



RAMA
UNIVERSITY

www.ramauniversity.ac.in

FACULTY OF ENGINEERING

DATA MINING & WAREHOUSEING
LECTURE-38

MR. DHIRENDRA

ASSISTANT PROFESSOR

RAMA UNIVERSITY

OUTLINE

- ❖ MEASURING THE CENTRAL TENDENCY
- ❖ MEASURING THE DISPERSION OF DATA
- ❖ BOXPLOT ANALYSIS
- ❖ BOXPLOT
- ❖ VISUALIZATION OF DATA DISPERSION: BOXPLOT ANALYSIS
- ❖ MCQ
- ❖ REFERENCES



Measuring the Central Tendency

- Mean $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
 - Weighted arithmetic mean

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

- Median: A holistic measure

- Middle value if odd number of values, or average of the middle two values otherwise
- estimated by interpolation

$$median = L_1 + \left(\frac{n/2 - (\sum f)l}{f_{median}} \right) c$$

- Mode

- Value that occurs most frequently in the data
- Unimodal, bimodal, trimodal
- Empirical formula:

$$mean - mode = 3 \times (mean - median)$$

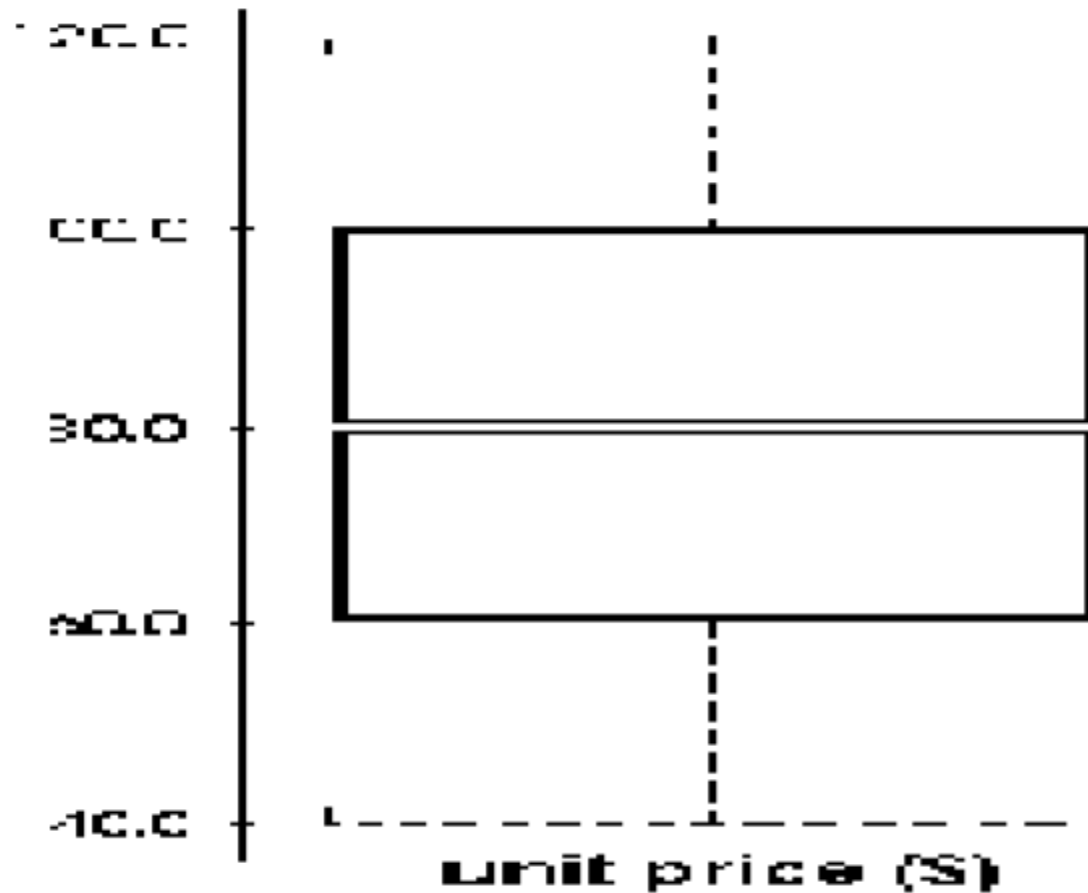
Measuring the Dispersion of Data

- Quartiles, outliers and boxplots
 - Quartiles: Q_1 (25th percentile), Q_3 (75th percentile)
 - Inter-quartile range: $IQR = Q_3 - Q_1$
 - Five number summary: min, Q_1 , M, Q_3 , max
 - Boxplot: ends of the box are the quartiles, median is marked, whiskers, and plot outlier individually
 - Outlier: usually, a value higher/lower than $1.5 \times IQR$
- Variance and standard deviation
 - Variance s^2 : (algebraic, scalable computation)
 - $$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right]$$
 - Standard deviation s is the square root of variance s^2

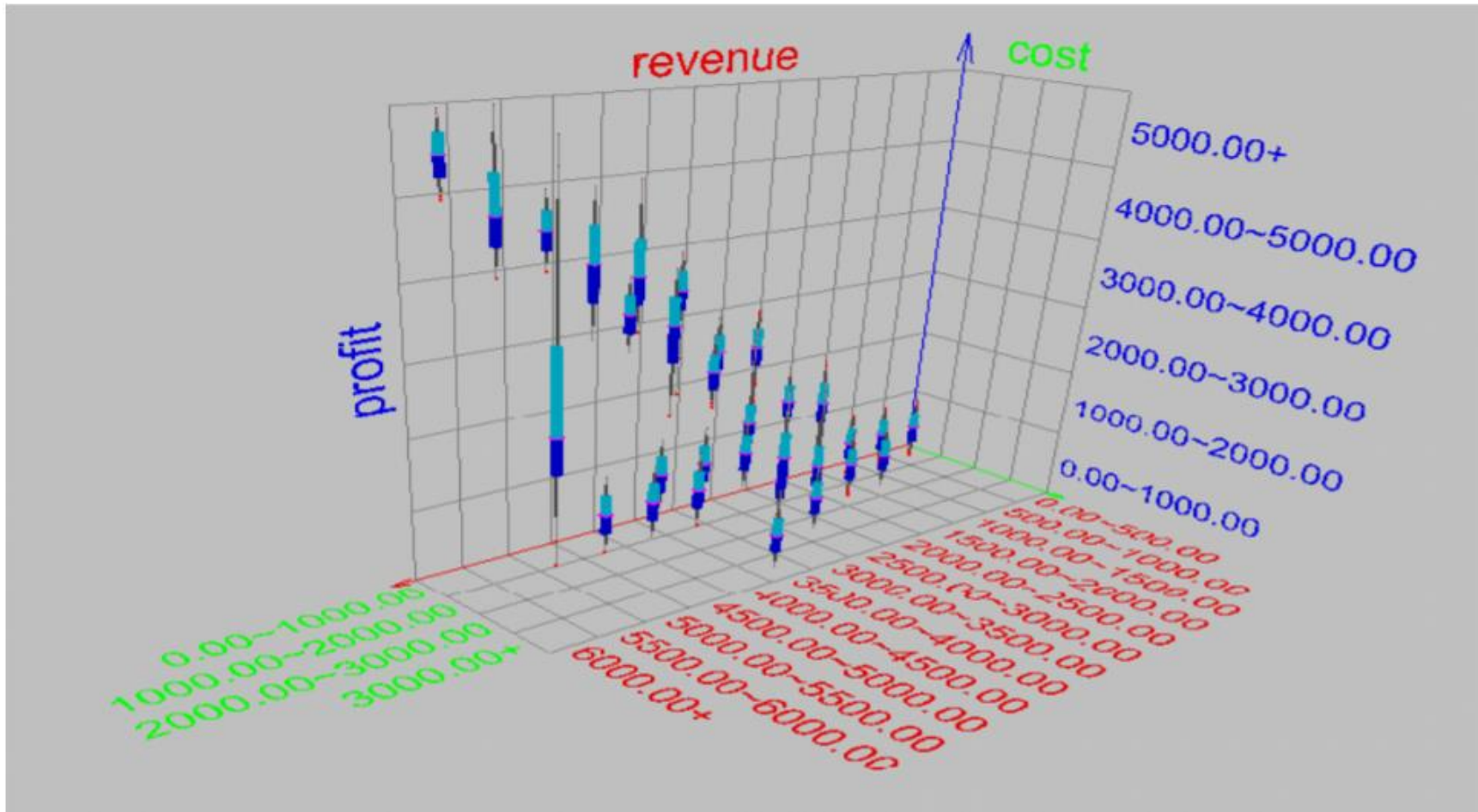
- Five-number summary of a distribution:
 - Minimum, Q1, M, Q3, Maximum
- Boxplot
 - Data is represented with a box
 - The ends of the box are at the first and third quartiles, i.e., the height of the box is IRQ
 - The median is marked by a line within the box
 - Whiskers: two lines outside the box extend to Minimum and Maximum

Boxplot

A boxplot



Visualization of Data Dispersion: Boxplot Analysis



Multiple Choice Question

1. Various visualization techniques are used in _____ step of KDD.

- a) selection
- b) transformaion
- c) data mining.
- d) interpretation.

2. Extreme values that occur infrequently are called as _____.

- a) outliers
- b) rare values.
- c) dimensionality reduction.
- d) All of the above.

3. Box plot and scatter diagram techniques are _____.

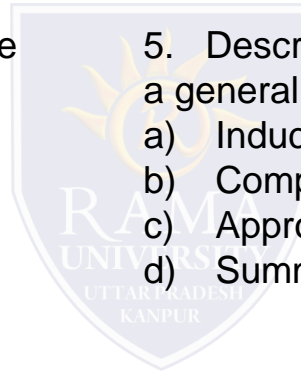
- a) Graphical
- b) Geometric
- c) Icon-based.
- d) Pixel-based.

4. _____ is used to proceed from very specific knowledge to more general information.

- a) Induction
- b) Compression.
- c) Approximation.
- d) Substitution.

5. Describing some characteristics of a set of data by a general model is viewed as _____

- a) Induction
- b) Compression
- c) Approximation
- d) Summarization



REFERENCES

- https://www.tutorialspoint.com/dwh/dwh_overview.htm
- <https://www.geeksforgeeks.org/>
- <http://myweb.sabanciuniv.edu/rdekharghani/files/2016/02/The-Morgan-Kaufmann-Series-in-Data-Management-Systems-Jiawei-Han-Micheline-Kamber-Jian-Pei-Data-Mining.-Concepts-and-Techniques-3rd-Edition-Morgan-Kaufmann-2011.pdf> DATA MINING BOOK WRITTEN BY Micheline Kamber
- <https://www.javatpoint.com/three-tier-data-warehouse-architecture>
- M.H. Dunham, “ Data Mining: Introductory & Advanced Topics” Pearson Education
- Jiawei Han, Micheline Kamber, “ Data Mining Concepts & Techniques” Elsevier
- Sam Anahory, Dennis Murray,” data warehousing in the Real World: A Practical Guide for Building Decision Support Systems, “ Pearson Education
- Mallach,” Data Warehousing System”, TMH
- R. Agrawal, A. Gupta, and S. Sarawagi. Modeling multidimensional databases. ICDE’97 S. Chaudhuri and U. Dayal. An overview of data warehousing and OLAP technology. ACM SIGMOD Record, 26:65-74, 1997
- S. Agarwal, R. Agrawal, P. M. Deshpande, A. Gupta, J. F. Naughton, R. Ramakrishnan, and S. Sarawagi. On the computation of multidimensional aggregates. VLDB’96 D. Agrawal, A. E. Abbadi, A. Singh, and T. Yurek. Efficient view maintenance in data warehouses. SIGMOD’97
- E. F. Codd, S. B. Codd, and C. T. Salley. Beyond decision support. Computer World, 27, July 1993.
- J. Gray, et al. Data cube: A relational aggregation operator generalizing group-by, cross-tab and sub-totals. Data Mining and Knowledge Discovery, 1:29-54, 1997.