

Week 4 - Descriptive Statistics in Mathematical Modelling

Introduction:

Welcome to this Module on Measures of Central Tendency and Measures of Variability in the context of mathematical modeling.

Descriptive statistics play a critical role in the context of mathematical modeling for several reasons. They provide a set of tools that help in summarizing and understanding the characteristics of a dataset, which is essential for the development, analysis, and validation of mathematical models. Here are some key reasons why descriptive statistics are useful in this context:

Understanding Data Distribution: Descriptive statistics help in understanding the distribution of the data. Measures such as mean, median, and mode give insights into the central tendency of the data, while variance and standard deviation offer information about the data's spread. Knowing the distribution helps in choosing the appropriate mathematical model that fits the data.

Identifying Outliers: Outliers can significantly affect the performance and accuracy of a mathematical model. Descriptive statistics enable the identification of outliers, which can then be analyzed to determine if they are errors or if they represent valuable extremes that should be included in the modeling process.

Data Summarization: In mathematical modeling, it is often necessary to present data in a summarized form to simplify analysis and interpretation. Descriptive statistics provide a concise description of the dataset through measures of central tendency and variability, making it easier to understand the overall characteristics of the data.

Model Assumption Verification: Many mathematical models make assumptions about the underlying distribution of the data (e.g., normality). Descriptive statistics can be used to verify these assumptions through graphical representations like histograms and Q-Q plots, ensuring that the chosen model is appropriate for the data.

Comparing Data Sets: Descriptive statistics are useful for comparing datasets, which is often required in modeling to understand changes over time or differences between groups. This comparison can inform model development and help in evaluating the model's effectiveness across different conditions or populations.

Feature Selection and Engineering: In the process of building a mathematical model, selecting the right features (variables) is crucial. Descriptive statistics can highlight relationships between variables through correlation coefficients and covariance, aiding in the selection of relevant features for the model.

Model Validation and Improvement: After a model has been developed, descriptive statistics are used to validate the model's performance by comparing predicted values against actual values. This process often involves analyzing residuals (the differences between observed and predicted values) to identify patterns that could suggest improvements to the model.

For the purposes of this model, we will be focusing on using descriptive statistics to understand data distributions, identify outliers and to summarize data.

Measures of Central Tendency:

Measures of central tendency help us find the “center” or typical value of a dataset.

1. Mean (Average):

The mean is calculated by adding up all the values in a dataset and dividing by the number of observations. It is sensitive to outliers and can be skewed by extreme values.

The formula for the mean () of a dataset with n observations is:

$$= (\sum x_i) / n$$

2. Median:

The median is the middle value when the data is arranged in ascending or descending order. It is not affected by extreme outliers and is a robust measure of central tendency. For datasets with an even number of observations, the median is the average of the two middle values.

3. Mode:

The mode is the value that appears most frequently in the dataset. Some datasets may have multiple modes (bimodal, trimodal, etc.), while others may have no mode (no value occurs more than once).

Measures of Variability:

Measures of variability quantify how spread out or dispersed the data points are from the central tendency.

1. Range:

The range is the difference between the maximum and minimum values in the dataset. It is simple to calculate but sensitive to outliers.

2. Variance and Standard Deviation:

Variance measures the average squared deviation of each data point from the mean.

Standard deviation is the square root of the variance.

A smaller standard deviation indicates less variability in the data, while a larger one suggests greater variability.

The formulas for variance ($\hat{\sigma}^2$) and standard deviation ($\hat{\sigma}$) are:

$$\begin{aligned}\hat{\sigma}^2 &= \frac{\sum (x_i - \bar{x})^2}{n} \\ \hat{\sigma} &= \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}\end{aligned}$$

3. Interquartile Range (IQR):

IQR is the range of the middle 50% of the data and is calculated as the difference between the third quartile (Q3) and the first quartile (Q1). It is less sensitive to outliers compared to the range.

Applications in Mathematical Modeling:

Measures of central tendency and variability are vital in mathematical modeling for various reasons:

- They help summarize and understand the data used to build models.
- They provide insights into the distribution and characteristics of the data.
- They help identify outliers and understand their impact on the model.
- They guide the selection of appropriate probability distributions for modeling uncertainty.

Conclusion:

Measures of central tendency (mean, median, mode) and variability (range, variance, standard deviation, IQR) are essential statistical tools for summarizing and understanding data in the context of mathematical modeling.

These measures assist in characterizing the central behavior and spread of data, aiding model development, and decision-making.