# Data Mining:
## Concepts and Techniques

— Chapter 2 —

# Chapter 2: Data Preprocessing

- General data characteristics

- Basic data description and exploration

- Measuring data similarity

- Data cleaning

- Data integration and transformation

- Data reduction

- Summary

# Types of Data Sets

- Record
  - Relational records
  - Data matrix, e.g., numerical matrix, crosstabs
  - Document data: text documents: term-frequency vector
  - Transaction data
- Graph
  - World Wide Web
  - Social or information networks
  - Molecular Structures
- Ordered
  - Spatial data: maps
  - Temporal data: time-series
  - Sequential Data: transaction sequences
  - Genetic sequence data

| | team | coach | play | ball | score | game | win | lost | timeout | season |
|---|---|---|---|---|---|---|---|---|---|---|
| Document 1 | 3 | 0 | 5 | 0 | 2 | 6 | 0 | 2 | 0 | 2 |
| Document 2 | 0 | 7 | 0 | 2 | 1 | 0 | 0 | 3 | 0 | 0 |
| Document 3 | 0 | 1 | 0 | 0 | 1 | 2 | 2 | 0 | 3 | 0 |

| TID | Items |
|---|---|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |

# Important Characteristics of Structured Data

- Dimensionality
  - Curse of dimensionality
- Sparsity
  - Only presence counts
- Resolution
  - Patterns depend on the scale
- Similarity
  - Distance measure

# Types of Attribute Values

- ## Nominal
  - E.g., profession, ID numbers, eye color, zip codes
- ## Ordinal
  - E.g., rankings (e.g., army, professions), grades, height in {tall, medium, short}
- ## Binary
  - E.g., medical test (positive vs. negative)
- ## Interval
  - E.g., calendar dates, body temperatures
- ## Ratio
  - E.g., temperature in Kelvin, length, time, counts

# Discrete vs. Continuous Attributes

- Discrete Attribute
  - Has only a finite or countably infinite set of values
  - E.g., zip codes, profession, or the set of words in a collection of documents
  - Sometimes, represented as integer variables
  - Note: Binary attributes are a special case of discrete attributes
- Continuous Attribute
  - Has real numbers as attribute values
  - Examples: temperature, height, or weight
  - Practically, real values can only be measured and represented using a finite number of digits
  - Continuous attributes are typically represented as floating-point variables

# Chapter 2: Data Preprocessing

- General data characteristics

- Basic data description and exploration

- Measuring data similarity

- Data cleaning

- Data integration and transformation

- Data reduction

- Summary

# Mining Data Descriptive Characteristics

- **Motivation**
  - To better understand the data: central tendency, variation and spread
- **Data dispersion characteristics**
  - median, max, min, quantiles, outliers, variance, etc.
- **Numerical dimensions** correspond to sorted intervals
  - Data dispersion: analyzed with multiple granularities of precision
  - Boxplot or quantile analysis on sorted intervals
- **Dispersion analysis on computed measures**
  - Folding measures into numerical dimensions
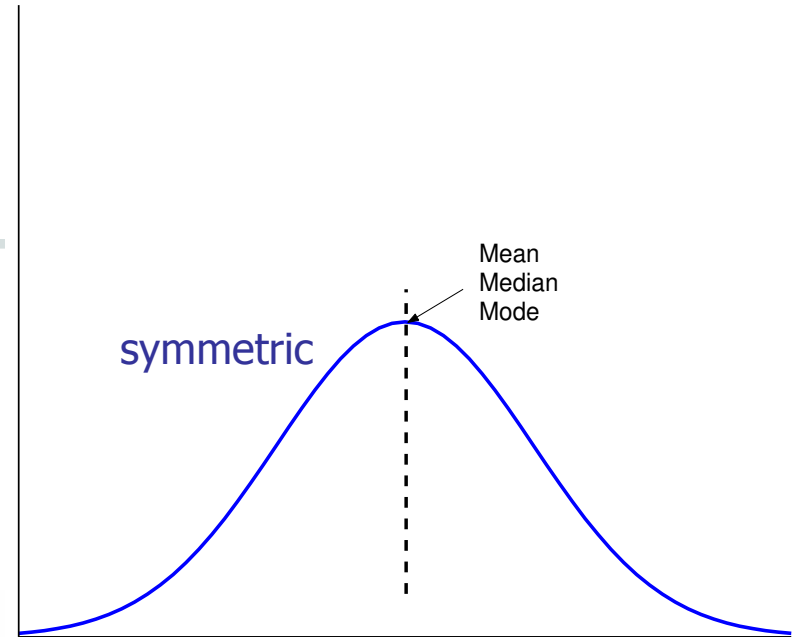  - Boxplot or quantile analysis on the transformed cube

# Measuring the Central Tendency

- Mean (algebraic measure) (sample vs. population): $\bar{x} = \frac{1}{n}\sum_{i=1}^{n} x_i$    $\mu = \frac{\sum x}{N}$
  - Weighted arithmetic mean:
  - Trimmed mean: chopping extreme values    $\bar{x} = \frac{\sum_{i=1}^{n} w_i x_i}{\sum_{i=1}^{n} w_i}$
- Median: A holistic measure
  - Middle value if odd number of values, or average of the middle two values otherwise
  - Estimated by interpolation (for *grouped data*):
- Mode    $median = L_1 + (\frac{N/2 - (\sum freq)l}{freq_{median}})width$
  - Value that occurs most frequently in the data
  - Unimodal, bimodal, trimodal
  - Empirical formula:    $mean - mode = 3 \times (mean - median)$

# Symmetric vs. Skewed Data

- Median, mean and mode of symmetric, positively and negatively skewed data



symmetric

Mean
Median
Mode



Mode    Mean

Median

positively skewed



Mean    Mode

Median

negatively skewed

# 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 x IQR
- Variance and standard deviation (*sample: s, population: σ*)
    - Variance: (algebraic, scalable computation)

$$s^2 = \frac{1}{n-1}\sum_{i=1}^{n}(x_i - \bar{x})^2 = \frac{1}{n-1}[\sum_{i=1}^{n}x_i^2 - \frac{1}{n}(\sum_{i=1}^{n}x_i)^2] \qquad \sigma^2 = \frac{1}{N}\sum_{i=1}^{n}(x_i - \mu)^2 = \frac{1}{N}\sum_{i=1}^{n}x_i^2 - \mu^2$$

    - Standard deviation *s (or σ)* is the square root of variance *$s^2$ (or $\sigma^2$)*
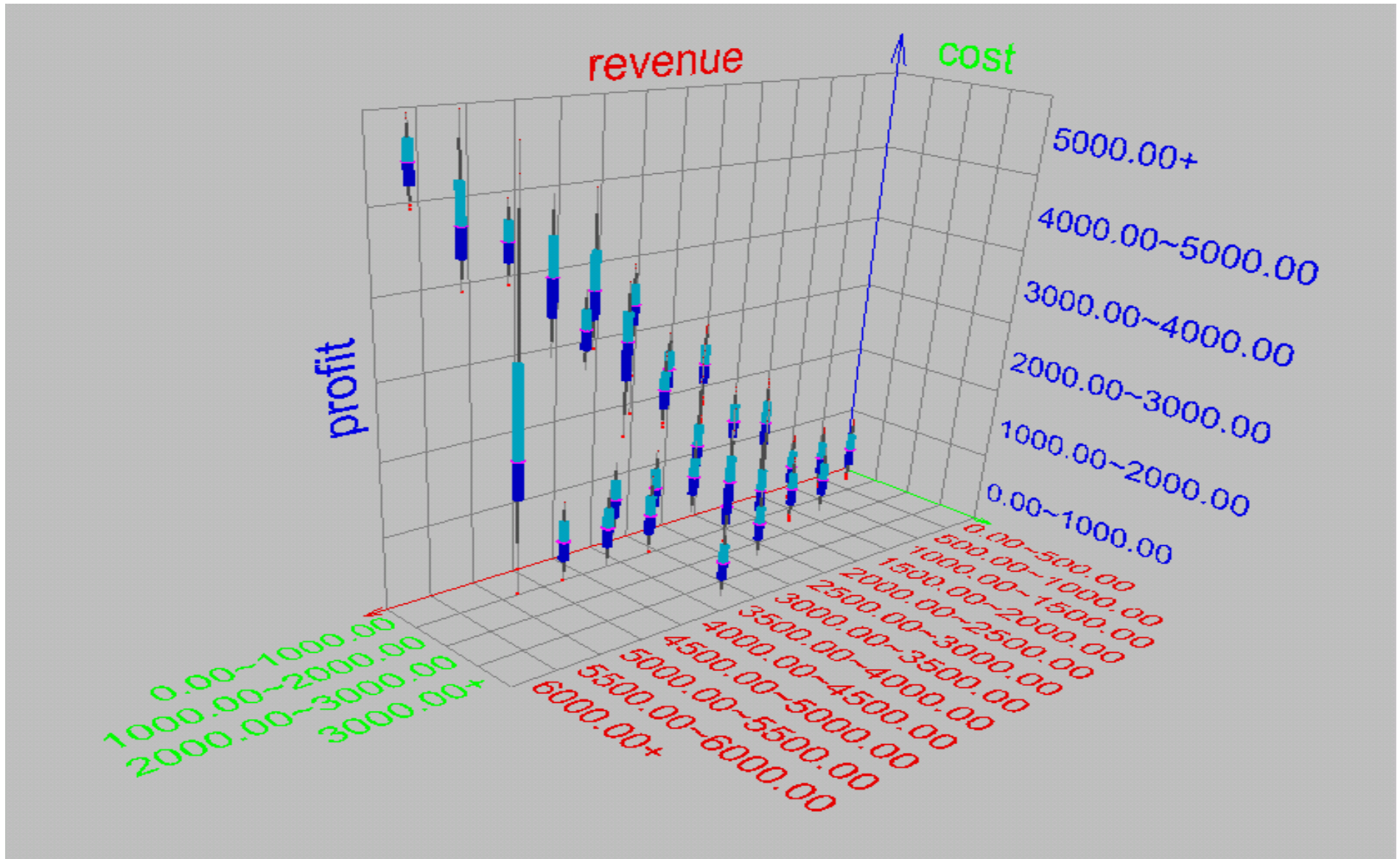
# Boxplot Analysis


unit price (S)

- **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 IQR

    - The median is marked by a line within the box

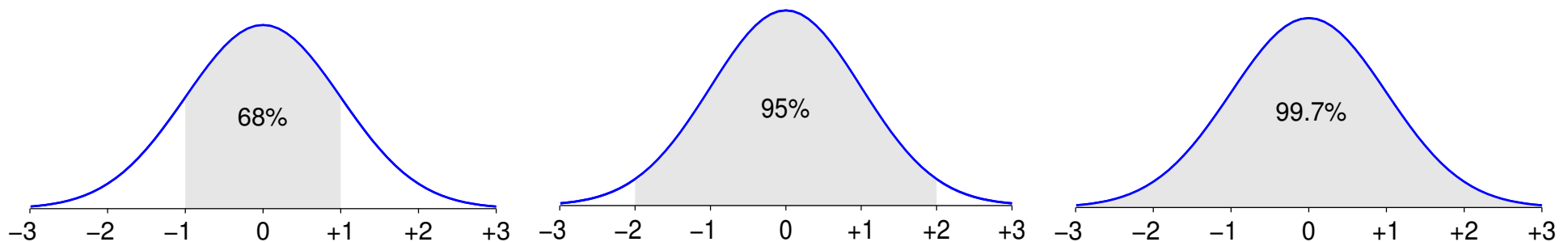    - Whiskers: two lines outside the box extend to Minimum and Maximum

# Visualization of Data Dispersion: 3-D Boxplots

# Properties of Normal Distribution Curve

- The normal (distribution) curve
  - From $\mu-\sigma$ to $\mu+\sigma$: contains about 68% of the measurements ($\mu$: mean, $\sigma$: standard deviation)
  - From $\mu-2\sigma$ to $\mu+2\sigma$: contains about 95% of it
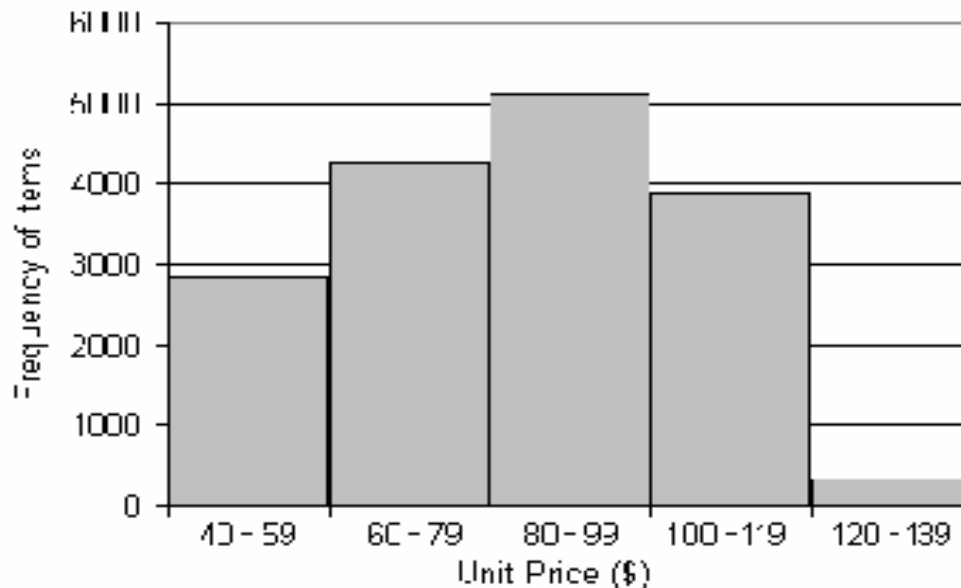  - From $\mu-3\sigma$ to $\mu+3\sigma$: contains about 99.7% of it

# Graphic Displays of Basic Statistical Descriptions

- Boxplot: graphic display of five-number summary
- Histogram: x-axis are values, y-axis repres. frequencies
- Quantile plot: each value $x_i$ is paired with $f_i$ indicating that approximately 100 $f_i$% of data are $\leq x_i$
- Quantile-quantile (q-q) plot: graphs the quantiles of one univariant distribution against the corresponding quantiles of another
- Scatter plot: each pair of values is a pair of coordinates and plotted as points in the plane
- Loess (local regression) curve: add a smooth curve to a scatter plot to provide better perception of the pattern of dependence
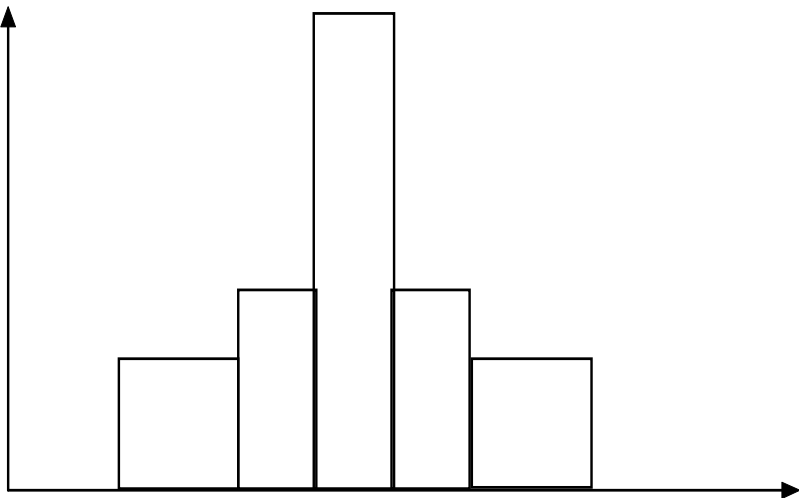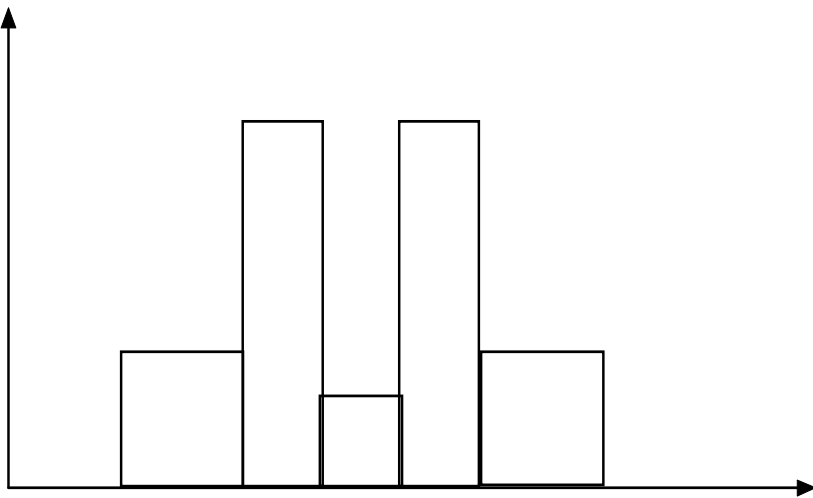
# Histogram Analysis

- 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
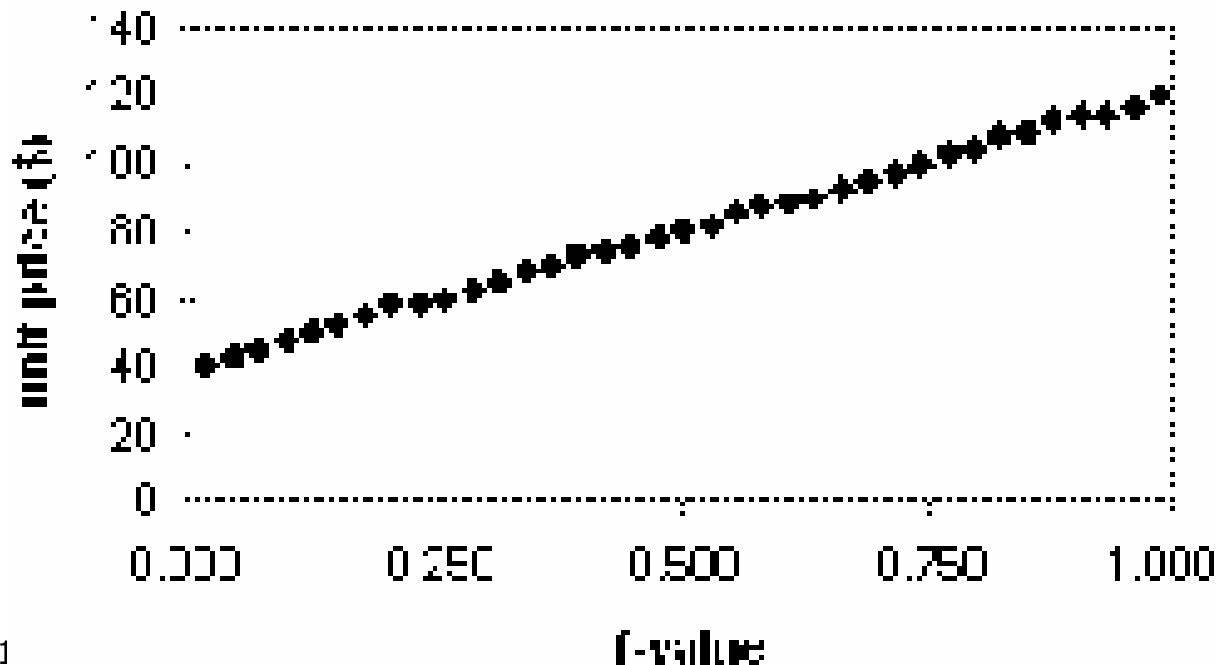
# Histograms Often Tells More than Boxplots



- The two histograms shown in the left may have the same boxplot representation
  - The same values for: min, Q1, median, Q3, max
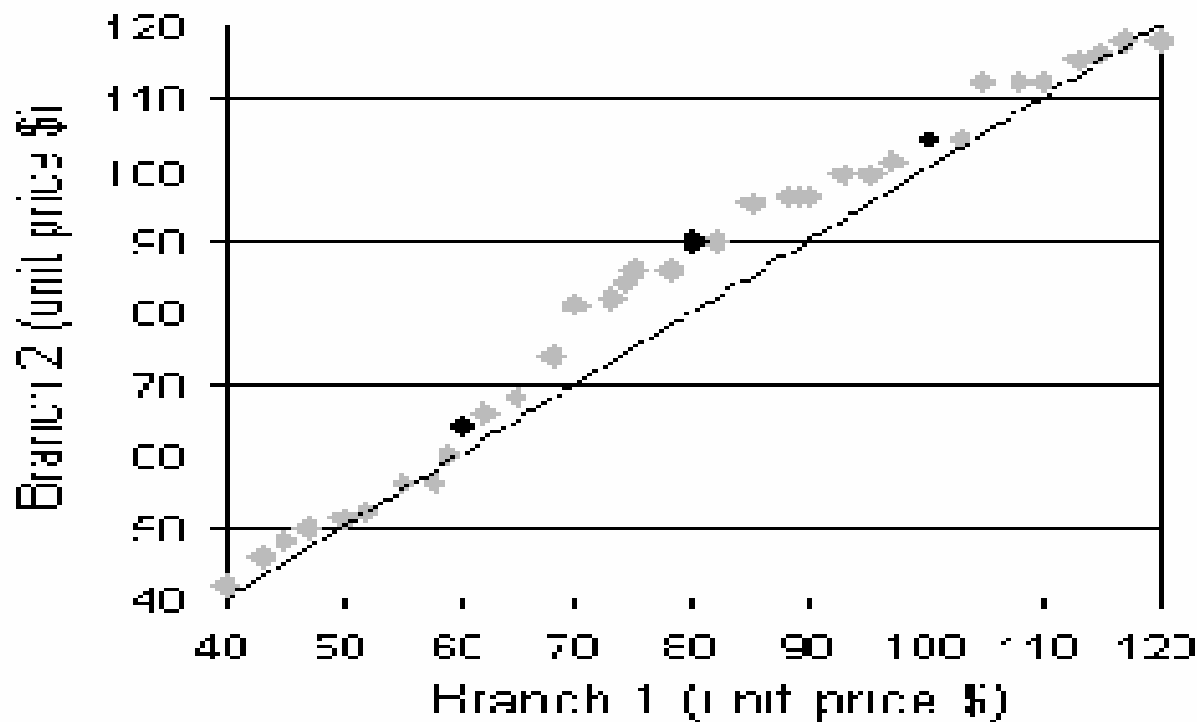- But they have rather different data distributions

# 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 $x_i$
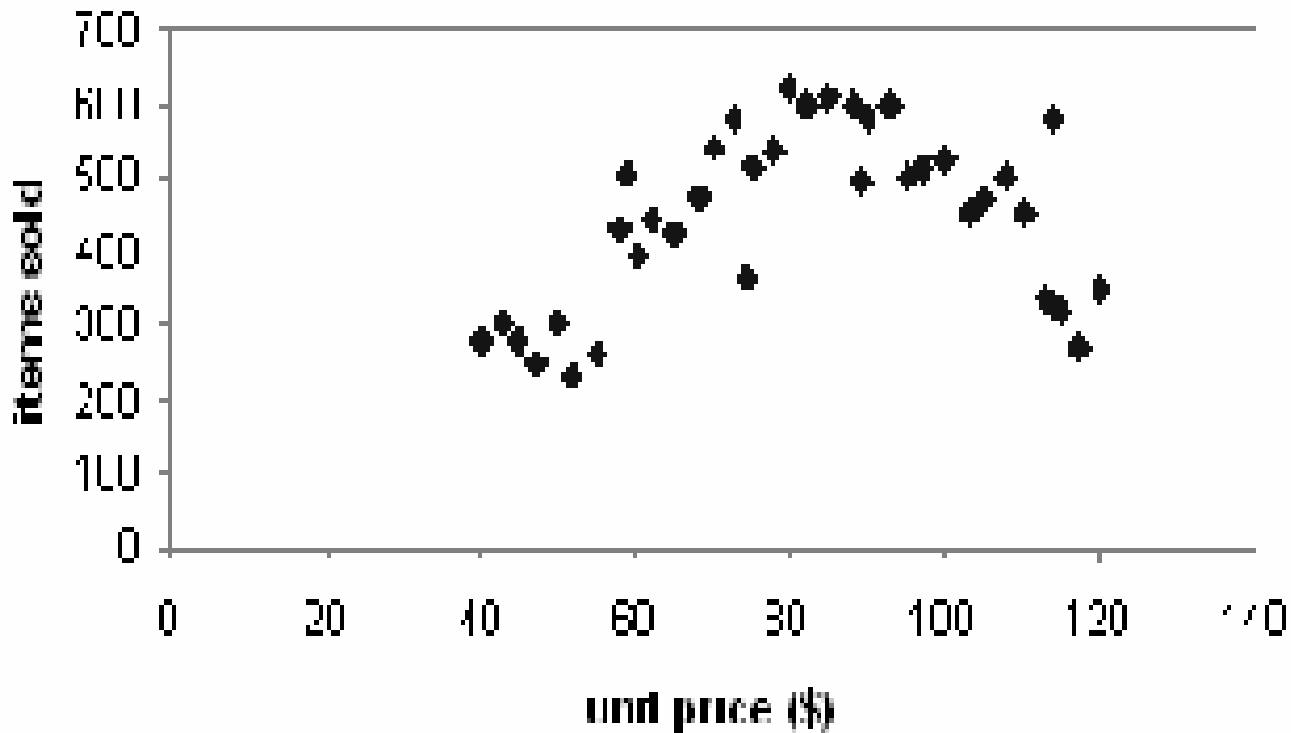
# Quantile-Quantile (Q-Q) Plot

- Graphs the quantiles of one univariate distribution against the corresponding quantiles of another

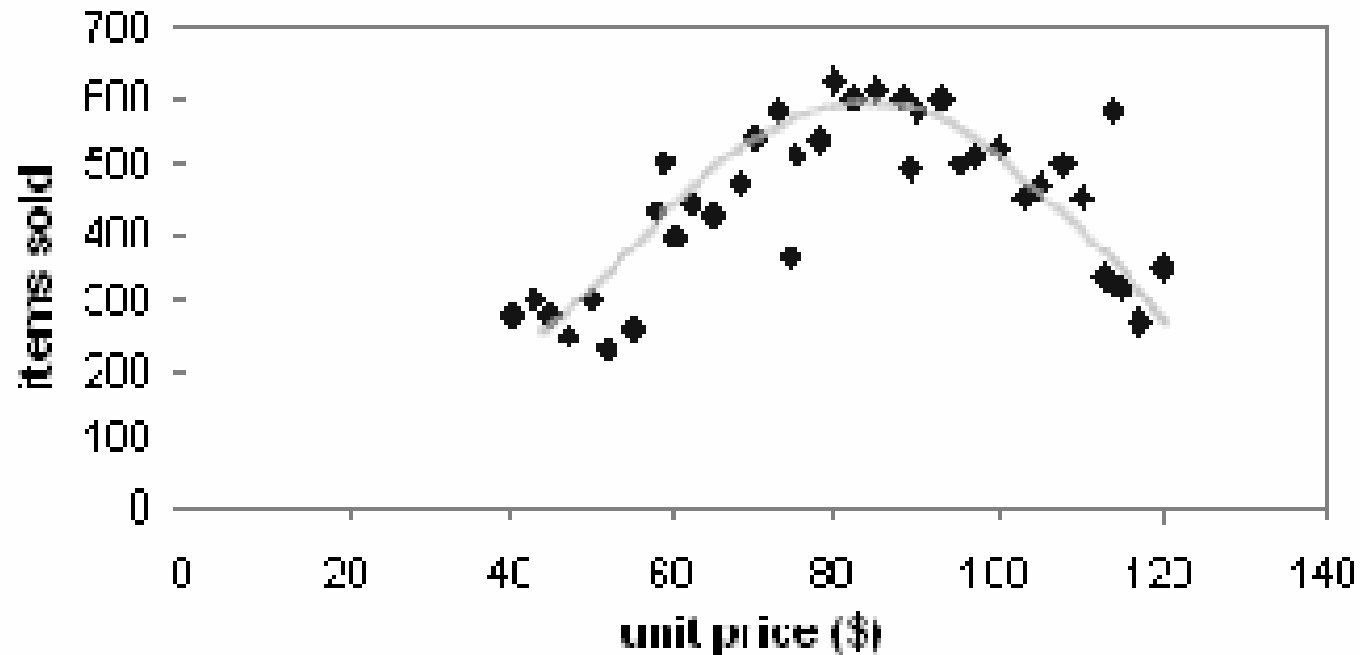- Allows the user to view whether there is a shift in going from one distribution to another

# Scatter plot

- Provides a first look at bivariate data to see clusters of points, outliers, etc
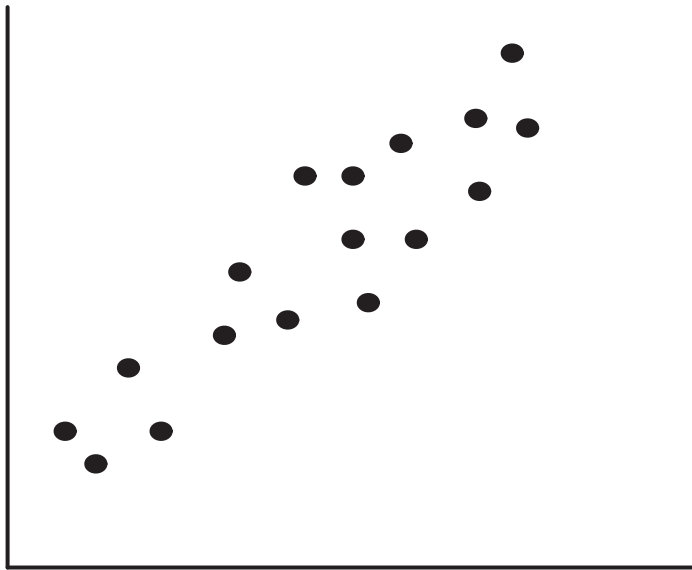- Each pair of values is treated as a pair of coordinates and plotted as points in the plane

# Loess Curve

- Adds a smooth curve to a scatter plot in order to provide better perception of the pattern of dependence
- Loess curve is fitted by setting two parameters: a smoothing parameter, and the degree of the polynomials that are fitted by the regression

# Positively and Negatively Correlated Data



- The left half fragment is positively correlated
- The right half is negative correlated

```
ERROR: invalidrestore
OFFENDING COMMAND: restore

STACK:

-savelevel-
-savelevel-
```