

Lecture (4 & 5)

understanding statistical description in Data Science.

- How data is spread across different values
- Gain insights from data
- Identify data patterns
- central value

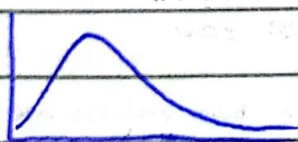
Data Distribution:

- ① symmetric distribution: (Both left and right parts are equal)
- ② Asymmetric distribution or skewed:

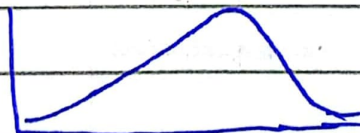
It can be positive and negative skewed

→ positive skew: Data spread at start, curve goes down slowly. Tail on right side

→ Negative skew: curve gradually increases, then abruptly closes. Tail on left side



positive skew



Negative skewed

→ Noise: Random unpredictable errors in the data

→ Bias (means data is leaning in one direction)

* Symmetric errors from incorrect assumptions made by model

→ statistical description tells us about an outliers

→ Through statistical description we get to know how to clean our data.

Main types of statistical Description

- Measure of central tendency
- Measure of dispersion (shows how values vary from each other)
- Measure of similarity

Measure of Central tendency

→ central tendency:

It's the value of feature/attribute that divide data into two halves

- Typical value indicates the most common or average value in the distribution (central value)

→ Measures of central tendency:

Mean (simple average) $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$ (con) Its sensitive to extreme values

Median (middle of sorted data) cons: computationally cost

Mode (most frequent)

Trimmed Mean (man 2%) (before calculation) It's modified mean where extreme values are removed

Midrange ($\frac{\min + \max}{2}$) ∴ easy to calculate → use for brief estimate

- If mean and median values are closer to each other then we can say it's symmetrically distributed.

- Mean, Median, mode also tells us whether the data is symmetric or skewed

median = mean = mode (symmetric Distribution)

mode < median < mean (positive skew)

mean < median < mode (Negative skew)