

STATISTICS

(1) DESCRIPTIVE STATISTICS

(2) INFERENCEAL STATISTICS

* (1) DESCRIPTIVE STATISTICS

① CENTRAL TENDENCY - MEAN, MEDIAN,
MODE

② DISPERSION - VARIATION / VARIANCE,
STANDARD DEVIATION (SD)

③ DATA EXPLORATION - Histograms, pie charts,
Bar charts.

* There are ~~3~~ Types of Descriptive
Statistics:

① MEASURES of central tendency

② MEASURES of variability (spread) -

① Range ② variance ③ standard deviation

* (1) DESCRIPTIVE STATISTICS

(A) CENTRAL TENDENCY

MEAN | MEDIAN | MODE

(1) MEAN :-

MEAN is the average of a set of numbers. It is calculated by adding up all the numbers and dividing by the total count of numbers.

Formula :-

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

where,

- * \bar{x} represents the mean.
- * $\sum_{i=1}^n x_i$ represents the sum of all the values in the set.
- * 'n' represents the total number of values in the set.

Example :-

In the set $[5, 10, 15, 20]$, the mean is calculated as:

$$\text{Mean} = \frac{5+10+15+20}{4} = \frac{50}{4} \\ = 12.5$$

(2) MEDIAN :-

The MEDIAN is the middle value in a sorted set of numbers. If there are an odd number of values, the median is the middle value.

If there are an even number of values, the median is the average of the two middle values.

Formula :-

(a) For an odd number of values.

$$\text{Median} = \left(\frac{n+1}{2} \right)^{\text{th}} \text{ value}$$

(b) For an even number of values:

$$\text{Median} = \left(\frac{n}{2} \right)^{\text{th}} \text{ value} + \left(\frac{n}{2} + 1 \right)^{\text{th}} \text{ value}$$

* 'n' if represent the total number of values in the set.

Example :-

* In the set [5, 10, 15, 20, 25].

The ~~middle~~ value 15 is the median.

* In the set [5, 10, 15, 20]
The median is -

$$\begin{aligned}\text{Median} &= \frac{10+15}{2} \\ &= 12.5\end{aligned}$$

(3) MODE :-

In Mode the value that appears most frequently in a set of numbers.

Example :-

a) In the set $[5, 10, 15, 20, 10, 5]$

The mode is 5 & 10, why because both appear twice, which is more than any other number.

b) In the set $[5, 10, 15, 20]$, there is ~~no~~ no mode as each number appears only once.

Summary :-

① Mean :-

The average of a set of numbers

② Median :-

The middle value in a sorted set
of numbers.

③ mode :-

The value that appears most
frequently in a set of numbers.

Question :- (B) DISPERSION

Q! Find the standard deviation & variance of the following data:

[5, 9, 8, 12, 6, 10, 6, 8]

(1) Variance :-

Variance is the average of the squared differences from the mean.

Formula :-

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

where,

- * σ^2 It represents the variance
- * x_i It represents each value in the dataset.
- * \bar{x} It represents the mean of the dataset.

* n it represents the total number of values in the dataset.

(2) standard deviation:

$$\text{Standard deviation}(\sigma) = \sqrt{\text{variance}}$$

* Calculation steps:

- (a) calculate the mean (\bar{x}).
- (b) subtract the mean from each value of find the deviation of each value.
- (c) square each deviation.
- (d) find the average of these squared deviations (this is the variance).
- (e) take the square root of the variance to get the standard deviation.

* Example dataset:

$$[5, 9, 8, 12, 6, 10, 6, 8]$$

① Calculate the Mean (\bar{x}):

$$\bar{x} = \frac{5+9+8+12+6+10+6+8}{8} = \frac{64}{8} = 8$$

② Calculate each deviation from the Mean and square it.

$$(5-8)^2 = (-3)^2 = 9$$

$$(9-8)^2 = (1)^2 = 1$$

$$(8-8)^2 = (0)^2 = 0$$

$$(12-8)^2 = (4)^2 = 16$$

$$(6-8)^2 = (-2)^2 = 4$$

$$(10-8)^2 = (2)^2 = 4$$

$$(6-8)^2 = (-2)^2 = 4$$

$$(8-8)^2 = (0)^2 = 0$$

⑥ Sum of squared deviations:

$$9+1+0+16+4+4+4+0 = 38$$

⑦ Calculate the variance:

$$\text{variance}(\sigma^2) = \frac{38}{8}$$

$$= 4.75_{11}$$

⑧ Calculate the standard deviation:

$$\text{Standard Deviation}(\sigma) = \sqrt{4.75} = 2.18_{11}$$

Conclusion

* variance is 4.75

* standard deviation ≈ 2.18

So, for the dataset [5, 9, 8, 12, 6, 10, 6, 8]
the variance is 4.75 and the standard deviation is approximately 2.18₁₁

(B) MEASURES of VARIABILITY (spread)

These measures indicate the spread or dispersion of the dataset.

① Range

② variance

③ standard deviation

(1) Range :

The difference between the highest and lowest values.

(2) variance :

The average of the squared differences from the mean.

(3) standard deviation :

The square root of the variance, indicating how much the values in the dataset deviate from the mean.

Example :-

~~Given~~

Consider the data set $[10, 20, 20, 30, 40]$

(1) Range :

$$\text{Range} = 40 - 10 \\ = 30,$$

(2) Variance :

(i) Calculate the mean : 24

(ii) Then, find the squared deviations
and their average :

$$(10 - 24)^2 = 196$$

$$(20 - 24)^2 = 16$$

$$(20 - 24)^2 = 16$$

$$(30 - 24)^2 = 36$$

$$(40 - 24)^2 = 256$$

(iii) Sum of square deviations :

$$= 196 + 16 + 16 + 36 + 256$$
$$= 520$$

(iv) variance $\sigma^2 = \frac{520}{5}$

$$= 104,$$

(v) standard deviation ?

$$\text{standard deviation} = \sqrt{104} \approx 10.2,$$

* OTHER DESCRIPTIVE STATISTICS

(1) Minimum and Maximum:

The smallest and largest values in the dataset.

(2) Quartiles:

values that divide the dataset into four equal parts.

(3) Interquartile Range (IQR):

The range of the middle 50% of the data, calculated as $Q_3 - Q_1$.

* Descriptive statistics provide a comprehensive summary of the dataset, allowing for a better understanding of the distribution, central value, & spread of the ~~data~~ data.

Example :

Consider the dataset : [10, 20, 20, 30, 40]

* Minimum : 10

* Maximum : 40

* Quartiles :

* Q_1 (25th percentile) : 20

* Q_2 (50th percentile or median) : 20

* Q_3 (75th percentile) : ~~30~~ 30

* Interquartile Range (IQR) :

$$IQR = Q_3 - Q_1 = 30 - 20 \\ = 10$$

Conclusion

* Central Tendency : Mean (24), Median (20), Mode (20)

* Variability : Range (30), variance (104), standard deviation ($\approx \pm 0.2$)

* Other statistics : Minimum (10), Maximum (40), Quartiles (20, 20, 30), IQR 10.

* DESCRIPTIVE STATISTICS

(C) DATA EXPLORATION:

1. pie chart
2. Histogram
3. Bar chart

1. Pie chart:

A pie chart is a circular graph divided into slices to illustrate numerical proportions.

Each slice represents a category's proportion relative to the whole.

Purpose:

* Representation:

Shows how individual parts contribute to a whole.

* Comparison:

Facilitates visual comparison of -

Different categories or parts.

Example:

Suppose we have survey data on favorite fruits.:

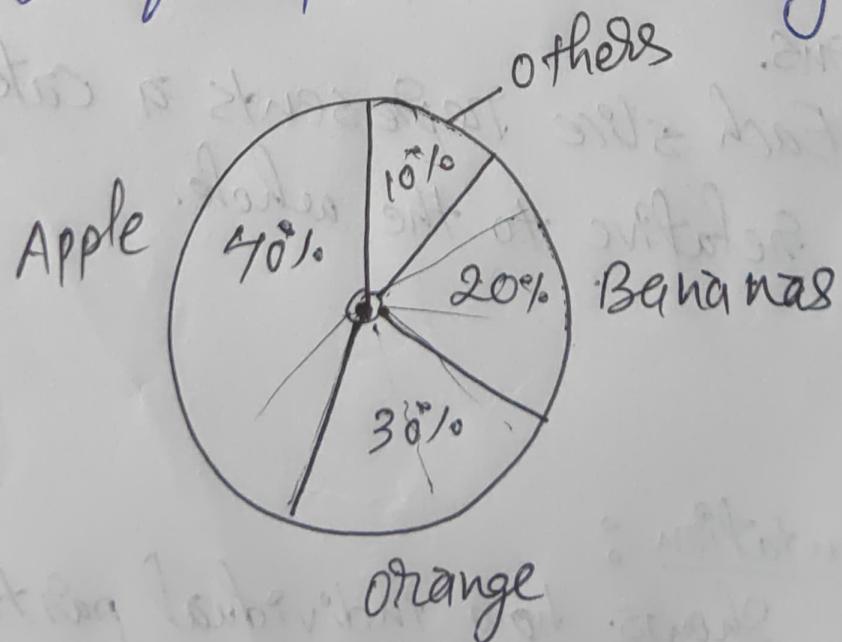
Apples (40%),

Oranges (30%),

Bananas (20%), &

Others (10%).

A pie chart visually displays these proportions, with each slice representing the percentage of respondents choosing each fruit.



2. Histogram:

A histogram is a graphical representation of the distribution of numerical data.

It consists of bars, where each bar represents a range (or bin) of values and the height represents the frequency of data points within that range.

Purpose:

* Distribution:

Shows the shape, spread and center of data distribution.

* Frequency:

visualize how often different values occur within specified intervals (bins).

Example:

* Example:

* Histograms are commonly used in statistics to display data distributions, such as exam scores grouped into bins like 0-10, 11-20 etc.

Each bar's height indicates how many students scored within that range.

3. Bar chart :

A bar chart represents categorical data with rectangular bars. The length or height of each bar is proportional to the frequency or value it represents.

Purpose :

* Comparison :

Compares different categories or groups.

* Frequency :

Shows absolute or relative frequencies of categorical data.

* Example :

* Consider sales data for different months. January (\$1000), February (\$1500), March (\$1200).

A bar chart displays these values as bars of corresponding heights, making it easy to compare sales across months visually.

* CONCLUSION:

- (1) pie chart shows proportions of a whole.
- (2) Histogram displays distribution of numerical data.
- (3) Bar chart compares categorical data with bars.

Usage in descriptive statistics:

* Descriptive statistics.

These charts are essential tools for summarizing and visualizing data distributions, frequency and proportions.

* Insights Generation.

They help analysts and researchers identify patterns, trends, outliers and relationships within datasets.

* Communication.

Effective communication of data insights to stakeholders, clients or decision makers.

STATISTICS

(2) INFERENTIAL STATISTICS

Inferential statistics involves making inferences about populations using data drawn from the population.

Instead of merely describing the data as in descriptive statistics, inferential statistics allow us to make predictions or generalizations about a larger group based on a sample of data.

* Key Concepts in Inferential Statistics

(1) population vs. sample

(2) parameter vs. statistic

(3) Hypothesis Testing

(4) Confidence Intervals

(5) Regression Analysis

(6) T-Tests and ANOVA

(1) population vs. sample.

* Population :-

The entire group that we want to draw conclusions about.

* Sample :-

A subset of the population used to collect data and make inferences about the population.

(2) parameter vs. statistic

* Parameter :-

A numerical characteristic of a population (eg, population mean).

* Statistic :-

A numerical characteristic of a sample (eg, sample mean).

(3) Hypothesis Testing:

* Null hypothesis (H_0):

The hypothesis that there is no effect or no difference.

* Alternative hypothesis (H_a):

The hypothesis that there is an effect or a difference.

* p-value:

The probability of observing the data, or something more extreme, if the null hypothesis is true.

* Significance level (α):

A threshold for determining whether the p-value is low enough to reject the null hypothesis (commonly 0.05).

(4) Confidence Intervals:

* A range of values that is likely to contain the population parameter with a certain level of confidence (e.g., 95%).

(5) Regression Analysis:

A statistical method for examining the relationship between two or more variables.

(6) T-tests and ANOVA:

* T-test: used to compare the means of two groups.

* ANOVA (Analysis of variance):

used to compare the means of three or more groups.

* EXAMPLE OF INFERENCEAL STATISTICS

Q: Scenario :-

A Company wants to know the average amount of time its employees spend on breaks during a workday.
Instead of asking all 1,000 employees, the company surveys a sample of 100 employees.

(1) Sample data collection :-

* Suppose the average break time from the sample of 100 employees is 30 minutes, 30 minutes with a standard deviation of 5 minutes.

(2) Estimate population mean :-

* The company uses the sample mean to estimate the population mean.

(3) Confidence Interval:

* Calculate the 95% confidence interval for the population mean break time.

* Formula :- $CI = \bar{x} \pm Z \left(\frac{\sigma}{\sqrt{n}} \right)$

* where,

* \bar{x} is the sample mean (30 minutes)

* Z is the ~~zero~~ z-score corresponding level (1.96 for 95%).

* σ is the sample standard deviation (5 minutes)

* n is the sample size (100)

* Calculations:

$$CI = 30 \pm 1.96 \left(\frac{5}{\sqrt{100}} \right)$$

$$CI = 30 \pm 1.96 (0.5)$$

$$CI = 30 \pm 0.98$$

$$CI = [29.02, 30.98]$$

* Interpretation :

we are 95% confident that the true mean break time for all employees is between 29.02 and 30.38 minutes.

(4) Hypothesis Testing :

- * Suppose the company has a policy that states the average break time should be 28 minutes.

* Null Hypothesis (H_0): The average break time is 28 minutes.

* Alternative Hypothesis (H_a): The average break time is not 28 minutes.

* Calculate the test statistic (t -score):

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

$$t = \frac{30 - 28}{\frac{s}{\sqrt{100}}}$$

* Compare the t-score to the critical value from the t-distribution table with $df = n - 1 = 99$.

* If the t-score exceeds the critical value, reject the null hypothesis.

(5) Conclusion :

If the p-value corresponding to the t-score is less than 0.05, we reject the null hypothesis and conclude that the average break time is significantly different from 28 minutes.

Summary

- * Inferential statistics allows us to make predictions or generalizations about a population based on a sample.
- * Key techniques include hypothesis testing, confidence intervals, and regression analysis.
- * Example : Estimating the average break time for all employees based on a sample and testing if it differs from a specified value.

~~safe~~

Inferential statistics provide powerful tools to make decisions and predictions about populations based on sample data, allowing for informed decision-making in various fields.