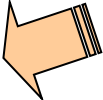# Data Warehousing and Data Mining

## —Xike Xie —

Slides are based on Prof. Han and Prof. Tan's works.

# Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types

- An Example of Clustering

- Basic Statistical Descriptions of Data

- Measuring Data Similarity and Dissimilarity

- Summary

# What is Data?

- Collection of data objects and their attributes

- An attribute is a property or characteristic of an object
  - Examples: eye color of a person, temperature, etc.
  - Attribute is also known as variable, field, characteristic, or feature
- A collection of attributes describe an object
  - Object is also known as record, point, case, sample, entity, or instance

Attributes

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

Objects

# 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
    - 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*

# Properties of Attribute Values

- The type of an attribute depends on which of the following properties it possesses:
    - Distinctness:          =  ≠
    - Order:            <  >
    - Addition:            +  -
    - Multiplication:        *  /

    - Nominal attribute: distinctness
    - Ordinal attribute: distinctness & order
    - Interval attribute: distinctness, order & addition
    - Ratio attribute: all 4 properties

| Attribute Type | Description | Examples | Operations |
|---|---|---|---|
| Nominal | The values of a nominal attribute are just different names, i.e., nominal attributes provide only enough information to distinguish one object from another. ($=$, $\neq$) | zip codes, employee ID numbers, eye color, sex: {*male, female*} | mode, entropy, contingency correlation, $\chi^2$ test |
| Ordinal | The values of an ordinal attribute provide enough information to order objects. ($<$, $>$) | hardness of minerals, {*good, better, best*}, grades, street numbers | median, percentiles, rank correlation, run tests, sign tests |
| Interval | For interval attributes, the differences between values are meaningful, i.e., a unit of measurement exists. ($+$, $-$) | calendar dates, temperature in Celsius or Fahrenheit | mean, standard deviation, Pearson's correlation, $t$ and $F$ tests |
| Ratio | For ratio variables, both differences and ratios are meaningful. ($*$, $/$) | temperature in Kelvin, monetary quantities, counts, age, mass, length, electrical current | geometric mean, harmonic mean, percent variation |

| Attribute Level | Transformation | Comments |
|---|---|---|
| Nominal | Any permutation of values | If all employee ID numbers were reassigned, would it make any difference? |
| Ordinal | An order preserving change of values, i.e., *new_value = f(old_value)* where *f* is a monotonic function. | An attribute encompassing the notion of good, better best can be represented equally well by the values {1, 2, 3} or by { 0.5, 1, 10}. |
| Interval | *new_value =a \* old_value + b* where a and b are constants | Thus, the Fahrenheit and Celsius temperature scales differ in terms of where their zero value is and the size of a unit (degree). |
| Ratio | *new_value = a \* old_value* | Length can be measured in meters or feet. |

# 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
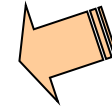  - Note: Binary attributes are a special case of discrete attributes
- **Continuous Attribute**
  - Has real numbers as attribute values
    - E.g., temperature, height, or weight
  - Practically, real values can only be measured and represented using a finite number of digits
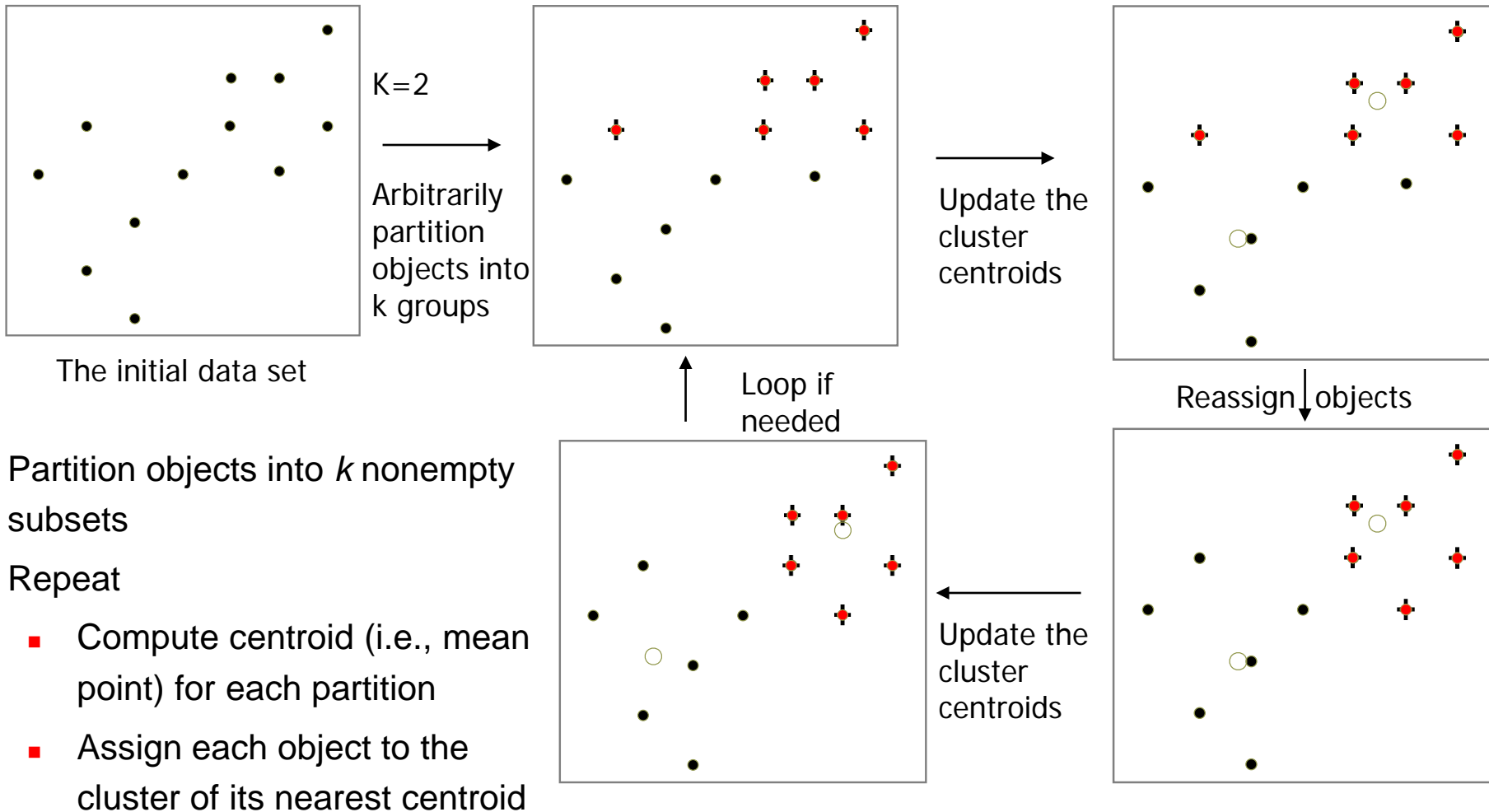  - Continuous attributes are typically represented as floating-point variables

# Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types

- An Example of Clustering

- Basic Statistical Descriptions of Data

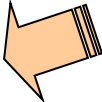- Measuring Data Similarity and Dissimilarity

- Summary

# An Example of *K-Means* Clustering



The initial data set

K=2

Arbitrarily partition objects into k groups

Update the cluster centroids

Reassign objects

Update the cluster centroids

Loop if needed

- Partition objects into *k* nonempty subsets
- Repeat
  - Compute centroid (i.e., mean point) for each partition
  - Assign each object to the cluster of its nearest centroid
- Until no change

# Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types

- An Example of Clustering

- Basic Statistical Descriptions of Data

- Measuring Data Similarity and Dissimilarity

- Summary

# Basic Statistical Descriptions of Data

- 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):

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

  $$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i \qquad \mu = \frac{\sum x}{N}$$

  - Weighted arithmetic mean:
  - Trimmed mean: chopping extreme values

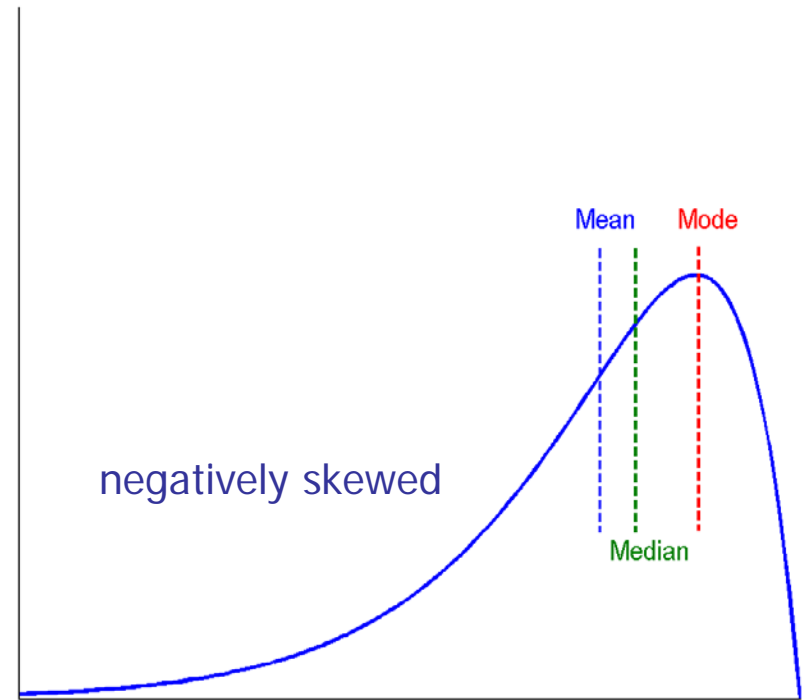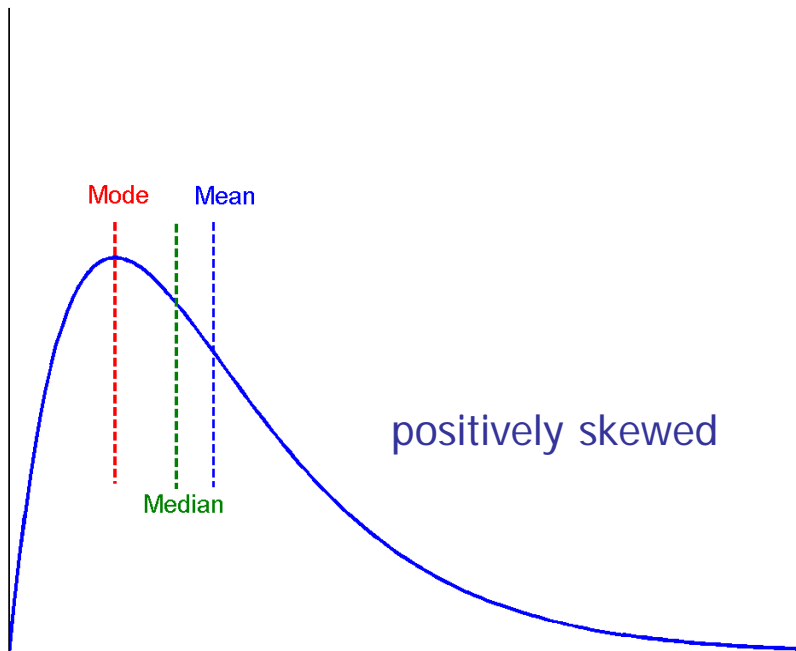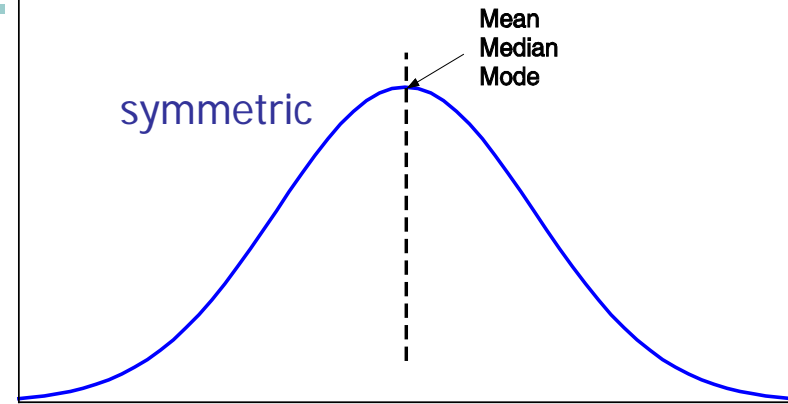  $$\bar{x} = \frac{\sum_{i=1}^{n} w_i x_i}{\sum_{i=1}^{n} w_i}$$

- Median:

  - Middle value if odd number of values, or average of the middle two values otherwise
  - Estimated by interpolation (for *grouped data*):

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

- Mode

  - Value that occurs most frequently in the data
  - Unimodal, bimodal, trimodal
  - Empirical formula: $mean - mode = 3 \times (mean - median)$

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

# Symmetric vs. Skewed Da

- 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**: $Q_1$ (25$^{th}$ percentile), $Q_3$ (75$^{th}$ percentile)

  - **Inter-quartile range**: IQR = $Q_3 - Q_1$

  - **Five number summary**: min, $Q_1$, median, $Q_3$, max

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

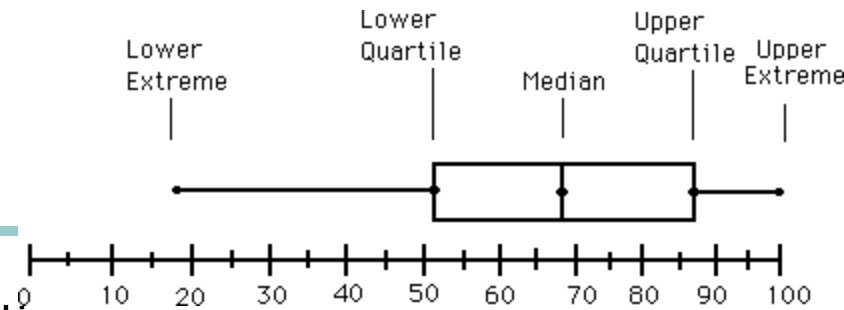  - **Outlier**: usually, a value higher/lower than 1.5 x IQR

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

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

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

  - **Standard deviation** *s (or σ)* is the square root of variance *s² (or σ²)*
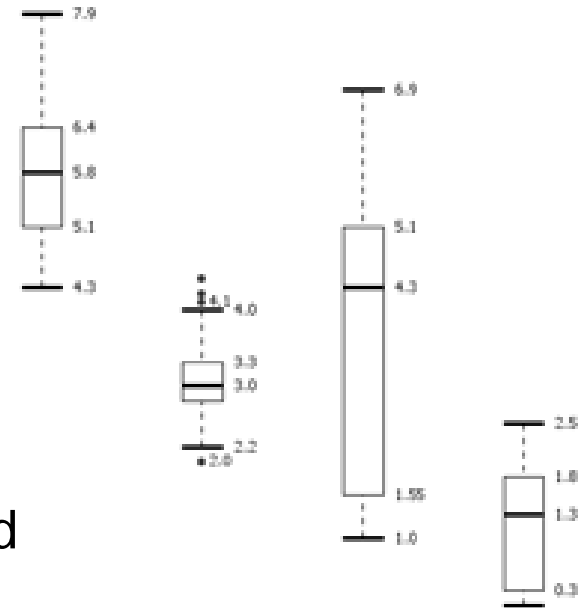
# Boxplot Analysis

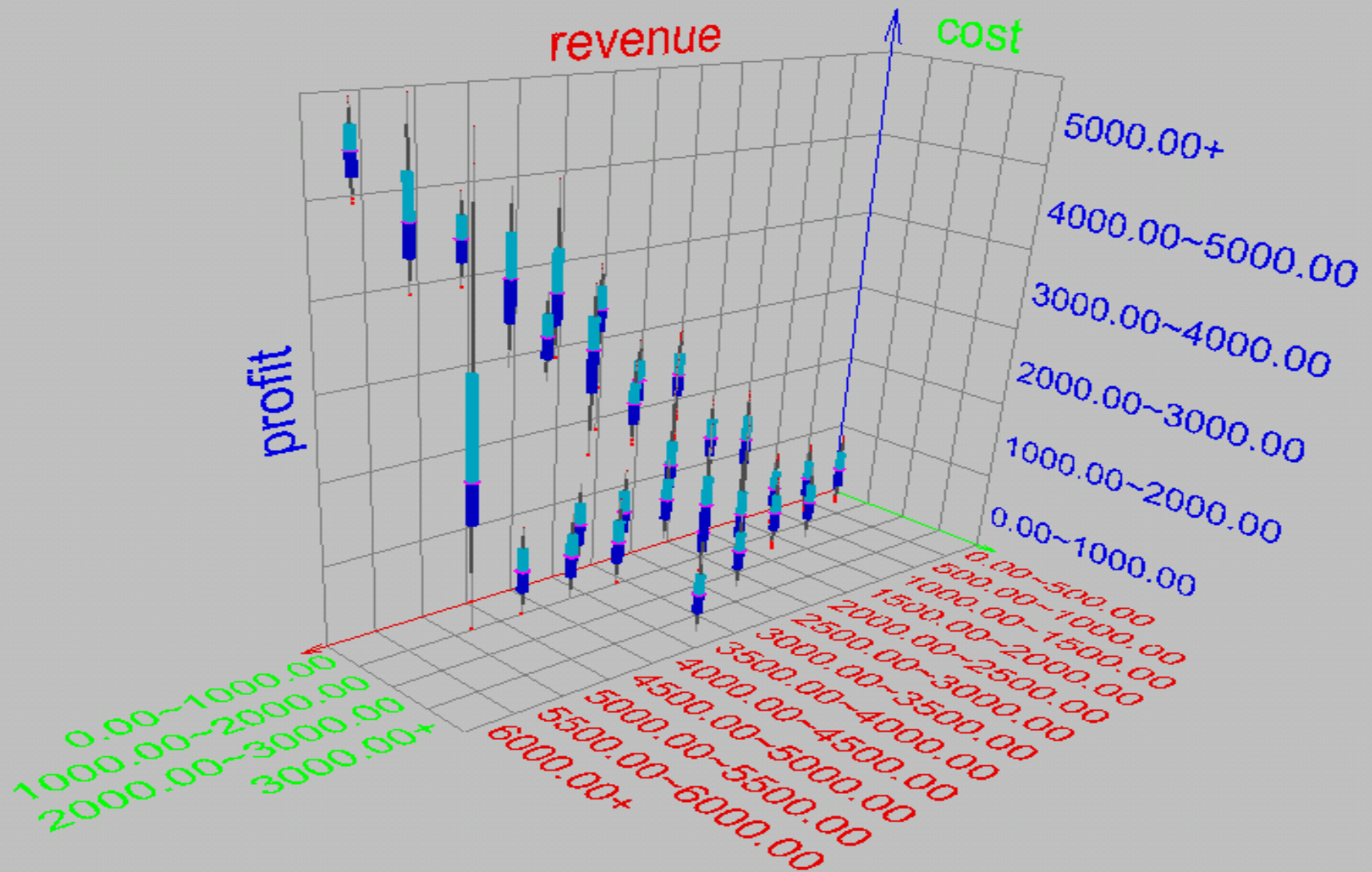- **Five-number summary** of a distribution
  - Minimum, Q1, Median, Q3, Maximum
- **Boxplot**
  - Data is represented with a box
  - The ends of the box are at the first and third quartiles, i.e., the height of the box is IQR
  - The median is marked by a line within the box
  - Whiskers: two lines outside the box extended to Minimum and Maximum
  - Outliers: points beyond a specified outlier threshold, plotted individually
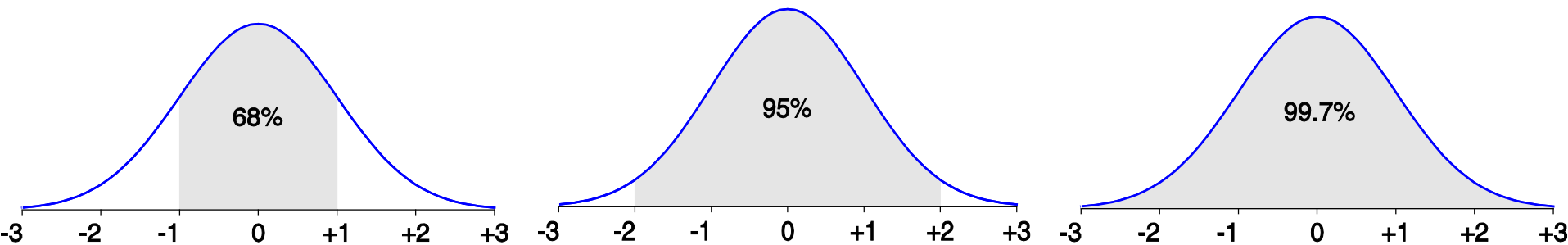
# Visualization of Data Dispersion: 3-D Boxplots

# 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
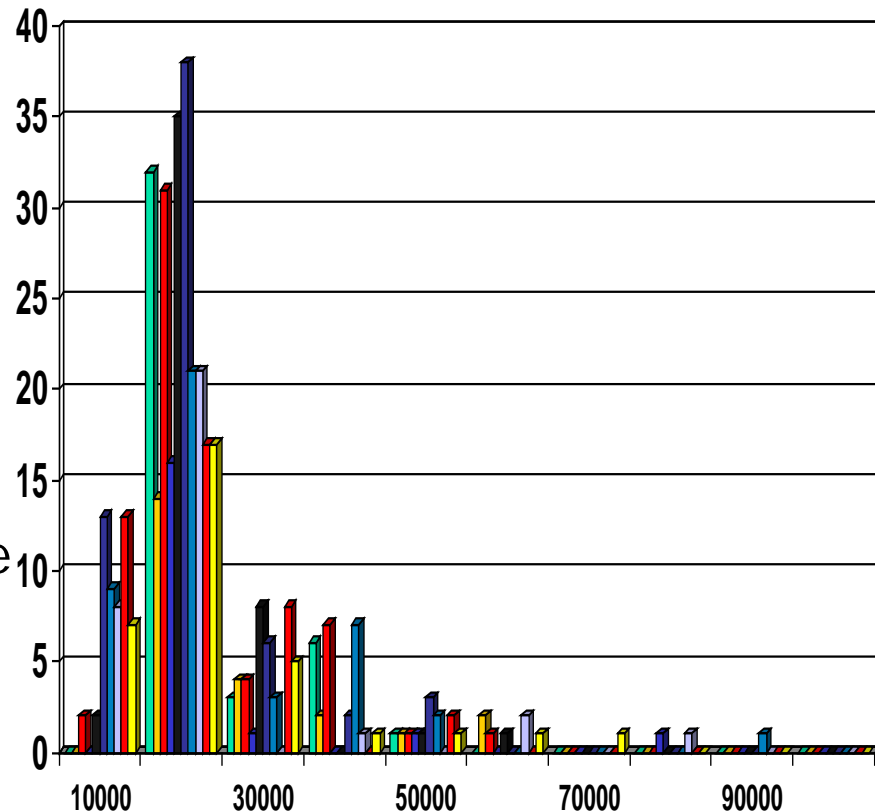  - From µ−3σ to µ+3σ: contains about 99.7% of it

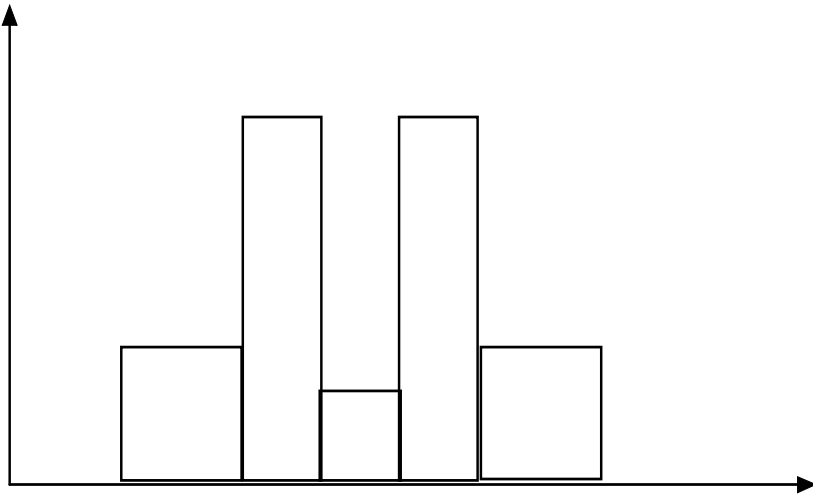# Graphic Displays of Basic Statistical Descriptions

- **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 (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
- 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

# Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types

- An Example of Data Clustering

- 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 with p dimensions
  - Two modes

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

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

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

# Proximity Measure for 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 variables

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

- <u>Method 2</u>: Use a large number of binary attributes

  - creating a new binary attribute for each of the $M$ nominal states

# Proximity Measure for Binary Attributes

- A contingency table for binary data

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

- Distance measure for symmetric binary variables:

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

- Distance measure for asymmetric binary variables:

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

- 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":

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

33

# Dissimilarity between Binary Variables

- Example

| Name | Gender | Fever | Cough | Test-1 | Test-2 | Test-3 | Test-4 |
|------|--------|-------|-------|--------|--------|--------|--------|
| Jack | M | Y | N | P | N | N | N |
| Mary | F | Y | N | P | N | P | N |
| Jim | M | Y | P | N | N | N | N |

- Gender is a symmetric attribute
- The remaining attributes are asymmetric binary
- Let the values Y and P be 1, and the value N 0

$$d(jack, mary) =$$

$$d(jack, jim) =$$

$$d(jim, mary) =$$

# Dissimilarity between Binary Variables

- Example

| Name | Gender | Fever | Cough | Test-1 | Test-2 | Test-3 | Test-4 |
|------|--------|-------|-------|--------|--------|--------|--------|
| Jack | M | Y | N | P | N | N | N |
| Mary | F | Y | N | P | N | P | N |
| Jim | M | Y | P | N | N | N | N |

- Gender is a symmetric attribute
- The remaining attributes are asymmetric binary
- Let the values Y and P be 1, and the value N 0

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

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

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

# Example:
# Data Matrix and Dissimilarity Matrix

## Data Matrix

| point | attribute1 | attribute2 |
|-------|-----------|-----------|
| x1 | 1 | 2 |
| x2 | 3 | 5 |
| x3 | 2 | 0 |
| x4 | 4 | 5 |

## Dissimilarity Matrix

## (with Euclidean Distance)

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

# Distance 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

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

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

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

- $h \to \infty$.  "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** | | | | |
| **x2** | | | | |
| **x3** | | | | |
| **x4** | | | | |

## 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 |

## Supremum

| **$L_\infty$** | **x1** | **x2** | **x3** | **x4** |
|----------------|--------|--------|--------|--------|
| **x1** | | | | |
| **x2** | | | | |
| **x3** | | | | |
| **x4** | | | | |

# 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 |

## Supremum

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

# Distance between data

## Comparison between Ordinal Attributes:

1. Questions:
   1. The grade of example have 5 levels such as A,B,C,D,E, how to judge the grade of two students in 5 courses
   2. There are 100 people participate in the election, and there are 20 juries, how to give the final results by integrating the results of these 20 juries?

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

2. Way Out:
   1. Convert the result of level into corresponding number
   2. Convert all different attributes which have different levels into [0,1]
   3. Use the way to calculate the distance after converting all ordinal attributes to numeric attributes

# Distance between data

## Ordinal Attributes: Complex Examples

1. In the given table, the three attributes is nominal, ordinal and numeric attribute, and we need to calculate the distance

2. Problem: we can calculate the distance for each individual attribute, but how to integrate them? Since all attributes have different physical meanings

3. Assuming fair=0, good=0.5, excellent=1

Table 2.2: A sample data table containing attributes of mixed type.

| object identifier | test-1 (nominal) | test-2 (ordinal) | test-3 (numeric) |
|---|---|---|---|
| 1 | code-A | excellent | 45 |
| 2 | code-B | fair | 22 |
| 3 | code-C | good | 64 |
| 4 | code-A | excellent | 28 |

$$Dist(test-1) = \begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

$$Dist(test-2) = \begin{bmatrix} 0 & & & \\ 1.0 & 0 & & \\ 0.5 & 0.5 & 0 & \\ 0 & 1.0 & 0.5 & 0 \end{bmatrix}$$

$$Dist(test-3) = \begin{bmatrix} 0 & & & \\ 0.55 & 0 & & \\ 0.45 & 1 & 0 & \\ 0.40 & 0.14 & 0.86 & 0 \end{bmatrix}$$

4. Convert the distance of each dimension into [0,1], and then calculate the weighted sum of all dimensions. Assuming the weight of each dimension is equivalent, we have

$$\mathcal{D} = (Dist(test-1) + Dist(test-2) + Dist(test-3))/3 = \begin{bmatrix} 0 & & & \\ 0.85 & 0 & & \\ 0.65 & 0.83 & 0 & \\ 0.13 & 0.71 & 0.79 & 0 \end{bmatrix}$$

**How to determine the weight of each dimension?**

# 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, ...
- Cosine measure: If $d_1$ and $d_2$ are two vectors (e.g., term-frequency vectors), then

$$\cos(d_1, d_2) = (d_1 \bullet d_2) / ||d_1|| \, ||d_2|| \, ,$$

where $\bullet$ indicates vector dot product, $||d||$: the length of vector $d$

# Example: Cosine Similarity

- $\cos(d_1, d_2) = (d_1 \bullet d_2)/||d_1|| \, ||d_2||$ ,
  where $\bullet$ indicates vector dot product, $||d||$: the length of vector $d$

- Ex: Find the **similarity** between documents 1 and 2.

$d_1 = (5, 0, 3, 0, 2, 0, 0, 2, 0, 0)$
$d_2 = (3, 0, 2, 0, 1, 1, 0, 1, 0, 1)$

$d_1 \bullet d_2 = 5*3+0*0+3*2+0*0+2*1+0*1+0*1+2*1+0*0+0*1 = 25$
$||d_1|| = (5*5+0*0+3*3+0*0+2*2+0*0+0*0+2*2+0*0+0*0)^{0.5} = (42)^{0.5}$
$\quad = 6.481$
$||d_2|| = (3*3+0*0+2*2+0*0+1*1+1*1+0*0+1*1+0*0+1*1)^{0.5} = (17)^{0.5}$
$\quad = 4.12$
$\cos(d_1, d_2) = 0.94$

# Chapter 2: Getting to Know Your Data

- Data Objects and Attribute Types

- Basic Statistical Descriptions of Data

- An Example of Data Clustering

- Measuring Data Similarity and Dissimilarity

- Summary

# Summary

- Data attribute types: nominal, binary, ordinal, interval-scaled, ratio-scaled

- Many types of data sets, e.g., numerical, text, graph, Web, image.

- Gain insight into the data by:

  - Basic statistical data description: central tendency, dispersion, graphical displays

  - Data visualization: map data onto graphical primitives

  - Measure data similarity

- Above steps are the beginning of data preprocessing.

- Many methods have been developed but still an active area of research.

# 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