - expected level of agreement

CORRELATION, REGRESSION AND SKEW

Correlation

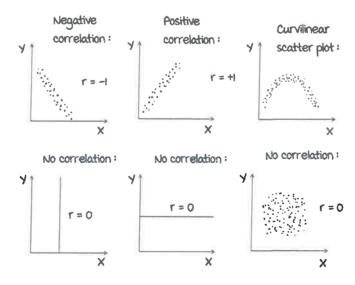
00:00:13.

Relation between a variables. Scatter plots are used.

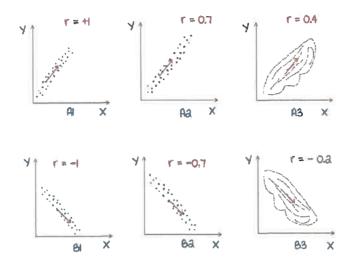


- Also known as Pearson-Karl correlation.
- Also Known as non-linear/
 Spearman correlation.
- · Represented by : r
- Represented by : p
- Range : -1 to +1
- -1: Perfect negative correlation.
- H: Perfect positive correlation.

r = 0 : No correlation.



Scatter plots



- +1: Perfect positive correlation (1 unit change in X axis = 1 unit change in Y axis).
- > 0.7: Strong positive correlation.
- 0.5 0.7: moderately positive correlation.
- < 0.5 : Weak correlation.
- < 0.3 : Very weak correlation.

Coefficient of determination (CD):

The percentage change in one variable which is accounted for by a unit change in another variable.

CD = ra in %.

Regression

00:18:26

Primarily refers to prediction.

Tupes

- 1. Linear: If variables as quantitiative.
- a. Logistic: If variables are qualitative.
 - 1. Univariate linear regression:
 - Eg: Predicting renal failure based on GFR.
 - a. Univariate logistic regression:
 - Eg: Predicting MI based on obesity levels.
 - 3. multivariate linear regression:
 - eg: Predicting the renal status based on serum Na, urea, creatinine and GFR levels.