



FACULTY OF ENGINEERING

DATA MINING & WAREHOUSEING
LECTURE-39

MR. DHIRENDRA
ASSISTANT PROFESSOR
RAMA UNIVERSITY

OUTLINE

- ❖ MEASURING THE CENTRAL TENDENCY
- ❖ HISTOGRAM ANALYSIS
- ❖ QUANTILE PLOT
- ❖ QUANTILE-QUANTILE (Q-Q) PLOT
- ❖ SCATTER PLOT
- ❖ MCQ
- ❖ REFERENCES



Measuring the Central Tendency

Variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left[\sum x_i^2 - \frac{1}{n} (\sum x_i)^2 \right]$$

- Standard deviation: the square root of the variance
 - Measures spread about the mean
 - It is zero if and only if all the values are equal
 - Both the deviation and the variance are algebraic

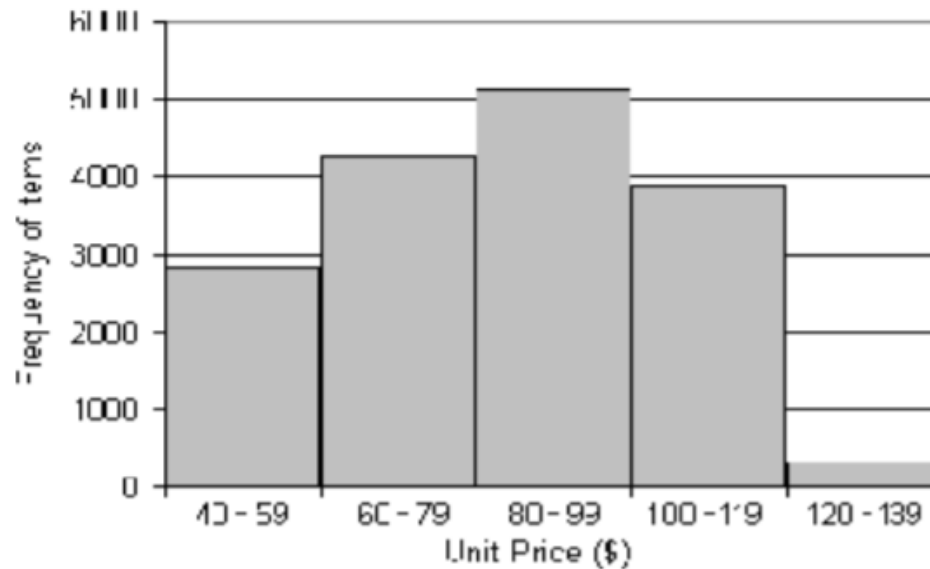
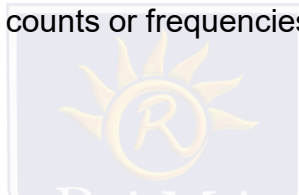


Histogram Analysis

Graph displays of basic statistical class descriptions

– Frequency histograms

- A univariate graphical method
- Consists of a set of rectangles that reflect the counts or frequencies of the classes present in the given data

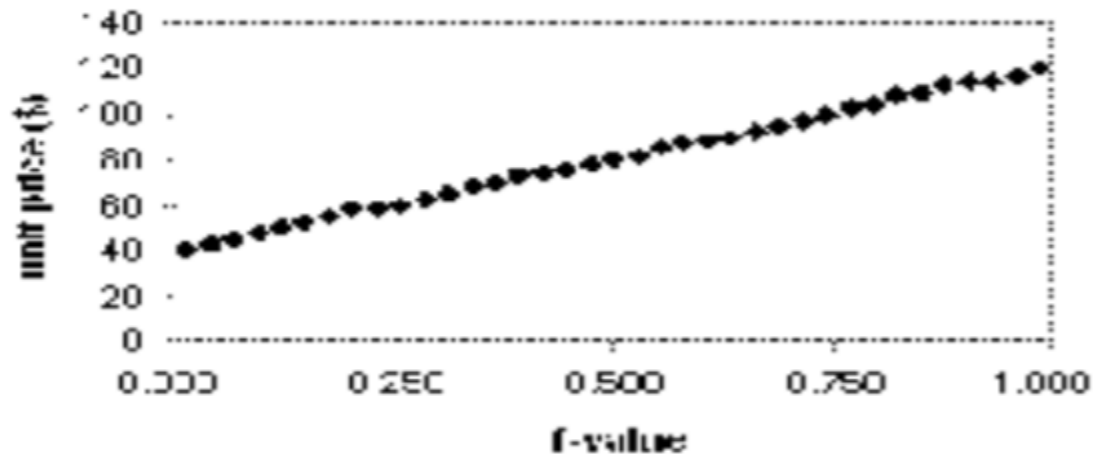
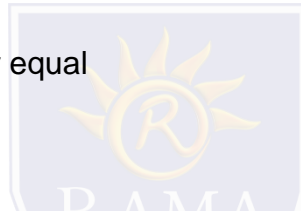


Quantile Plot

Displays all of the data (allowing the user to assess both the overall behavior and unusual occurrences)

- Plots quantile information

– For a data x_i data sorted in increasing order, f_i indicates that approximately 100 $f_i\%$ of the data are below or equal to the value x_i

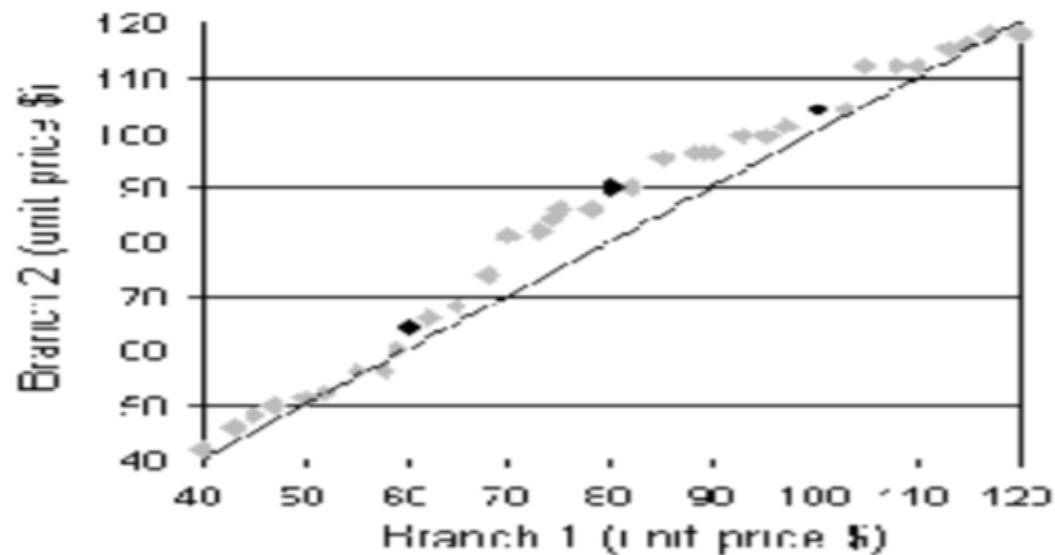


Quantile-Quantile (Q-Q) Plot

Graphs the quantiles of one univariate distribution

against the corresponding quantiles of another

- Allows the user to view where there is a shift in going from one distribution to another



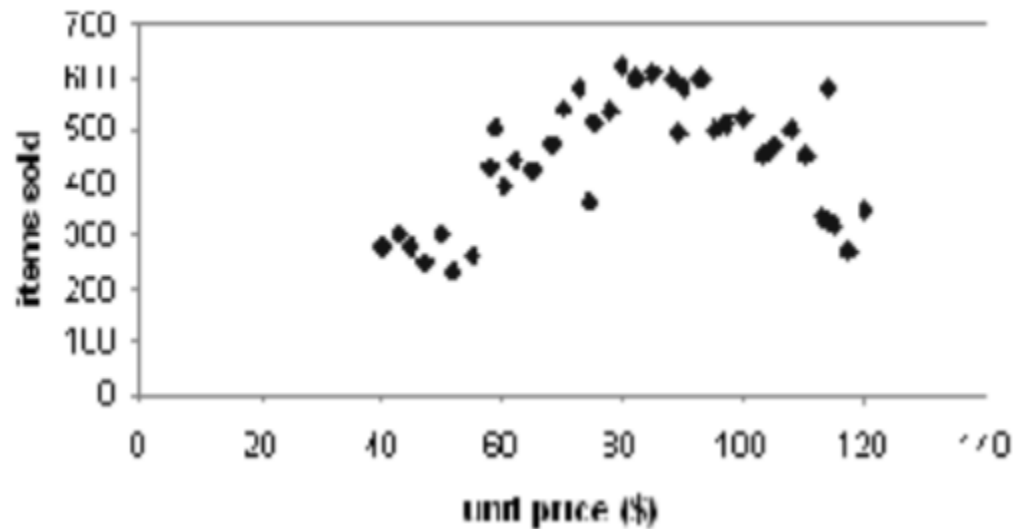
Scatter plot

Provides a first look at bivariate data to see clusters

of points, outliers, etc

- Each pair of values is treated as a pair of coordinates

and plotted as points in the plane



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 _____.

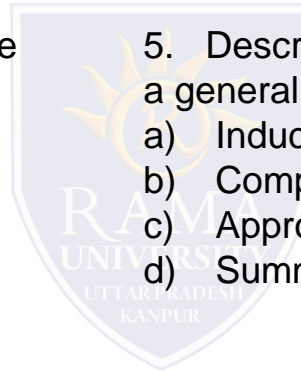
- a) Graphical
- b) Geometric
- c) Icon-based.
- d) Pixel-based.

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/rdehkharghani/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, Denniss 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.