



Data

Dr. N. Abdolvand

Data Mining

Data Mining: Concepts and Techniques

— Chapter 2 —

Jiawei Han, Micheline Kamber, and Jian Pei
University of Illinois at Urbana-Champaign
Simon Fraser University
©2013 Han, Kamber, and Pei. All rights reserved.

Chapter 2: Getting to Know Your Data

- ◆ Data Objects and Attribute Types
- ◆ Basic Statistical Descriptions of Data
- ◆ Data Visualization
- ◆ Measuring Data Similarity and Dissimilarity
- ◆ Summary

3

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:
 - Video data:

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

<i>TID</i>	<i>Items</i>
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

4

Important Characteristics of Structured Data

- ◆ Dimensionality
 - Curse of dimensionality
- ◆ Sparsity
 - Only presence counts
- ◆ Resolution
 - Patterns depend on the scale
- ◆ Distribution
 - Centrality and dispersion

5

Data Objects

- ◆ Data sets are made up of data objects.
- ◆ A data object represents an entity.
- ◆ Examples:
 - sales database: customers, store items, sales
 - medical database: patients, treatments
 - university database: students, professors, courses
- ◆ Also called samples , examples, instances, data points, objects, tuples.
- ◆ Data objects are described by attributes.
- ◆ Database rows -> data objects; columns -> attributes.

6

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

7

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)
 - Symmetric binary: both outcomes equally important
 - e.g., gender
 - Asymmetric binary: 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

8

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*

9

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

10

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

11

Measuring the Central Tendency

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

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

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

 - Weighted arithmetic mean:

$$\bar{x} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$
 - Trimmed mean: chopping extreme values
- ◆ Median:
 - Middle value if odd number of values, or average of the middle two values otherwise
 - Estimated by interpolation (for grouped data):

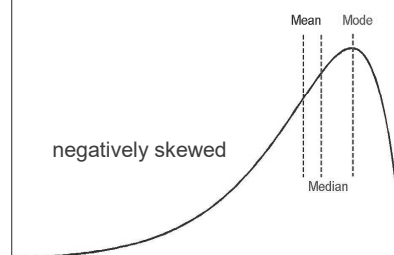
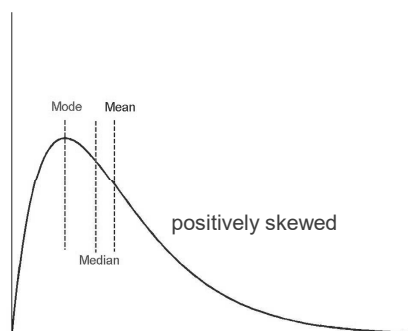
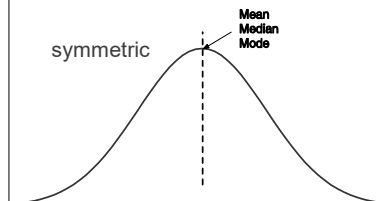
age	frequency
1-5	200
6-15	450
16-20	300
21-50	1500
51-80	700
81-110	44

$median = L_1 + \left(\frac{n/2 - (\sum freq)_l}{freq_{median}} \right) width$ Median interval →
- ◆ Mode
 - Value that occurs most frequently in the data
 - Unimodal, bimodal, trimodal
 - Empirical formula: $mean - mode = 3 \times (mean - median)$

12

Symmetric vs. Skewed Data

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



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
 - **Boxplot:** ends of the box are the quartiles; median is marked; add whiskers, and plot outliers individually
 - **Outlier:** usually, a value higher/lower than $1.5 \times IQR$

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

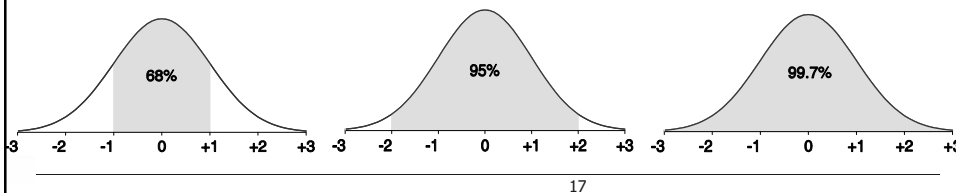
- **Variance:** (algebraic, scalable computation)

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

- **Standard deviation s (or σ)** is the square root of variance s^2 (or σ^2)

Properties of Normal Distribution Curve

- ◆ The normal (distribution) curve
 - From $\mu - \sigma$ to $\mu + \sigma$: contains about 68% of the measurements (μ : mean, σ : 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 univariate 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

19

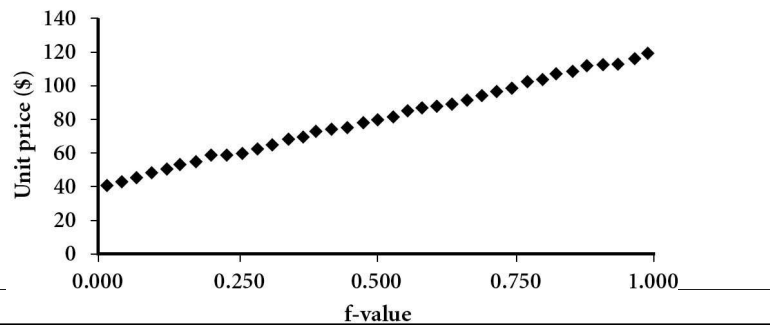
Histograms Often Tell More than Boxplots

- 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

20

Quantile Plot

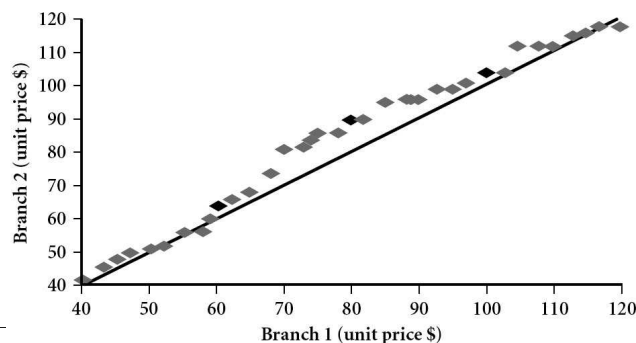
- ◆ Displays all of the data (allowing the user to assess both the overall behavior and unusual occurrences)
- ◆ Plots **quantile** information
 - For a data x_i , data sorted in increasing order, f_i indicates that approximately $100 f_i\%$ of the data are below or equal to the value x_i



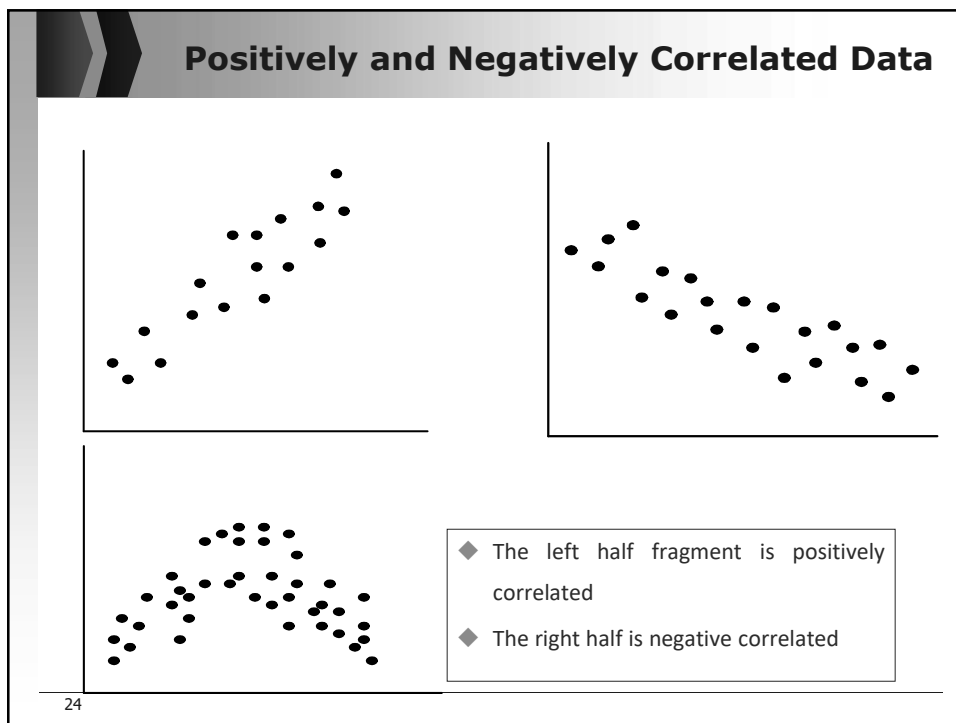
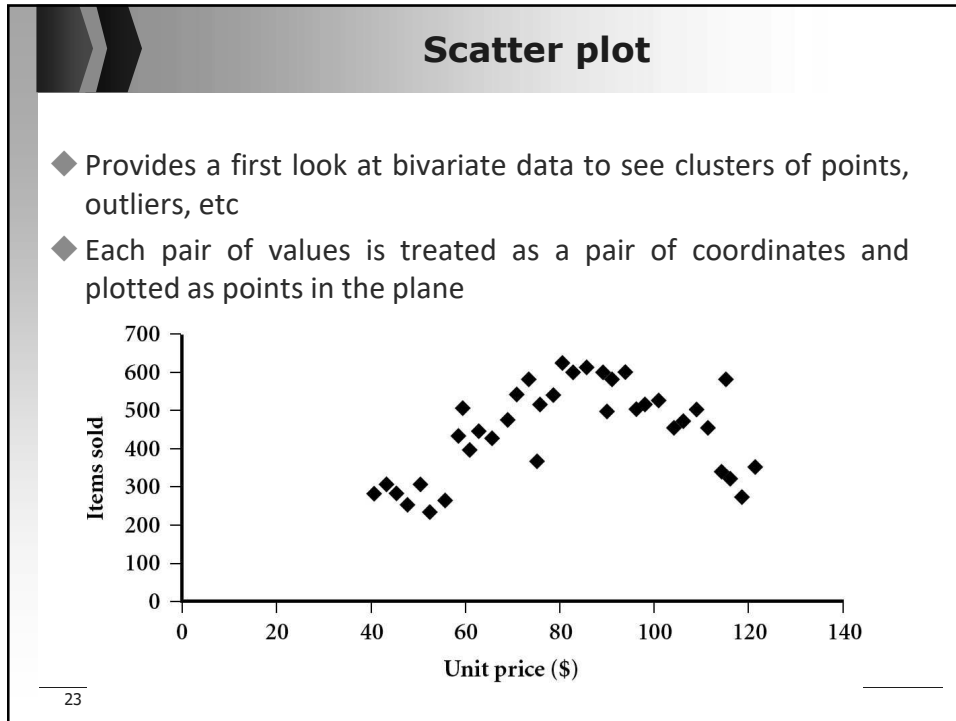
21

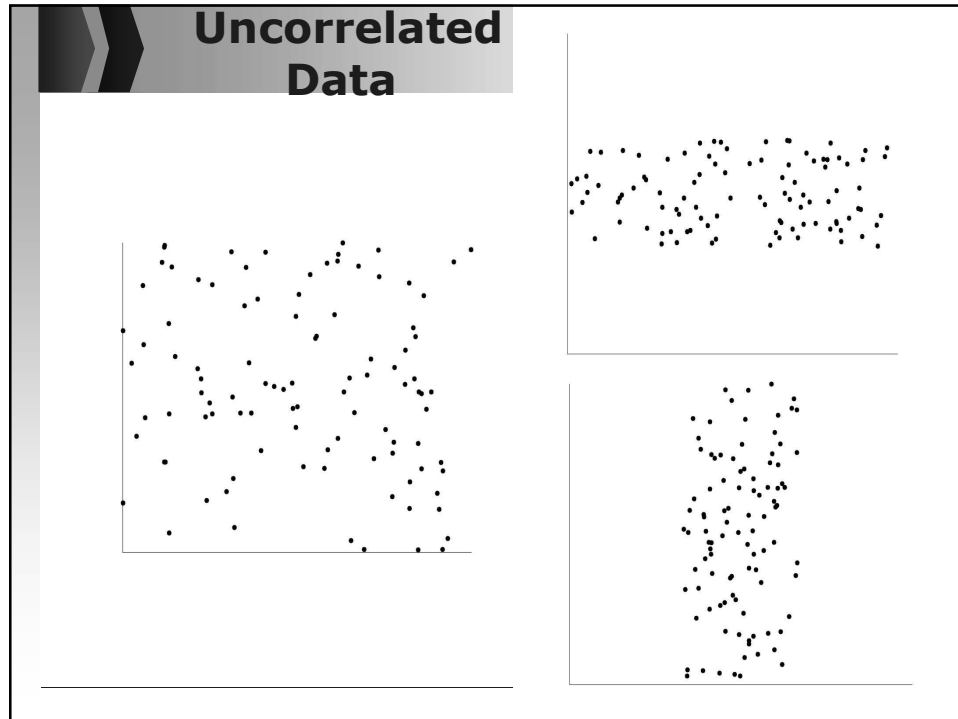
Quantile-Quantile (Q-Q) Plot

- ◆ Graphs the quantiles of one univariate distribution against the corresponding quantiles of another
- ◆ View: Is there 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.



22






Data Visualization


- ◆ Why data visualization?
 - Gain insight into an information space by mapping data onto graphical primitives
 - Provide qualitative overview of large data sets
 - Search for patterns, trends, structure, irregularities, relationships among data
 - Help find interesting regions and suitable parameters for further quantitative analysis
 - Provide a visual proof of computer representations derived
- ◆ Categorization of visualization methods:
 - Pixel-oriented visualization techniques
 - Geometric projection visualization techniques
 - Icon-based visualization techniques
 - Hierarchical visualization techniques
 - Visualizing complex data and relations

Pixel-Oriented Visualization Techniques

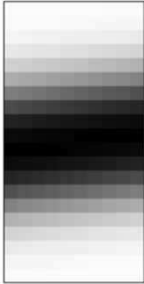
- ◆ For a data set of m dimensions, create m windows on the screen, one for each dimension
- ◆ The m dimension values of a record are mapped to m pixels at the corresponding positions in the windows
- ◆ The colors of the pixels reflect the corresponding values



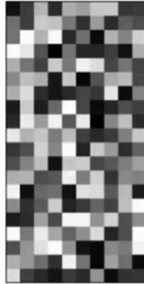
(a) Income



(b) Credit Limit



(c) transaction volume

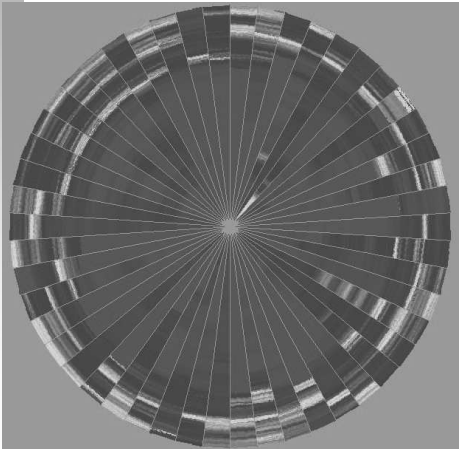


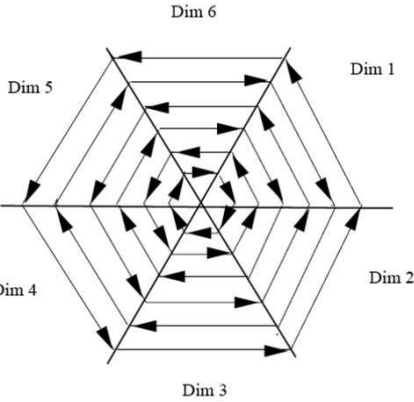
(d) age

27

Laying Out Pixels in Circle Segments

- ◆ To save space and show the connections among multiple dimensions, space filling is often done in a circle segment





(b) Laying out pixels in circle segment

Representing about 265,000 50-dimensional Data Items with the 'Circle Segments' Technique

28

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

29

Data Matrix and Dissimilarity Matrix

◆ Data matrix

- n data points with p dimensions
- Two modes

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

◆ Dissimilarity matrix

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

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

30

Proximity Measure for Nominal Attributes

- ◆ Can take 2 or more states, e.g., red, yellow, blue, green (generalization of a binary attribute)
- ◆ Method 1: Simple matching
 - m : # of matches, p : total # of variables
$$d(i, j) = \frac{p - m}{p}$$
- ◆ Method 2: Use a large number of binary attributes
 - creating a new binary attribute for each of the M nominal states

31

Proximity Measure for Binary Attributes

- ◆ A contingency table for binary data

		Object j		sum
		1	0	
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”:

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

32

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

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

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

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

33

Standardizing Numeric Data

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

- X: raw score to be standardized, μ : mean of the population, σ : standard deviation
- the distance between the raw score and the population mean in units of the standard deviation
- negative when the raw score is below the mean, "+" when above

- ### ◆ An alternative way: Calculate the mean absolute deviation

$$s_f = \frac{1}{n} (|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

where

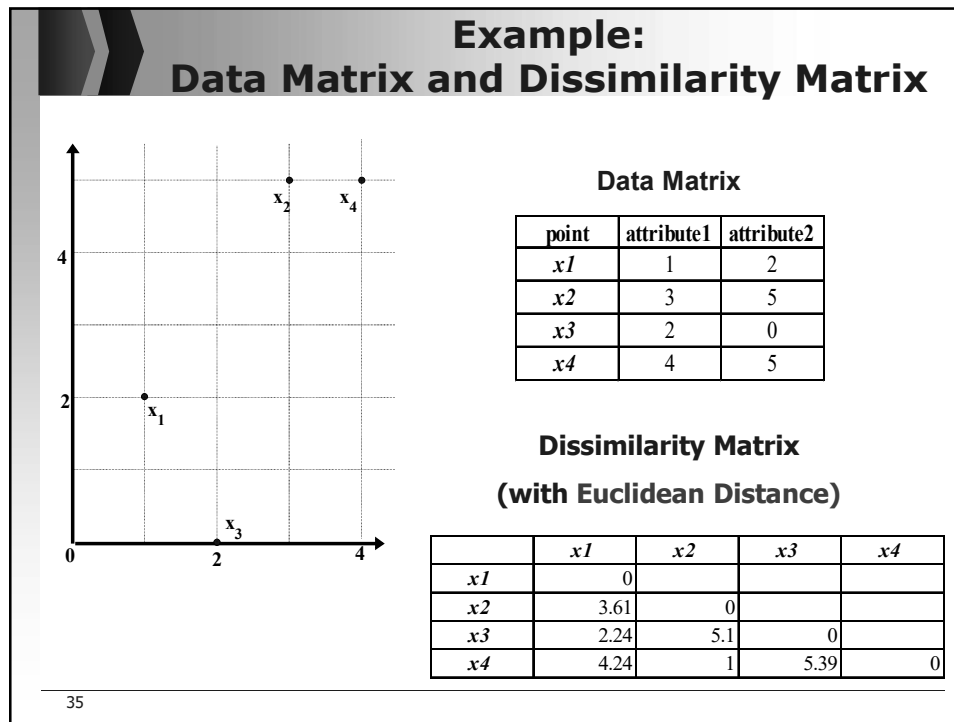
$$m_f = \frac{1}{n} (x_{1f} + x_{2f} + \dots + x_{nf})$$

- standardized measure (z-score):

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

- ### ◆ Using mean absolute deviation is more robust than using standard deviation

34



Distance on Numeric Data: Minkowski Distance

- ◆ *Minkowski distance*: A popular distance measure

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

where $i = (x_{i1}, x_{i2}, \dots, x_{ip})$ and $j = (x_{j1}, x_{j2}, \dots, x_{jp})$ 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

36

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

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

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

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

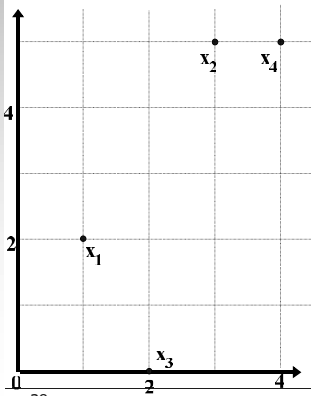
- ◆ $h \rightarrow \infty$. "supremum" (L_{\max} norm, L_∞ norm) distance.
 - This is the maximum difference between any component (attribute) of the vectors

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

37

Example: Minkowski Distance

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



38

Dissimilarity Matrices

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_∞	x1	x2	x3	x4
x1	0			
x2	3	0		
x3	2	5	0	
x4	3	1	5	0

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 $r_{if} \in \{1, \dots, M_f\}$
 - map the range of each variable onto $[0, 1]$ by replacing i -th object in the f -th variable by

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

- compute the dissimilarity using methods for interval-scaled variables

39

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

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

- f is binary or nominal:
 - $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$, or $d_{ij}^{(f)} = 1$ otherwise
- f is numeric: use the normalized distance
- f is ordinal
 - Compute ranks r_{if} and
 - Treat z_{if} as interval-scaled

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

40

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



Questions?