Università degli Studi di Milano

Master Degree in Computer Science

# Information Management course

Teacher: Alberto Ceselli
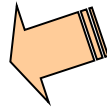
Lecture 02: 07/10/2015

# Data Mining:

## Concepts and Techniques

**— Chapter 2 —**

Jiawei Han, Micheline Kamber, and Jian Pei

University of Illinois at Urbana-Champaign

Simon Fraser University

# Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types

- Basic Statistical Descriptions of Data

- Measuring Data Similarity and Dissimilarity

- 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 and network
  - World Wide Web
  - Social or information networks
  - Molecular Structures
- Ordered
  - Video data: sequence of images
  - Temporal data: time-series
  - Sequential Data: transaction sequences
  - Genetic sequence data
- Spatial, image and multimedia:
  - Spatial data: maps
  - Image data: .bmp
  - Video data: .avi

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 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
    (the volume of the space grows fast with the number of dimensions, and the available data becomes sparse)

- Sparsity
  - Only presence counts

- Resolution
  - Patterns depend on the scale

- Distribution
  - Centrality and dispersion

# Data Objects

- Data sets are made up of data objects.

- A **data object** represents an entity (also called *samples , examples, instances, data points, objects, tuples ...)*

- Examples:
  - sales database:  customers, store items, sales
  - medical database: patients, treatments
  - university database: students, professors, courses

- Data objects are described by **attributes** (also called *variables, dimensions, features ...*)


- In databases: rows -> data objects; columns ->attributes.

# Attributes

- **Attribute (**or **dimensions, features, variables**): a data field, representing a characteristic or feature of a data object.
  - *E.g., customer _ID, name, address*
- Types:
  - Nominal
  - Binary
  - Ordinal
  - Numeric: quantitative
    - Interval-scaled
    - Ratio-scaled

# Attribute Types

- **Nominal:** categories, states, or "names of things"
  - *Hair_color = {auburn, black, blond, brown, grey, red, white}*
  - marital status, occupation, ID numbers, zip codes
- **Binary**
  - Nominal attribute with only 2 states (0 and 1)
  - <u>Symmetric binary</u>: both outcomes equally important
    - e.g., gender
  - <u>Asymmetric binary</u>: outcomes not equally important.
    - e.g., medical test (positive vs. negative)
    - Convention: assign 1 to most important outcome (e.g., HIV positive)
- **Ordinal**
  - Values have a meaningful order (ranking) but magnitude between successive values is not known.
  - *Size = {small, medium, large},* grades, army rankings

# **Numeric Attribute Types**

- Quantity (integer or real-valued)
- **Interval**
  - Measured on a scale of **equal-sized units**
  - Values have order
    - E.g., *temperature in C˚or F˚, calendar dates*
  - No true zero-point
- **Ratio**
  - Inherent **zero-point**
  - We can speak of values as being an order of magnitude larger than the unit of measurement (10 K˚ is twice as high as 5 K˚).
    - e.g., *temperature in Kelvin, length, counts, monetary quantities*

# Discrete vs. Continuous Attributes (ML view)

- **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
    - E.g., 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: Getting to Know Your Data

- Data Objects and Attribute Types

- Basic Statistical Descriptions of Data

- Measuring Data Similarity and Dissimilarity

- Summary

# Basic Statistical Descriptions of Data

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

# Measuring the Central Tendency

- Mean (algebraic measure) (sample vs. population):

  Note: $n$ is sample size and $N$ is population size.

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

  - Weighted arithmetic mean

  - Sensitive to outliers: trimmed mean (chopping extreme values)

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

- Median:

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

  - Estimated by interpolation (for *grouped data*):

| age | frequency |
| --- | --- |
| 1–5 | 200 |
| 6–15 | 450 |
| 16–20 | 300 |
| 21–50 | 1500 |
| 51–80 | 700 |
| 81–110 | 44 |

$$median = L_1 + \left( \frac{\frac{n}{2} - (\sum freq)_l}{freq_{median}} \right) width$$

Lower boundary of the median interval

Sum of freq. of intervals preceding the median

# values in the dataset

Freq. of the median interval

13

# Measuring the Central Tendency

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

  $$mean - mode \simeq 3 \times (mean - median)$$
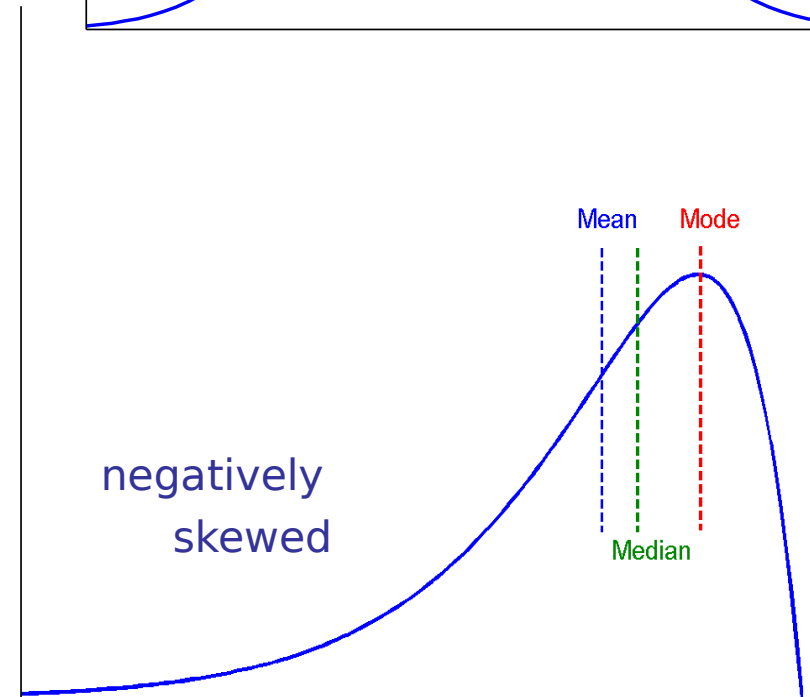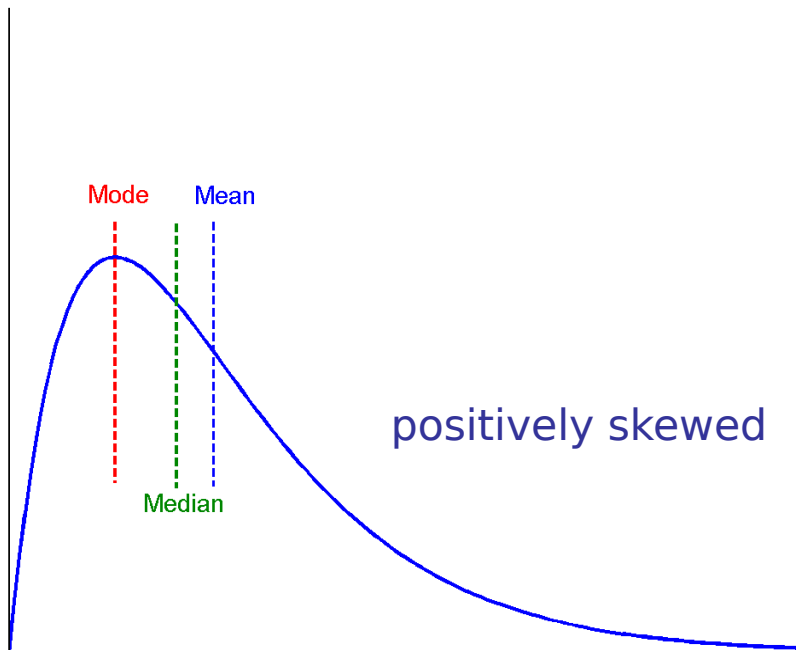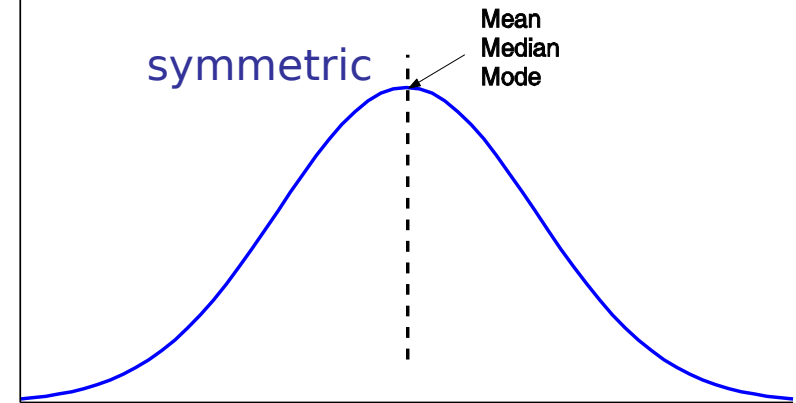
  Mean: 58

  Median: (52+56)/2 = 54

  Mode: 52 and 70 (bimodal)

  Midrange: (30+110) /2 = 70

| Employed | Salary |
|---|---|
| 1 | 30 |
| 2 | 36 |
| 3 | 47 |
| 4 | 50 |
| 5 | 52 |
| 6 | 52 |
| 7 | 56 |
| 8 | 60 |
| 9 | 63 |
| 10 | 70 |
| 11 | 70 |
| 12 | 110 |

# Symmetric vs. Skewed Data

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

symmetric

Mean
Median
Mode

positively skewed

Mode  Mean

Median

negatively skewed

Mean   Mode

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, $Q_1$, median, $Q_3$, max (nice for skewed distributions)

    - **Boxplot**: ends of the box are the quartiles; median is marked; add whiskers, and plot outliers individually

    - **Outlier**: usually, a value higher/lower than 1.5 x IQR

- Variance and standard deviation (*sample: s, population: σ*)

    - **Variance**: (algebraic, scalable computation)

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

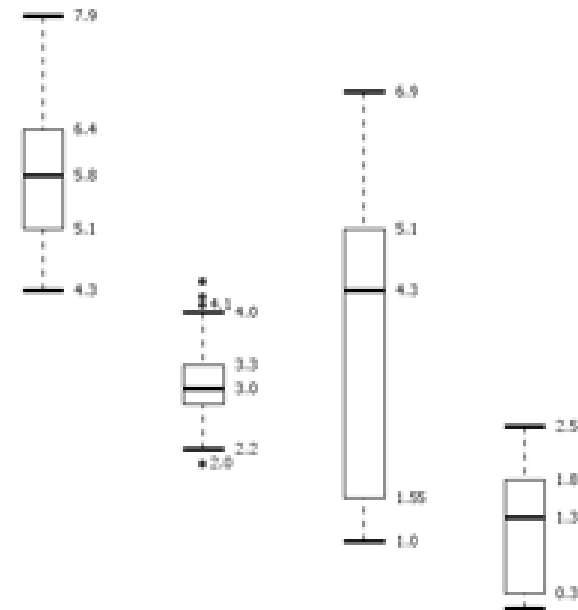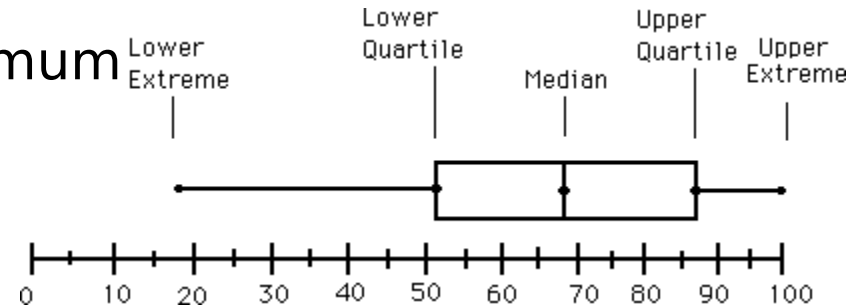    - **Standard deviation** *s (or σ)* is the square root of variance

# Boxplot Analysis
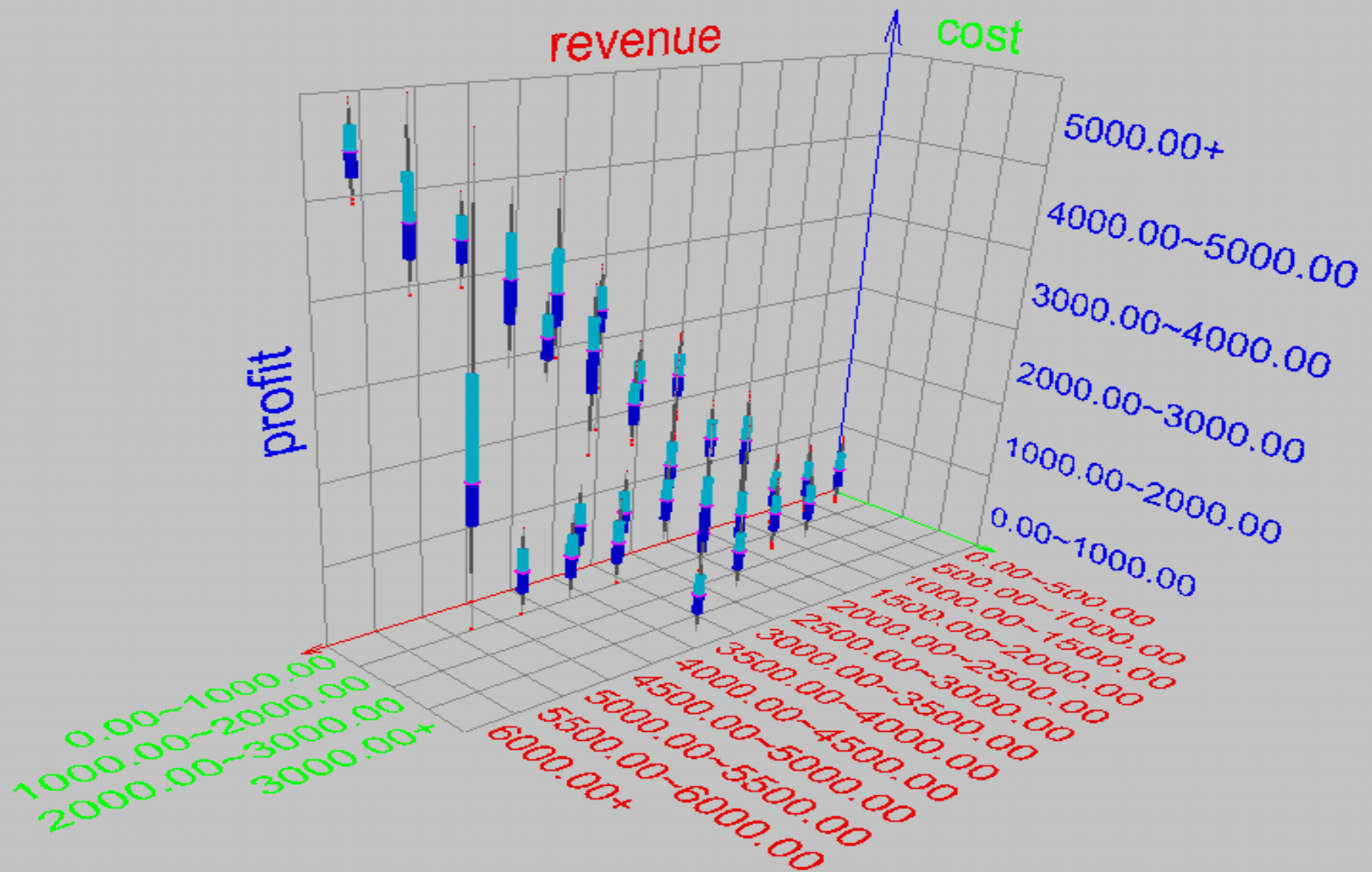
- **Five-number summary** of a distribution
  - Minimum, Q1, Median, Q3, Maximum
- **Boxplot**
  - Data is represented with a box
  - The ends of the box are at the first and third quartiles, i.e., the height of the box is IQR
  - The median is marked by a line within the box
  - Whiskers: two lines outside the box extended to Minimum and Maximum
  - Outliers: points beyond a specified outlier threshold, plotted individually

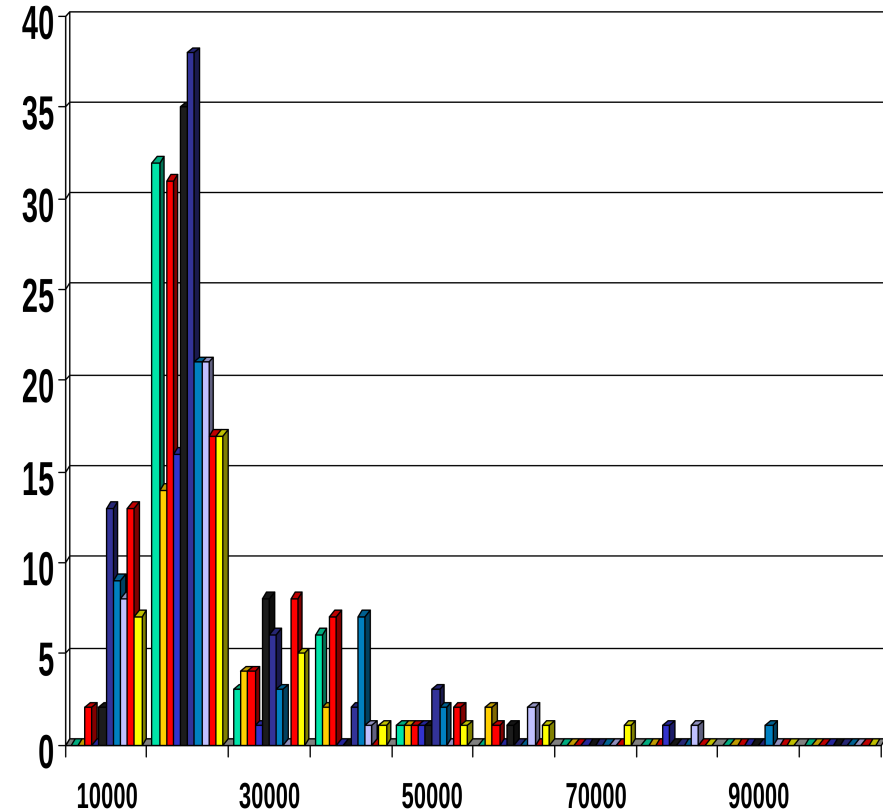# Visualization of Data Dispersion: 3-D Boxplots

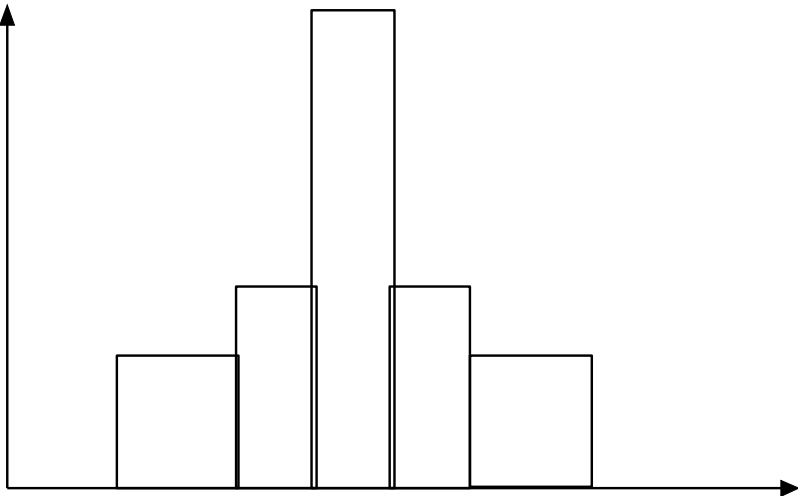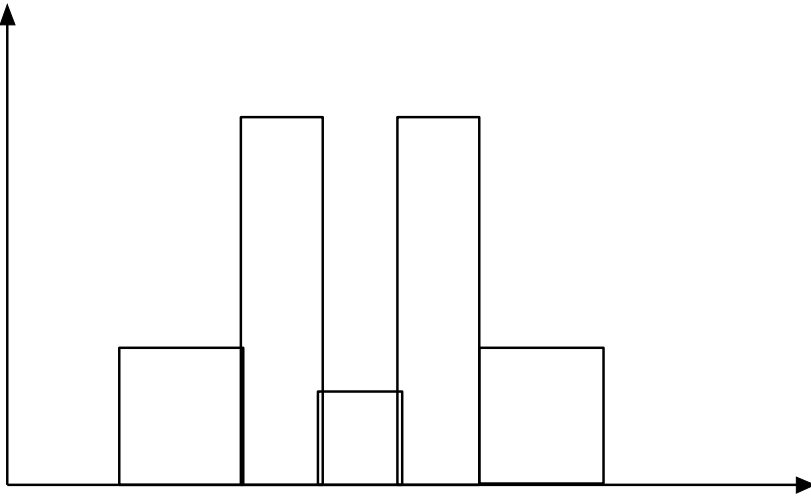# Graphic Displays of Basic Statistical Descriptions

- **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

# Histogram Analysis

- Histogram: Graph display of tabulated frequencies, shown as bars

- It shows what proportion of cases fall into each of several categories

- Differs from a bar chart in that it is the *area* of the bar that denotes the value, not the height as in bar charts, a crucial distinction when the categories are not of uniform width

- The categories are usually specified as non-overlapping intervals of some variable. The categories (bars) must be adjacent
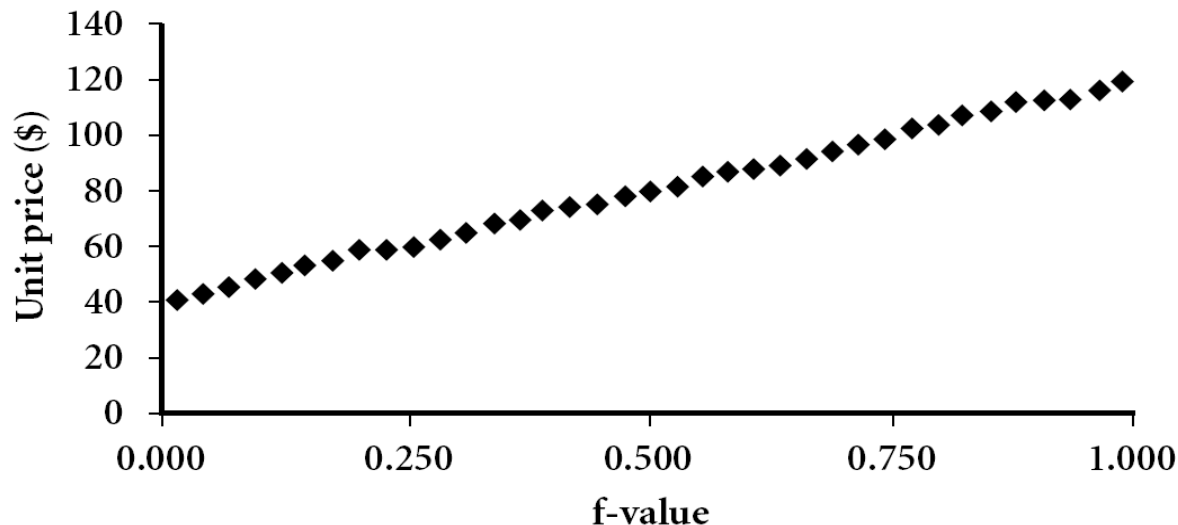
# Histograms Often Tell More than Boxplots



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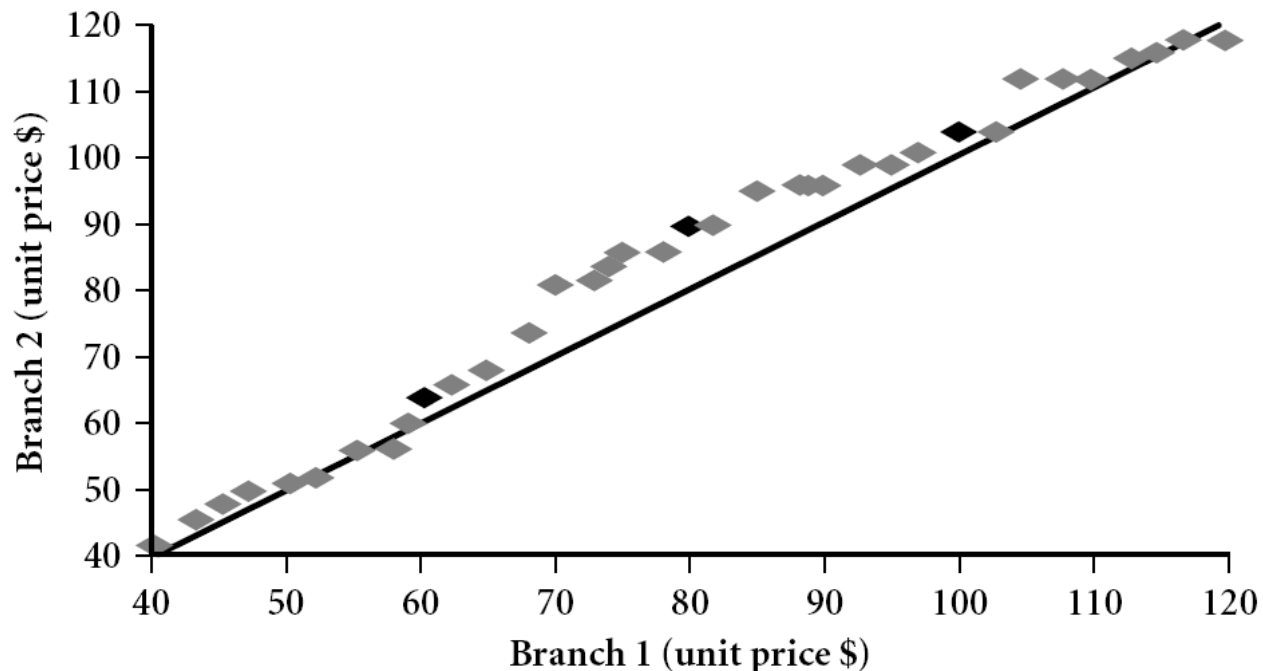- But they have rather different data distributions

# Quantile Plot

- Displays all of the data (assess both the overall behavior and unusual occurrences)
- Plots **quantile** information
  - Select an attribute $x_i$; *sort* data by non-decreasing $x_i$ value; plot it equally spaced on the x axis
  - v(f) indicates the value s.t. a fraction f of data has value at most v(f)
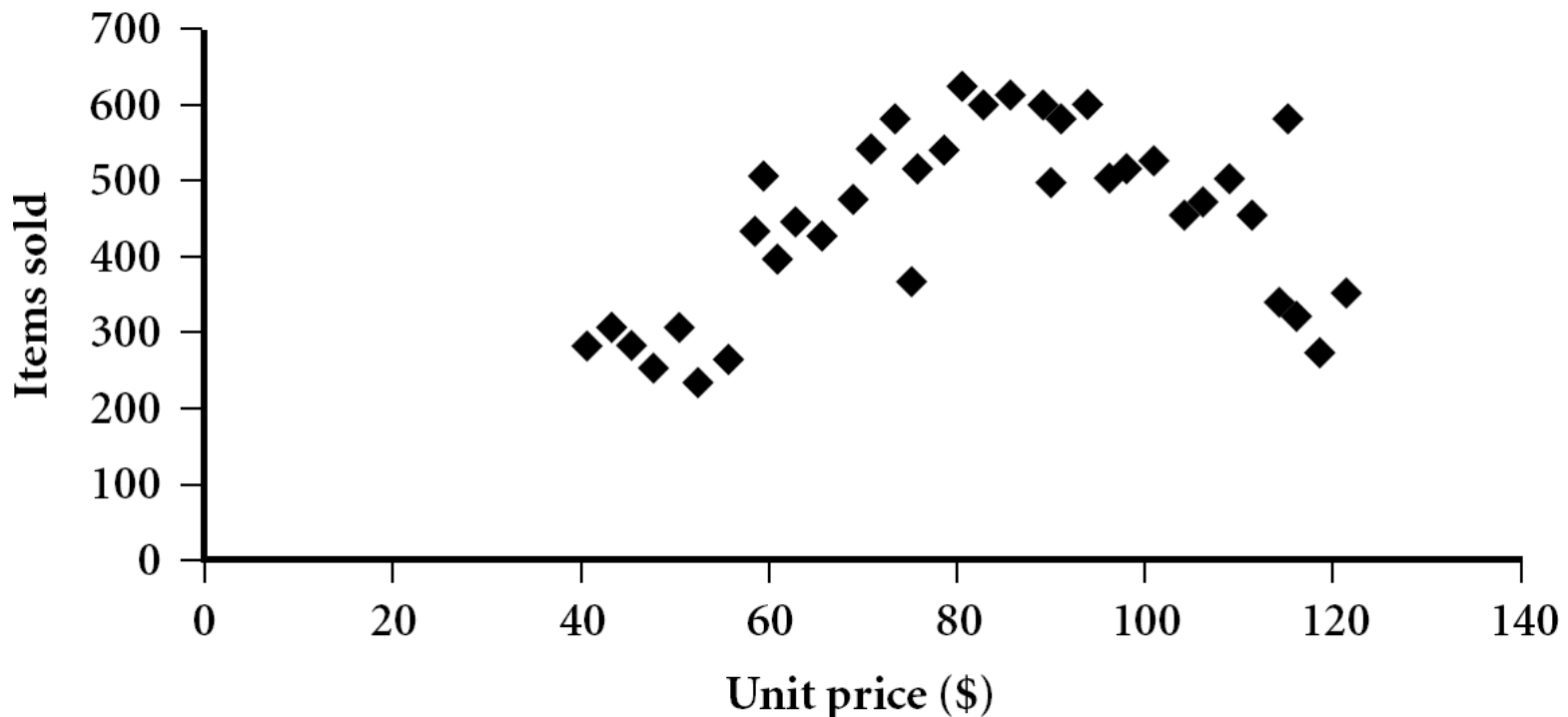
# Quantile-Quantile (Q-Q) Plot

- Graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- View: Is there is a shift in going from one distribution to another?
- Example shows unit price of items sold at Branch 1 vs. Branch 2 for each quantile.  Unit prices of items sold at Branch 1 tend to be lower than those at Branch 2.
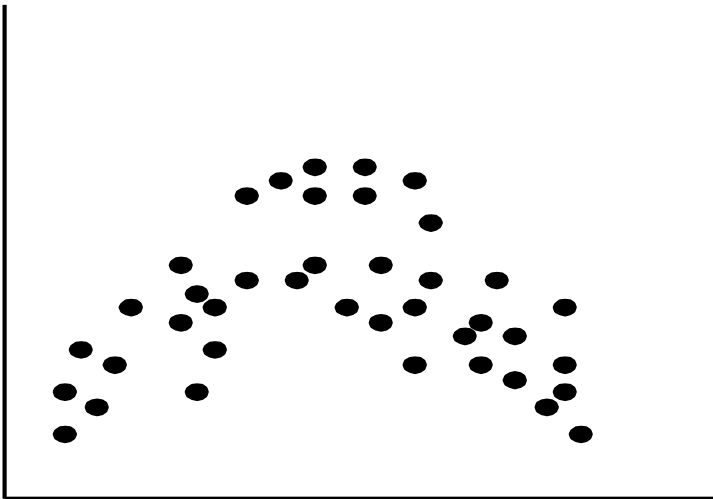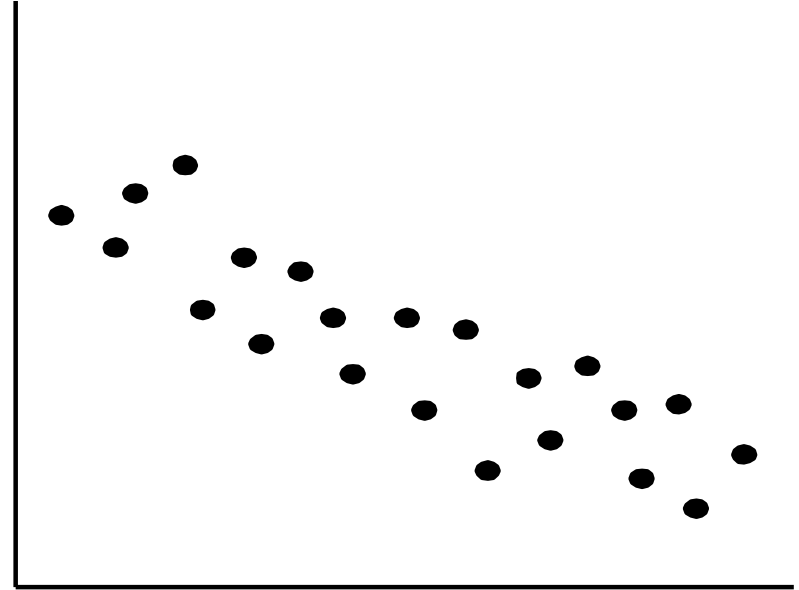
# Scatter plot

- 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

# Positively and Negatively Correlated Data

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

# Uncorrelated Data

# Scatterplot Matrices



Matrix of scatterplots (x-y-diagrams) of the k-dim. data [total of (k2/2-k) scatterplots]

# Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types

- Basic Statistical Descriptions of Data

- Measuring Data Similarity and Dissimilarity

- Summary

# Similarity and Dissimilarity

- **Similarity**
  - Numerical measure of how alike two data objects are
  - Value is higher when objects are more alike
  - Often falls in the range [0,1]
- **Dissimilarity** (e.g., distance)
  - Numerical measure of how different two data objects are
  - Lower when objects are more alike
  - Minimum dissimilarity is often 0
  - Upper limit varies
- **Proximity** refers to a similarity or dissimilarity

# Data Matrix and Dissimilarity Matrix

- ## Data matrix
  - n data points (objects) with p dimensions (features)

$$\begin{bmatrix} x_{11} & \dots & x_{1f} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{if} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nf} & \dots & x_{np} \end{bmatrix}$$

- ## Dissimilarity matrix
  - n data points, but registers only the distance
  - A triangular matrix

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

# Proximity Measures for Binary Attributes

- A contingency table for binary data

Number of attributes for which both data objects have value 1

Data object $j$

Data object $i$

|  | 1 | 0 | sum |
|---|---|---|---|
| 1 | $q$ | $r$ | $q+r$ |
| 0 | $s$ | $t$ | $s+t$ |
| sum | $q+s$ | $r+t$ | $p$ |

# Proximity Measures for Binary Attributes

- … but we can do the same for attributes (transpose)

Number of data objects for which both attributes have value 1

Attribute $j$

Attribute $i$

|  | 1 | 0 | sum |
|---|---|---|---|
| 1 | $q$ | $r$ | $q + r$ |
| 0 | $s$ | $t$ | $s + t$ |
| sum | $q + s$ | $r + t$ | $p$ |

# Proximity Measures for Binary Attributes

Data object $j$

- A contingency table for binary data

Data object $i$

|  | 1 | 0 | sum |
|---|---|---|---|
| 1 | $q$ | $r$ | $q+r$ |
| 0 | $s$ | $t$ | $s+t$ |
| sum | $q+s$ | $r+t$ | $p$ |

- Distance measure for symmetric bin. vars (0 and 1 equally important):

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

- Distance measure for asymm. bin. vars (1 more important – e.g. diseases):

$$d(i, j) = \frac{r + s}{q + r + s}$$

- Jaccard coefficient (*similarity* measure for *asymmetric* binary variables):

$$sim_{Jaccard}(i, j) = \frac{q}{q + r + s}$$

- Note: Jaccard coefficient is the same as "coherence":

$$coherence(i, j) = \frac{sup(i, j)}{sup(i) + sup(j) - sup(i, j)} = \frac{q}{(q + r) + (q + s) - q}$$

# Dissimilarity between Binary Attributes

- Example

| Name | Gender | Fever | Cough | Test-1 | Test-2 | Test-3 |
|------|--------|-------|-------|--------|--------|--------|
| Jack | M | Y | N | P | N | N |
| Mary | F | Y | N | P | N | P |
| Jim | M | Y | P | N | N | N |

- Gender is a symmetric attribute (let's discard it!)
- The remaining attributes are asymmetric binary
- Let the values Y and P be 1, and the value N 0

$$d(jack, mary) = \frac{0 + 1}{2 + 0 + 1} = 0.33$$

$$d(jack, jim) = \frac{1 + 1}{1 + 1 + 1} = 0.67$$

$$d(jim, mary) = \frac{1 + 2}{1 + 1 + 2} = 0.75$$

|   | 1 | 0 | sum |
|---|---|---|-----|
| 1 | $q$ | $r$ | $q + r$ |
| 0 | $s$ | $t$ | $s + t$ |
| sum | $q + s$ | $r + t$ | $p$ |

$$d(i, j) = \frac{r + s}{q + r + s}$$

# Proximity Measures for Categorical (or "nominal") Attributes

- Can take 2 or more states, e.g., red, yellow, blue, green (generalization of a binary attribute)

- <u>Method 1</u>: Simple matching
  - $m$: # of matches, $p$: total # of attributes

$$d(i,j) = \frac{p-m}{p}$$

- <u>Method 2</u>: Use a large number of binary attributes
  - creating a new binary attribute for each of the $M$ categories

# Proximity on Numeric Data: Minkowski Distance

- *Minkowski distance*: A popular distance measure

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \cdots + |x_{ip} - x_{jp}|^h}$$

  where $i = (x_{i1}, x_{i2}, \ldots, x_{ip})$ and $j = (x_{j1}, x_{j2}, \ldots, x_{jp})$ are two $p$-dimensional data objects, and $h$ is the order (the distance so defined is also called L-$h$ norm)

- Properties

  - $d(i, j) > 0$ if $i \neq j$, and $d(i, i) = 0$ (Positive definiteness)

  - $d(i, j) = d(j, i)$ (Symmetry)

  - $d(i, j) \leq d(i, k) + d(k, j)$ (Triangle Inequality)

- A distance that satisfies these properties is a metric

# Special Cases of Minkowski Distance

- $h = 1$:  Manhattan (city block, $L_1$ norm) distance
  - E.g., the Hamming distance: the number of bits that are different between two binary vectors

$$d(i,j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \ldots + |x_{ip} - x_{jp}|$$

- $h = 2$:  ($L_2$ norm) Euclidean distance

$$d(i,j) = \sqrt{\left( |x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \ldots + |x_{ip} - x_{jp}|^2 \right)}$$

- $h \rightarrow \infty$.  "supremum" ($L_{max}$ norm, $L_\infty$ norm) distance.
  - This is the maximum difference between any component (attribute) of the vectors

$$d(i,j) = \lim_{h \to \infty} \left( \sum_{f=1}^{p} |x_{if} - x_{jf}|^h \right)^{\frac{1}{h}} = \max_{f}^{p} |x_{if} - x_{jf}|$$

# Example: Minkowski Distance

| point | attribute 1 | attribute 2 |
|---|---|---|
| x1 | 1 | 2 |
| x2 | 3 | 5 |
| x3 | 2 | 0 |
| x4 | 4 | 5 |

## Manhattan ($L_1$)

| L | x1 | x2 | x3 | x4 |
|---|---|---|---|---|
| x1 | 0 | | | |
| x2 | 5 | 0 | | |
| x3 | 3 | 6 | 0 | |
| x4 | 6 | 1 | 7 | 0 |

## Euclidean ($L_2$)

| L2 | x1 | x2 | x3 | x4 |
|---|---|---|---|---|
| x1 | 0 | | | |
| x2 | 3.61 | 0 | | |
| x3 | 2.24 | 5.1 | 0 | |
| x4 | 4.24 | 1 | 5.39 | 0 |

**Dissimilarity Matrices**

## Supremum ($L_{inf}$)

| $L_\infty$ | x1 | x2 | x3 | x4 |
|---|---|---|---|---|
| x1 | 0 | | | |
| x2 | 3 | 0 | | |
| x3 | 2 | 5 | 0 | |
| x4 | 3 | 1 | 5 | 0 |

# Standardizing Numeric Data

- Z-score: $z = \dfrac{x - \mu}{\sigma}$

  - X: raw data, μ: mean of the population, σ: standard deviation
  - the distance between the raw score and the population mean in units of the standard deviation
  - <0 when the raw score is below the mean, >0 when above
- An alternative way: Calculate the mean absolute deviation

$$s_f = \frac{1}{n}\left(\left|x_{1f} - m_f\right| + \left|x_{2f} - m_f\right| + \ldots + \left|x_{nf} - m_f\right|\right)$$
where

$$m_f = \frac{1}{n}(x_{1f} + x_{2f} + \ldots + x_{nf}).$$

  - standardized measure (*z-score*): $z_{if} = \dfrac{x_{if} - m_f}{s_f}$

- mean absolute deviation is more robust than std dev

# **Properties of Normal Distribution Curve**

- The normal (distribution) curve
  - From μ–σ to μ+σ: contains about 68% of the measurements  (μ: mean, σ: standard deviation)
  -  From μ–2σ to μ+2σ: contains about 95% of it
  - From μ–3σ to μ+3σ: contains about 99.7% of it

# Ordinal Variables

- An ordinal variable can be discrete or continuous

- Order is important, e.g., rank

- Can be treated like interval-scaled

  - replace $x_{if}$ by their rank $r_{if} \in \{1, ..., M_f\}$

  - map (normalize) the range of each variable onto [0, 1] by replacing $x_{if}$ by

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

  - compute the dissimilarity using distance measures for numeric attributes

# Attributes of Mixed Type

- A database may contain all attribute types
  - Nominal, symmetric binary, asymmetric binary, numeric, ordinal
- One may use a weighted formula to combine their effects

$$d(i,j) = \frac{\Sigma_{f=1}^{p} \delta_{ij}^{(f)} d_{ij}^{(f)}}{\Sigma_{f=1}^{p} \delta_{ij}^{(f)}}$$

- Choice of $\delta_{ij}^{(f)}$
  - Set $\delta_{ij}^{(f)} = 0$ if
    - $x_{if}$ or $x_{jf}$ is missing
    - $x_{if} = x_{jf} = 0$ and f is asymmetric binary
  - Set $\delta_{ij}^{(f)} = 1$ otherwise

# Attributes of Mixed Type

$$d(i,j) = \frac{\sum_{f=1}^{p} \delta_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^{p} \delta_{ij}^{(f)}}$$

- *Choice of $d_{ij}^{(f)}$*
  - *when f is binary or nominal:*
    $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$, $d_{ij}^{(f)} = 1$ otherwise
  - *when f is numeric:* use the normalized distance
  - *when f is ordinal*

    - Compute ranks $r_{if}$ and $z_{if} = \dfrac{r_{if} - 1}{M_f - 1}$

    - Treat $z_{if}$ as interval-scaled

# Cosine Similarity

- A **document** can be represented by thousands of attributes, each recording the *frequency* of a particular word (such as keywords) or phrase in the document.

| Document | team | coach | hockey | baseball | soccer | penalty | score | win | loss | season |
|----------|------|-------|--------|----------|--------|---------|-------|-----|------|--------|
| Document1 | 5 | 0 | 3 | 0 | 2 | 0 | 0 | 2 | 0 | 0 |
| Document2 | 3 | 0 | 2 | 0 | 1 | 1 | 0 | 1 | 0 | 1 |
| Document3 | 0 | 7 | 0 | 2 | 1 | 0 | 0 | 3 | 0 | 0 |
| Document4 | 0 | 1 | 0 | 0 | 1 | 2 | 2 | 0 | 3 | 0 |

- Other vector objects: gene features in micro-arrays, …
- Applications: information retrieval, biologic taxonomy, gene feature mapping, …
- Issue: very long and **sparse**

- Treat documents as vectors, and compute a **cosine similarity**

# Cosine Similarity

- Cosine measure: If x and y are two vectors (e.g., term-frequency vectors), then

$$\cos(x, y) = (x \bullet y) / ||x|| \, ||y||$$

  where
  - • indicates vector dot product, $\quad x \cdot y = \sum_{i=1}^{p} x_i \, y_i$
  - $||x||$: the L2 norm (length) of vector x $\quad ||x|| = \sqrt{x_1^2 + x_2^2 + \ldots + x_p^2}$

- Remark: when attributes are binary valued:
  - • indicates the number of shared features
  - $||x|| \, ||y||$ is the geometric mean between the number of features of x and the number of features of y:

        sqrt(a) * sqrt(b) = sqrt( a * b )

  - cos (x, y) measures relative possession of common features

# Example: Cosine Similarity

- $\cos(x, y) = (x \cdot y) / \|x\| \|y\|$

- Ex: Find the **similarity** between documents x and y.

  **x =** (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)
  **y =** (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)

  x • y = 5*3+0*0+3*2+0*0+2*1+0*1+0*1+2*1+0*0+0*1= = 25
  $\|x\| = (5*5+0*0+3*3+0*0+2*2+0*0+0*0+2*2+0*0+0*0)^{0.5} =$ = 6.481
  $\|y\| = (3*3+0*0+2*2+0*0+1*1+1*1+0*0+1*1+0*0+1*1)^{0.5} =$ = 4.12
  $\cos(x, y) = 25 / (6.481 * 4.12) = 0.94$

# References

- W. Cleveland, Visualizing Data, Hobart Press, 1993
- T. Dasu and T. Johnson.  Exploratory Data Mining and Data Cleaning. John Wiley, 2003
- U. Fayyad, G. Grinstein, and A. Wierse. Information Visualization in Data Mining and Knowledge Discovery, Morgan Kaufmann, 2001
- L. Kaufman and P. J. Rousseeuw. Finding Groups in Data: an Introduction to Cluster Analysis. John Wiley & Sons, 1990.
- H. V. Jagadish et al., Special Issue on Data Reduction Techniques. Bulletin of the Tech. Committee on Data Eng., 20(4), Dec. 1997
- D. A. Keim. Information visualization and visual data mining, IEEE trans. on Visualization and Computer Graphics, 8(1), 2002
- D. Pyle. Data Preparation for Data Mining. Morgan Kaufmann, 1999
- S.  Santini and R. Jain," Similarity measures", IEEE Trans. on Pattern Analysis and Machine Intelligence, 21(9), 1999
- E. R. Tufte. The Visual Display of Quantitative Information, 2nd ed., Graphics Press, 2001
- C. Yu et al., Visual data mining of multimedia data for social and behavioral studies, Information Visualization, 8(1), 2009