

# 6. DATA ANALYSIS & DATA MANAGEMENT



# Copyright Notice

2

- The majority of the information is from:
- Data Mining: Concepts and Techniques
  - ▶ Third Edition
  - ▶ ©2011 Han, Kamber, and Pei

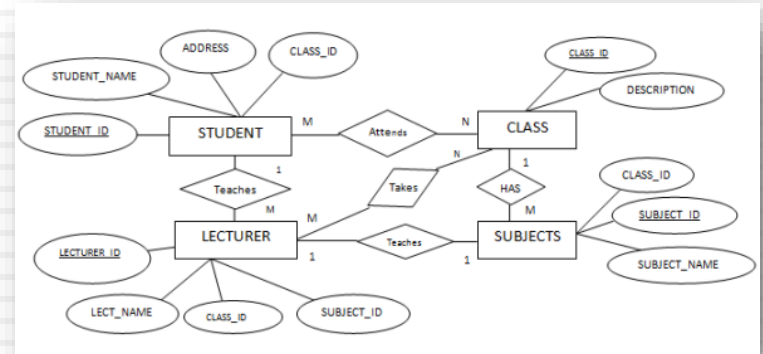
# 6.1: Data Objects and Attribute Types

## 6.1: Data Objects and Attribute Types

## 6.2: Basic Statistical Descriptions of Data

## 6.3: Data Visualization

## 6.4: Data Mining



# Learning Objectives

4

- List types of data sets
- Describe important characteristics of structured data
- Define data object
- Understand attributes and its types

# Types of Data Sets (1)

- Record
  - Relational records
  - Data matrix, e.g., numerical matrix, crosstabs
  - Document data: text documents: term-frequency vector
  - Transaction data

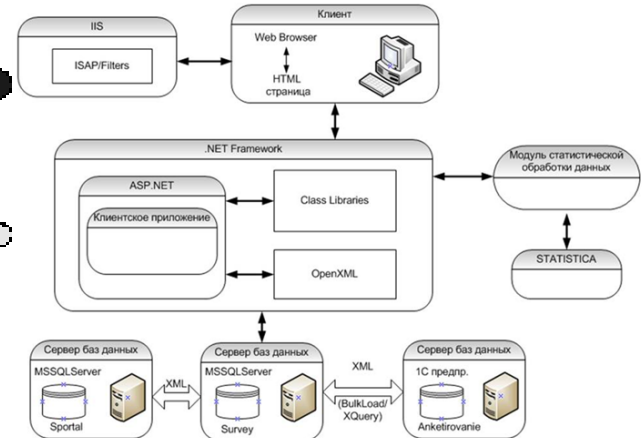
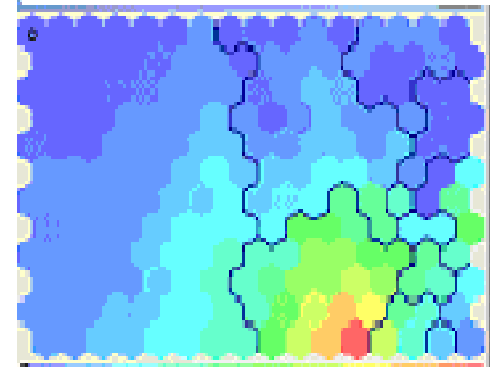
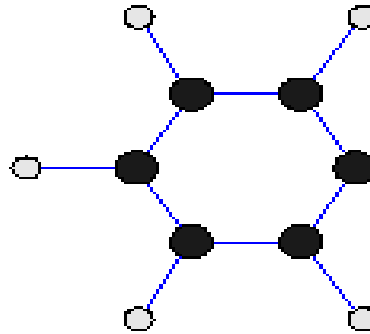
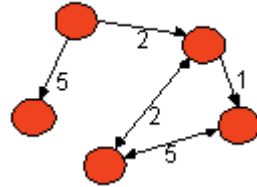
	team	coach	play	ball	score	game	win	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0

<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

# Types of Data Sets (2)

6

- 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:



# Important Characteristics of Structured Data

7

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

# Data Objects

8

- 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.



# Attributes

9

- **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

10

- **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

# Numeric Attribute Types

11

- 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

## □ **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

# Summary

13

- Data attribute types
  - ▣ nominal, binary, ordinal, interval-scaled, ratio-scaled
- Many types of data sets
  - ▣ e.g., numerical, text, graph, Web, image

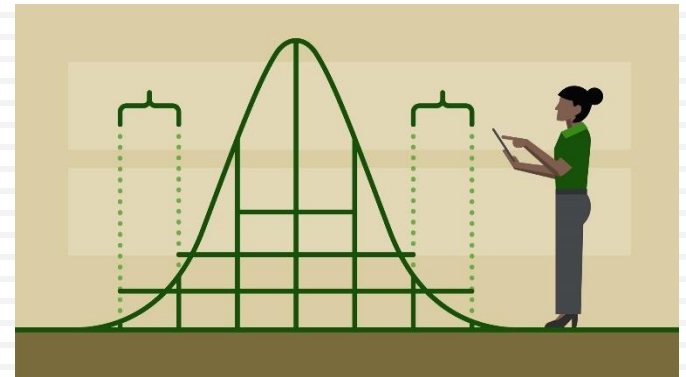
## 6.2: Basic Statistical Descriptions of Data

6.1: Data Objects and Attribute Types

**6.2: Basic Statistical Descriptions of Data**

6.3: Data Visualization

6.4: Data Mining



# Learning Objectives

15

- Define measuring the central tendency
- Describe measuring the dispersion of data
- Understand properties of normal dispersion curve
- List graphic displays of basic statistical descriptions

# Measuring the Central Tendency

16

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

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

- Weighted arithmetic mean:
- 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*):

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

□ Mode

- Value that occurs most frequently in the data
- Unimodal, bimodal, trimodal

□ Empirical formula:  $mean - mode = 3 \times (mean - median)$

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

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

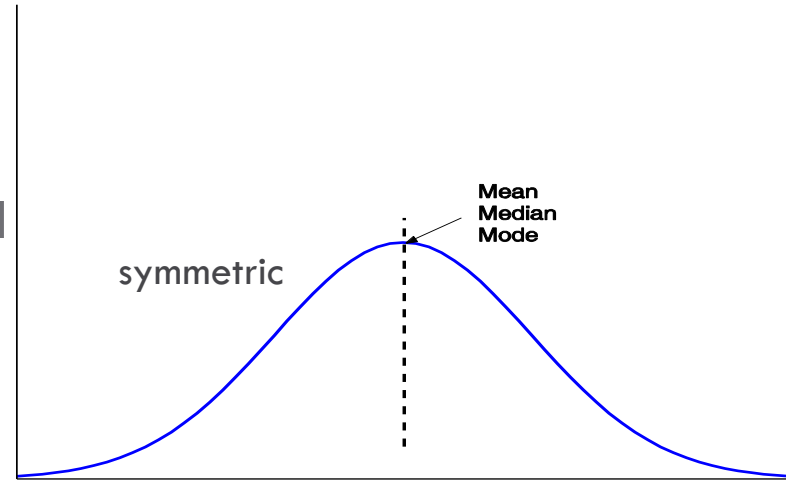
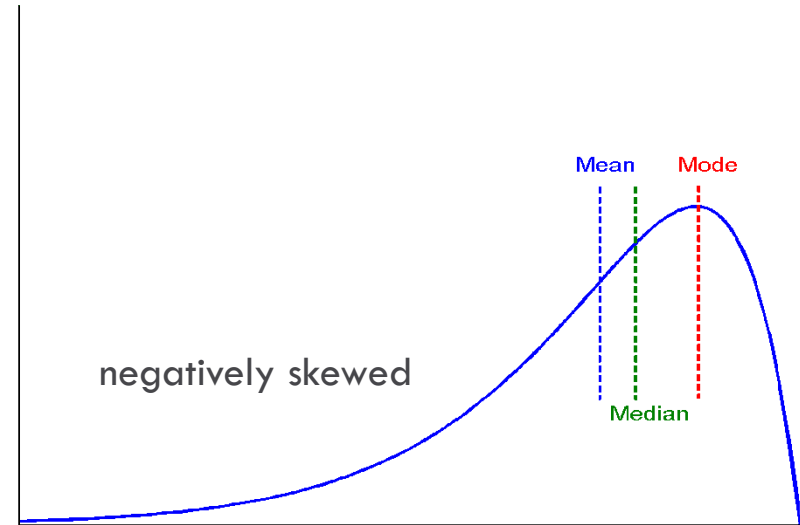
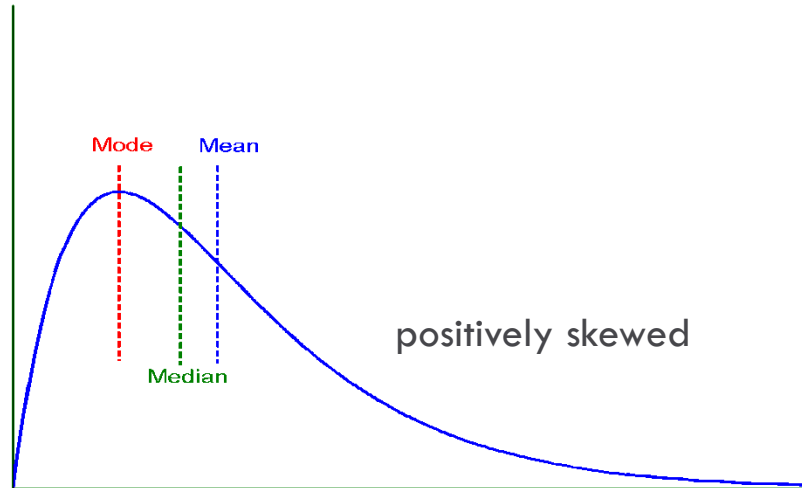
age	frequency
1-5	200
6-15	450
16-20	300
21-50	1500
51-80	700
81-110	44



# Symmetric vs. Skewed Data

17

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



# Measuring the Dispersion of Data

18

- Quartiles, outliers and boxplots
  - ▣ **Quartiles:**  $Q_1$  (25<sup>th</sup> percentile),  $Q_3$  (75<sup>th</sup> 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:*

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

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{N} \sum_{i=1}^n x_i^2 - \mu^2$$

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

# Boxplot Analysis

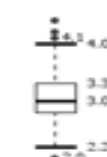
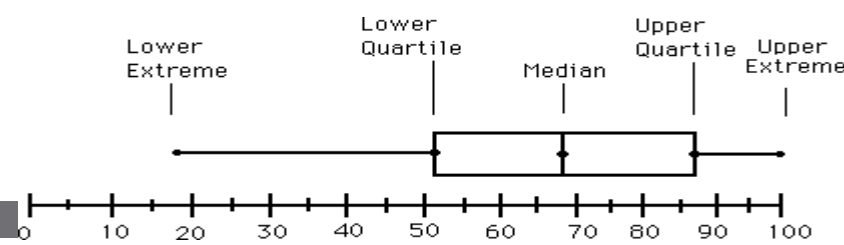
19

- **Five-number summary** of a distribution

- Minimum, Q1, Median, Q3, Maximum

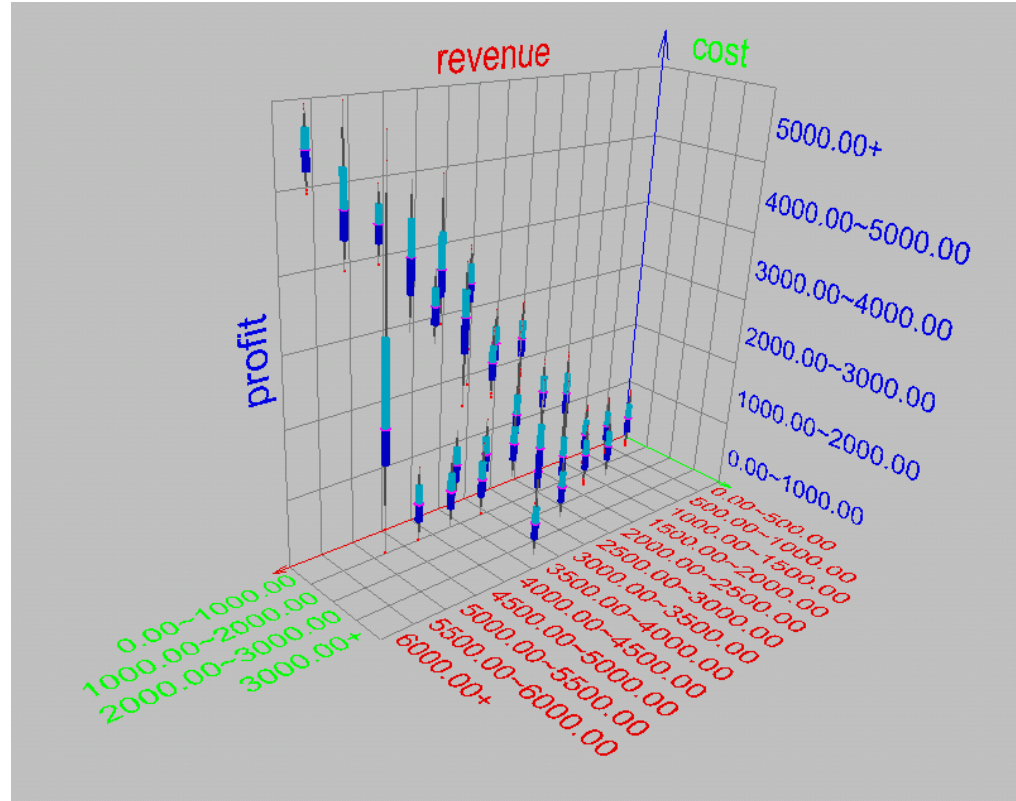
- **Boxplot**

- 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

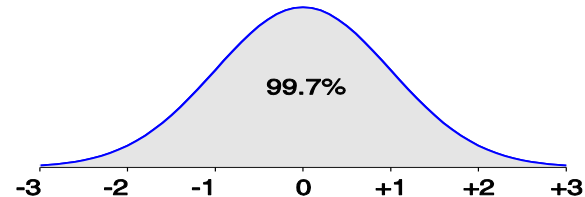
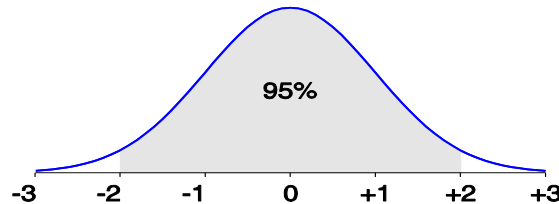
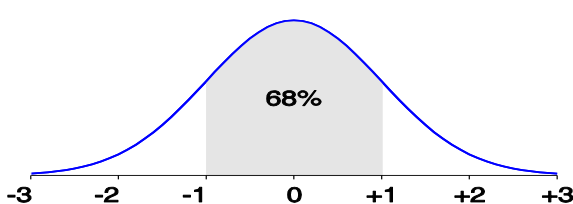
20



# Properties of Normal Distribution Curve

21

- 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

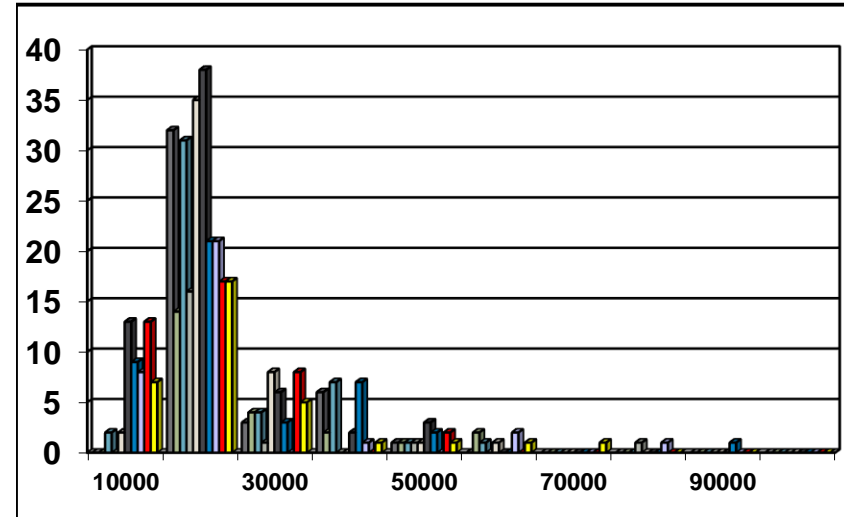
22

- **Boxplot:** graphic display of five-number summary
- **Histogram:** x-axis contains values, y-axis represents 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

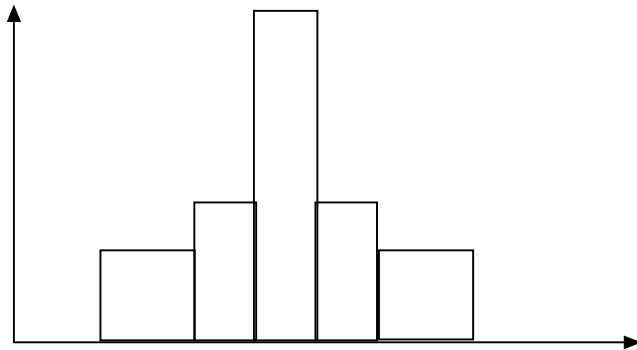
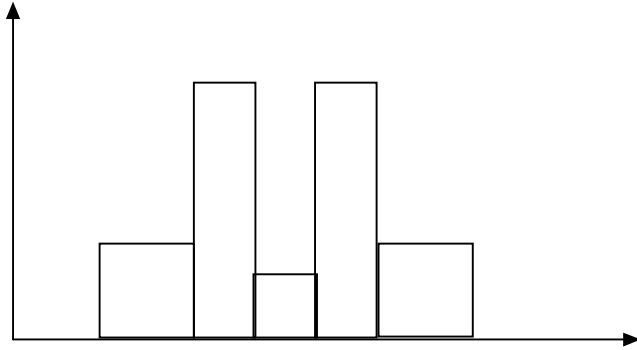
23

- **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

24



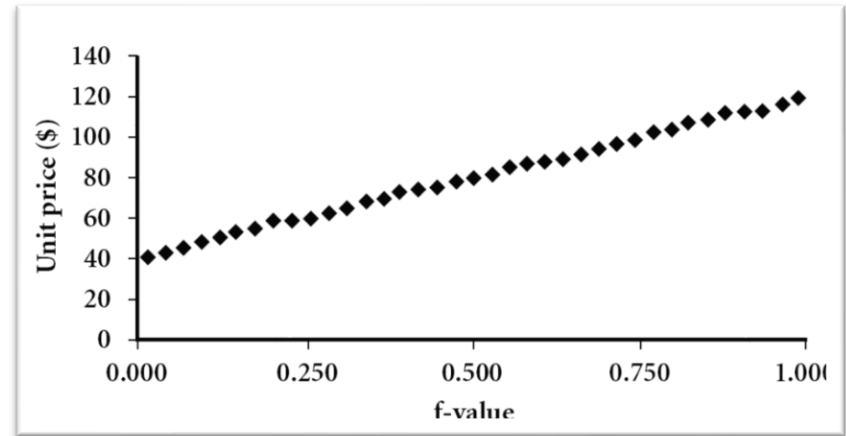
- The two histograms shown in the left may have the same boxplot representation
  - ▣ The same values for: min, Q1, median, Q3, max
- But they have rather different data distributions



# Quantile Plot

25

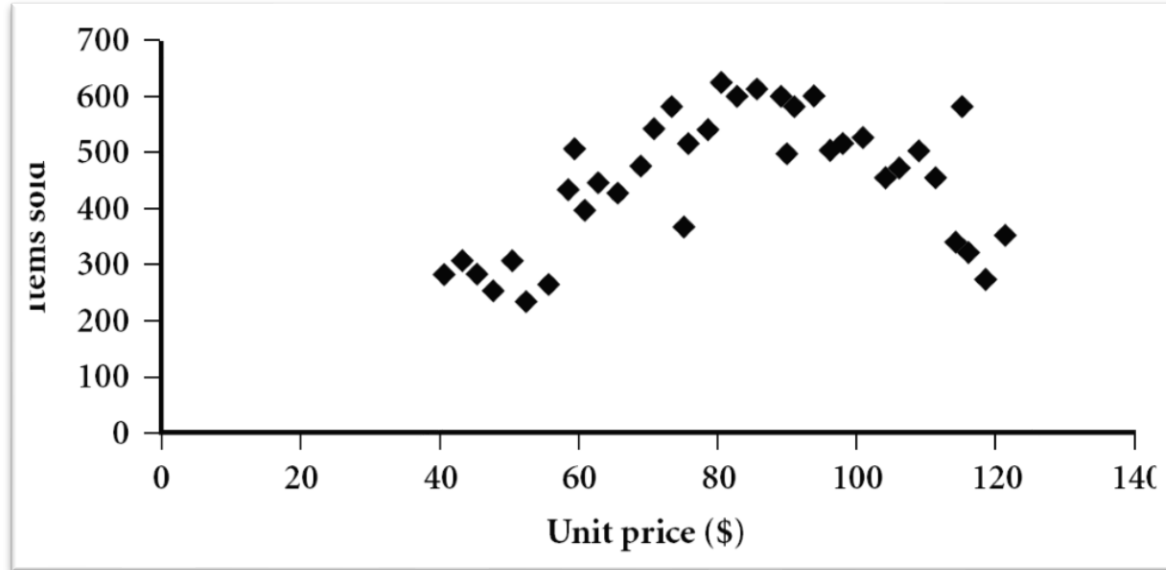
- Displays all of the data (allowing the user to assess both the overall behavior and unusual occurrences)
- Plots quantile information
  - For a data  $x_i$  data sorted in increasing order,  $f_i$  indicates that approximately  $100 f_i\%$  of the data are below or equal to the value  $x_i$



# Scatter plot

26

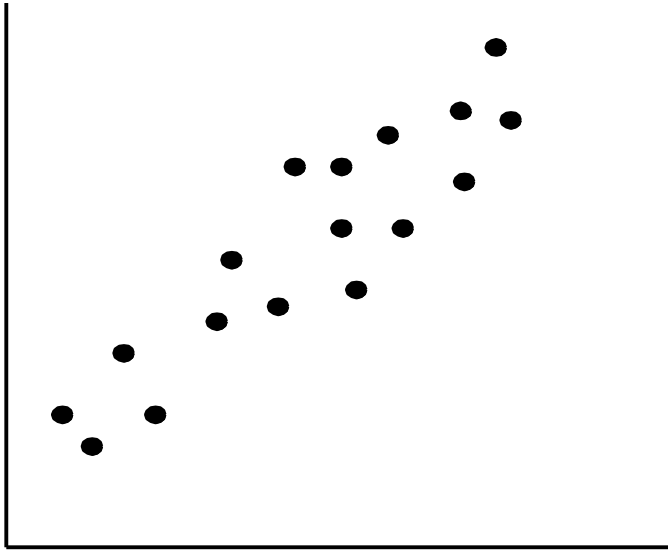
- Provides a first look at bivariate data to see clusters of points, outliers, etc.
- Each pair of values is treated as a pair of coordinates and plotted as points in the plane



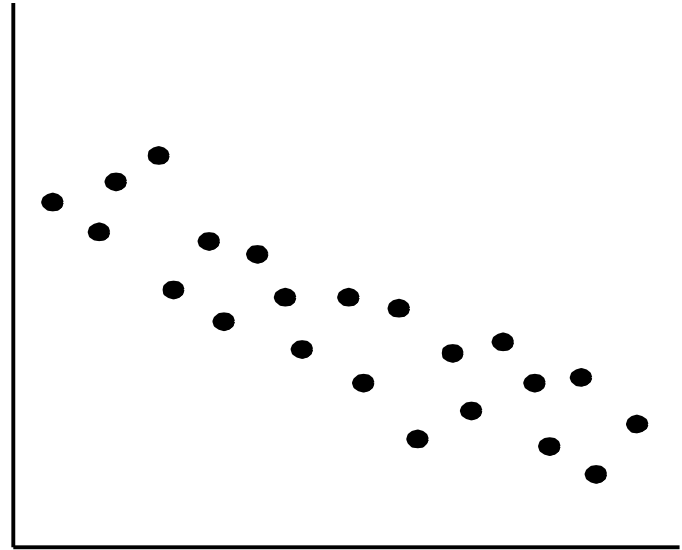
# Positively and Negatively Correlated Data

27

Positively correlated data



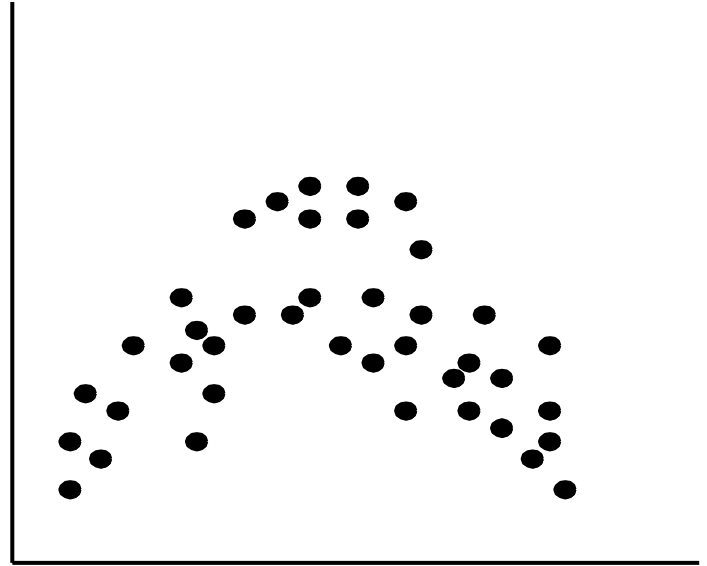
Negative correlated



# Curvilinear relationships

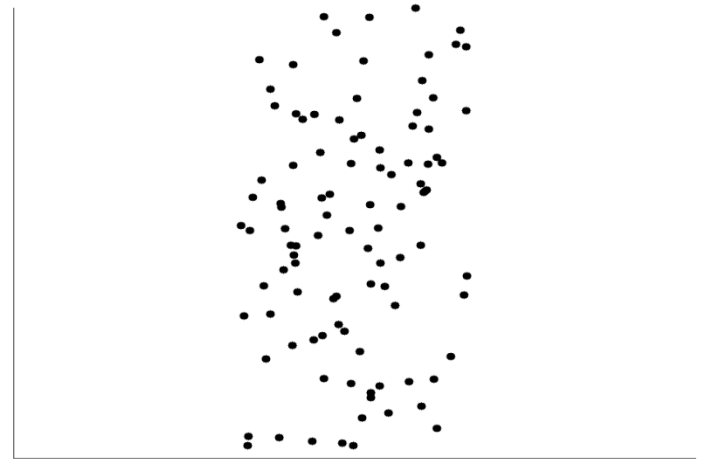
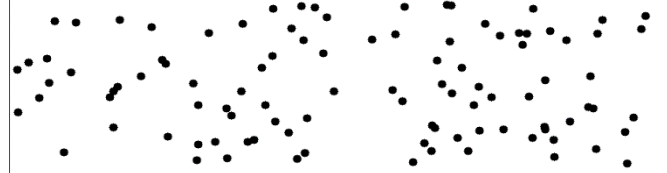
28

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



# Uncorrelated Data

29



# Summary

30

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

## 6.3: Data Visualization

6.1: Data Objects and Attribute Types

6.2: Basic Statistical Descriptions of Data

**6.3: Data Visualization**

6.4: Data Mining



# Learning Objectives

32

- Describe data visualization techniques:
  - pixel-oriented,
  - geometric projection,
  - icon-based,
  - hierarchical,
  - visualizing complex data and relations.



# Data Visualization

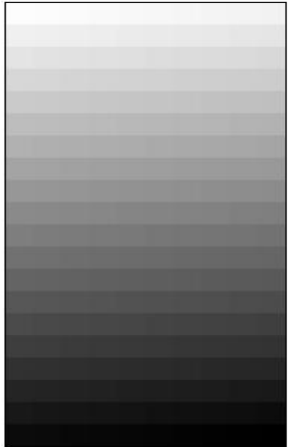
33

- Categorization of visualization methods:
  - Pixel-oriented visualization techniques
  - Geometric projection visualization techniques
  - Icon-based visualization techniques
  - Hierarchical visualization techniques
  - Visualizing complex data and relations

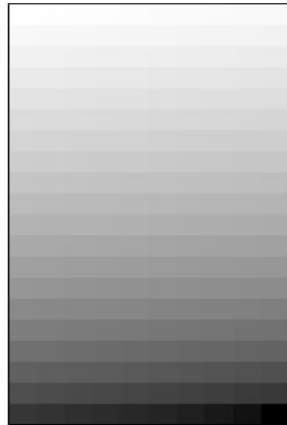
# Pixel-Oriented Visualization Techniques

34

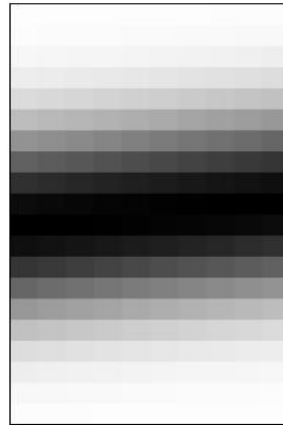
- For a data set of  $m$  dimensions, create  $m$  windows on the screen, one for each dimension
- The  $m$  dimension values of a record are mapped to  $m$  pixels at the corresponding positions in the windows
- The colors of the pixels reflect the corresponding values



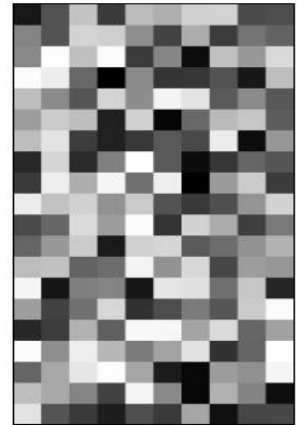
(a) Income



(b) Credit Limit



(c) transaction volume

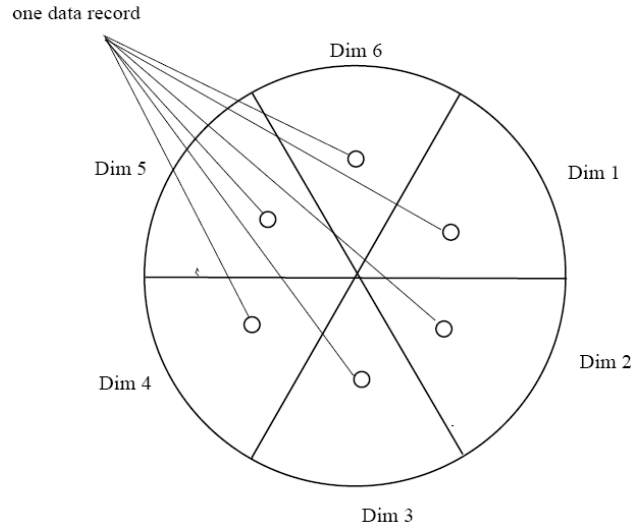


(d) age

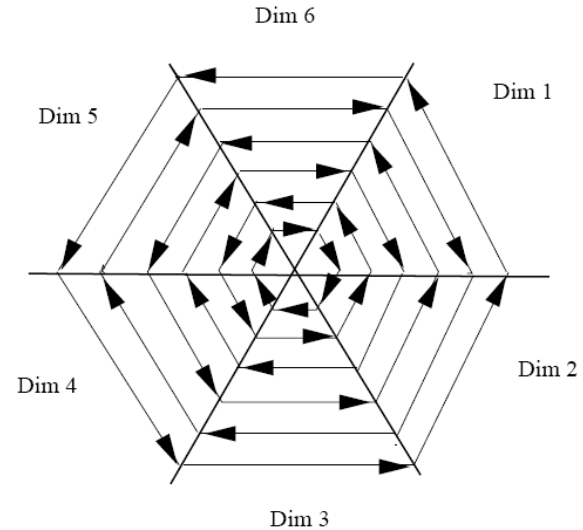
# Laying Out Pixels in Circle Segments

35

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



(a) Representing a data record in circle segment



(b) Laying out pixels in circle segment

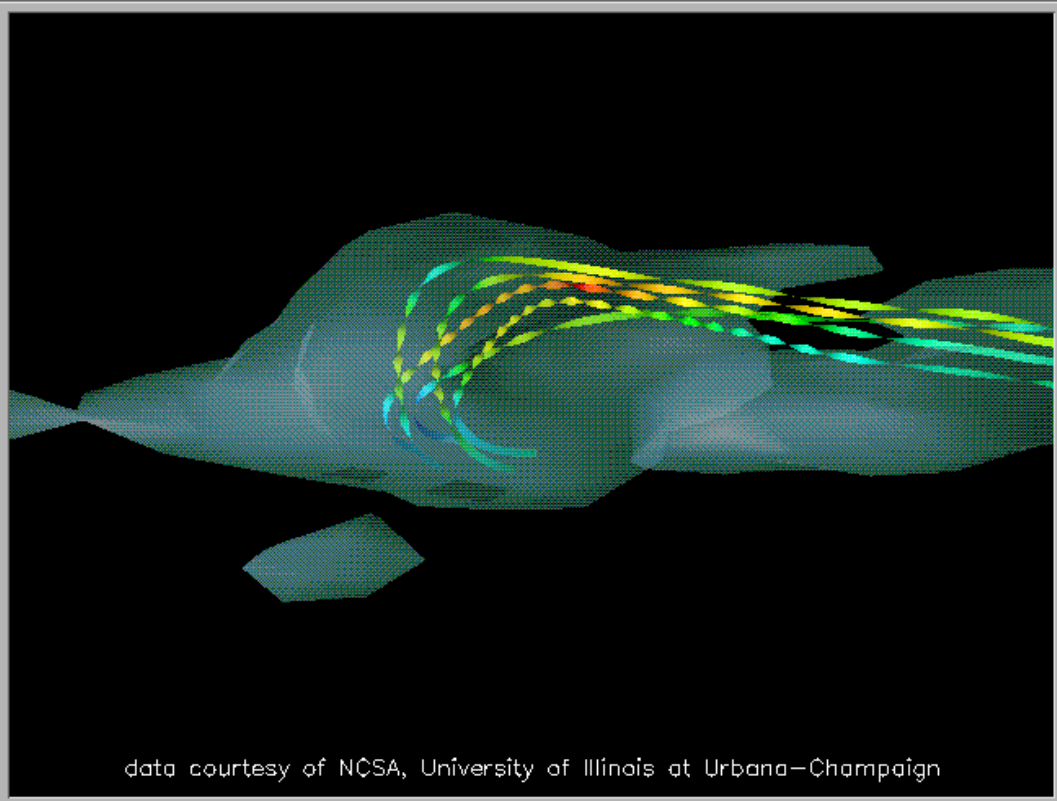
# Geometric Projection Visualization Techniques

36

- Visualization of geometric transformations and projections of the data
- Methods
  - ▣ Direct visualization
  - ▣ Scatterplot and scatterplot matrices
  - ▣ Landscapes
- Methods (2)
  - ▣ Projection pursuit technique: Help users find meaningful projections of multidimensional data
  - ▣ Prosection views
  - ▣ Hyperslice
  - ▣ Parallel coordinates

# Direct Data Visualization

37

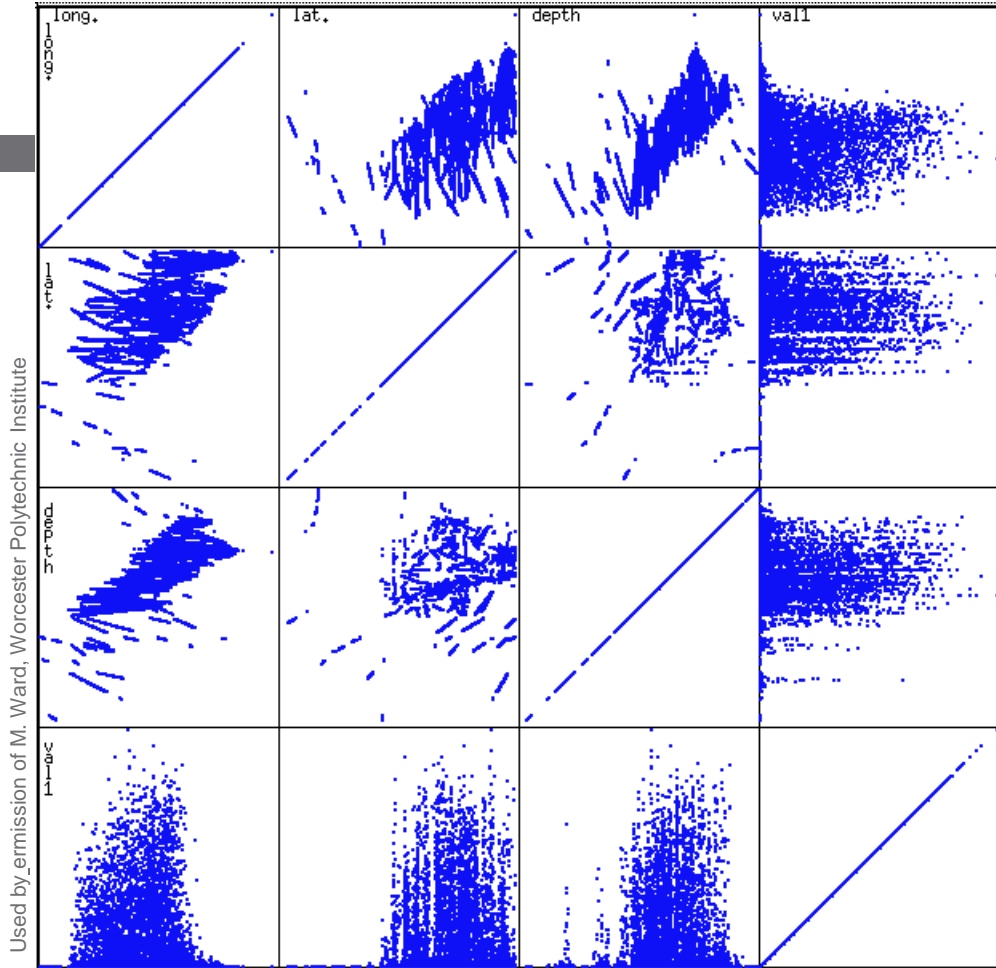


Ribbons with Twists Based  
on Vorticity

# Scatterplot Matrices

38

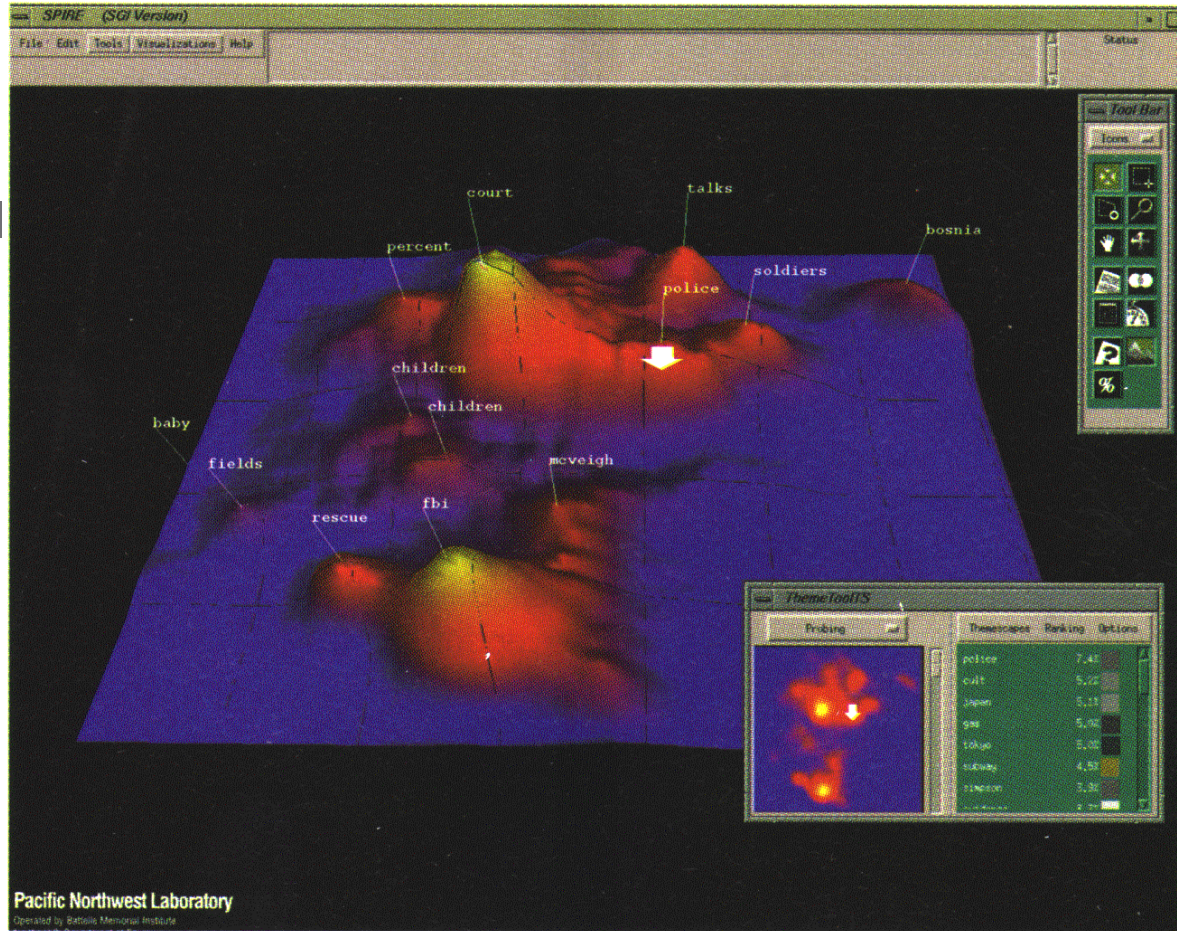
- Matrix of scatterplots (x-y-diagrams) of the k-dim. data [total of  $(k^2/2-k)$  scatterplots]



# Landscapes

39

- Visualization of the data as perspective landscape
- The data needs to be transformed into a (possibly artificial) 2D spatial representation which preserves the characteristics of the data

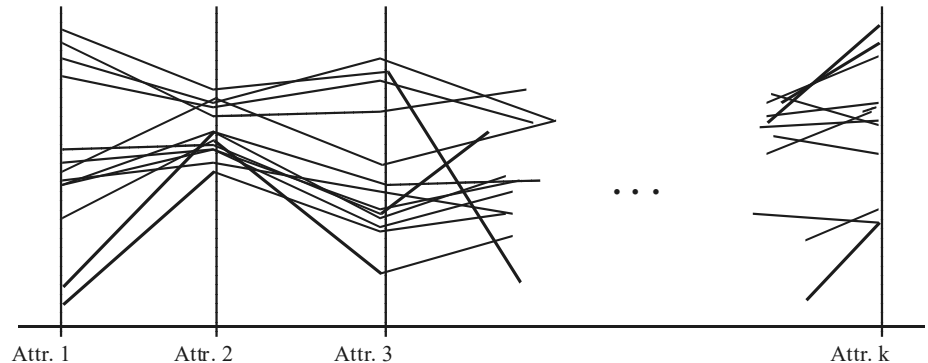


news articles visualized as a landscape

# Parallel Coordinates

40

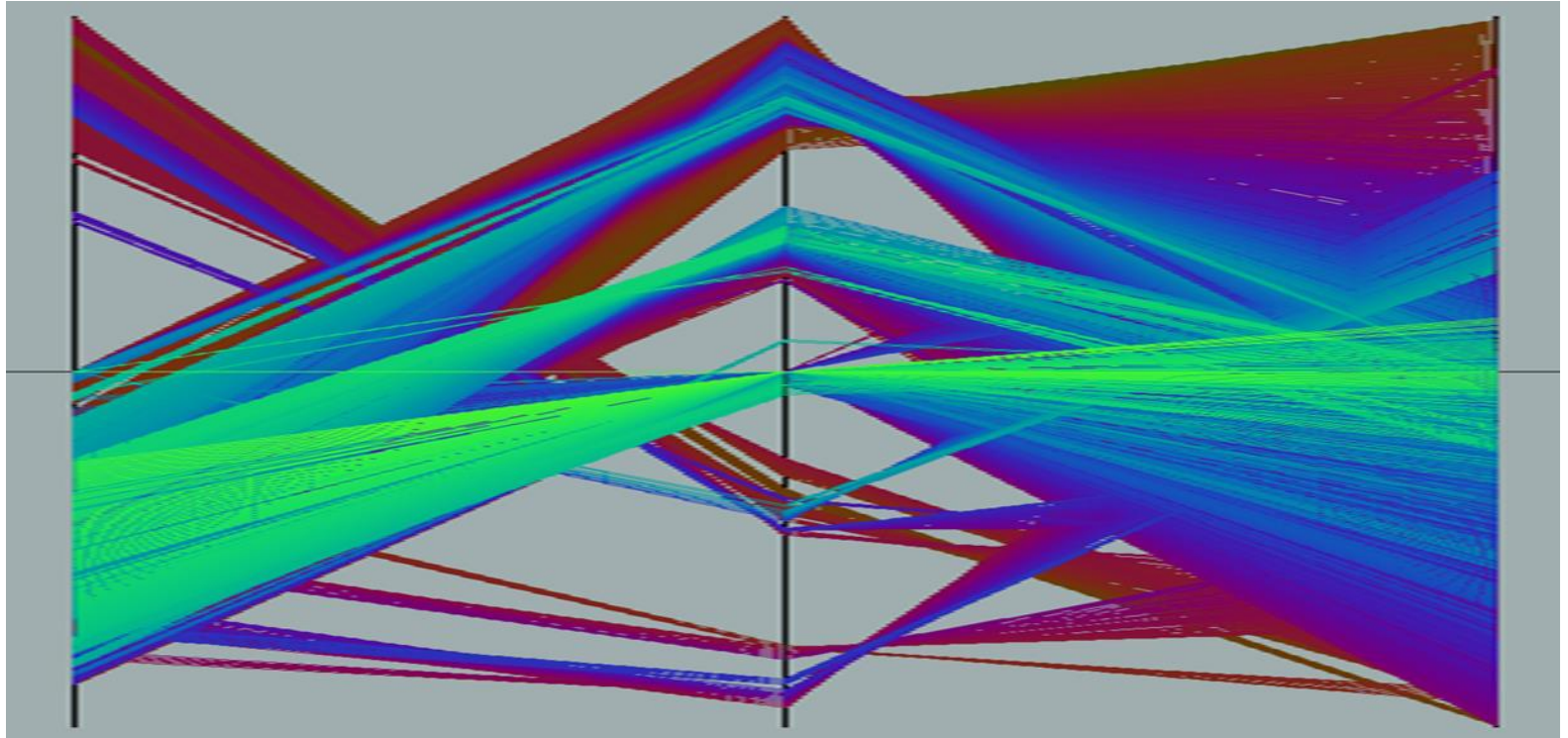
- n equidistant axes which are parallel to one of the screen axes and correspond to the attributes
- The axes are scaled to the [minimum, maximum]: range of the corresponding attribute
- Every data item corresponds to a polygonal line which intersects each of the axes at the point which corresponds to the value for the attribute





# Parallel Coordinates of a Data Set

41



# Icon-Based Visualization Techniques

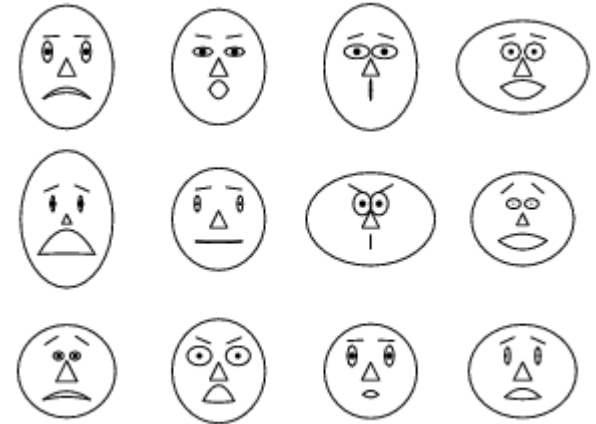
42

- Visualization of the data values as features of icons
- Typical visualization methods
  - ▣ Chernoff Faces
  - ▣ Stick Figures
- General techniques
  - ▣ Shape coding: Use shape to represent certain information encoding
  - ▣ Color icons: Use color icons to encode more information
  - ▣ Tile bars: Use small icons to represent the relevant feature vectors in document retrieval

# Chernoff Faces

43

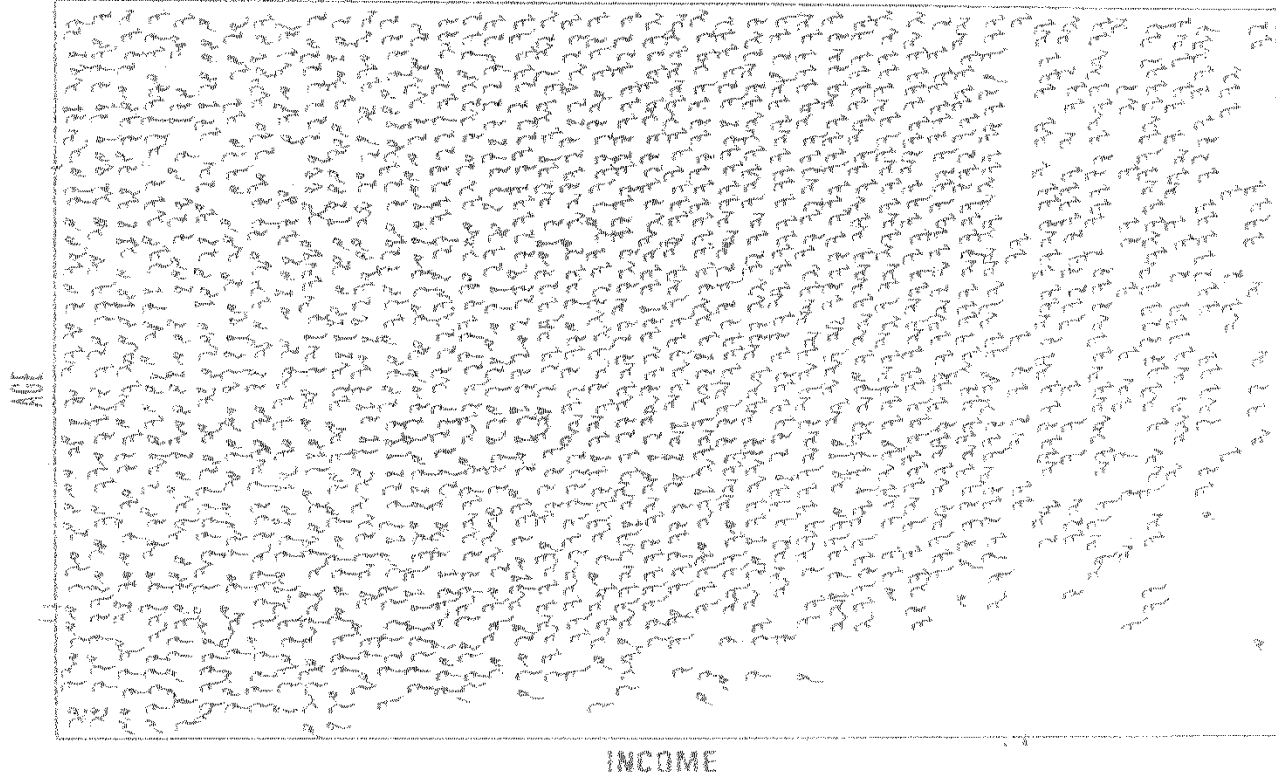
- A way to display variables on a two-dimensional surface, e.g., let  $x$  be eyebrow slant,  $y$  be eye size,  $z$  be nose length, etc.
- The figure shows faces produced using 10 characteristics--head eccentricity, eye size, eye spacing, eye eccentricity, pupil size, eyebrow slant, nose size, mouth shape, mouth size, and mouth opening): Each assigned one of 10 possible values, generated using [Mathematica](#) (S. Dickson)



# Stick Figure

44

used by permission of G. Grinstein, University of Massachusetts  
at Lowell



A census data figure showing age, income, gender, education, etc.

A 5-piece stick figure (1 body and 4 limbs w. different angle/length)

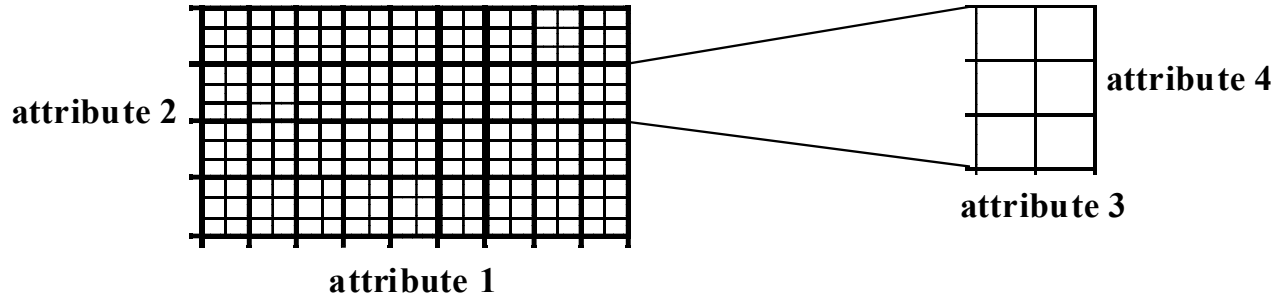
# Hierarchical Visualization Techniques

45

- Visualization of the data using a hierarchical partitioning into subspaces
- Methods
  - Dimensional Stacking
  - Worlds-within-Worlds
  - Tree-Map
  - Cone Trees
  - InfoCube

# Dimensional Stacking

46

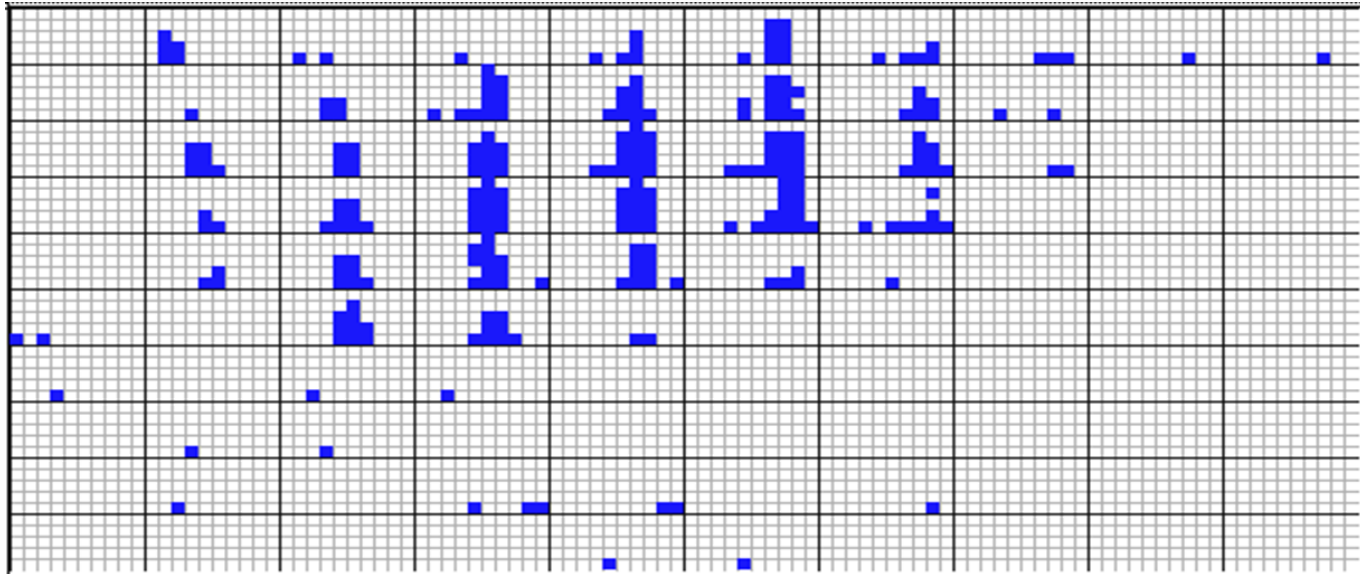


- Partitioning of the n-dimensional attribute space in 2-D subspaces, which are 'stacked' into each other
- Partitioning of the attribute value ranges into classes. The important attributes should be used on the outer levels.
- Adequate for data with ordinal attributes of low cardinality
- But, difficult to display more than nine dimensions
- Important to map dimensions appropriately

# Dimensional Stacking

47

- Visualization of oil mining data with longitude and latitude mapped to the outer x-, y-axes and ore grade and depth mapped to the inner x-, y-axes



# Worlds-within-Worlds

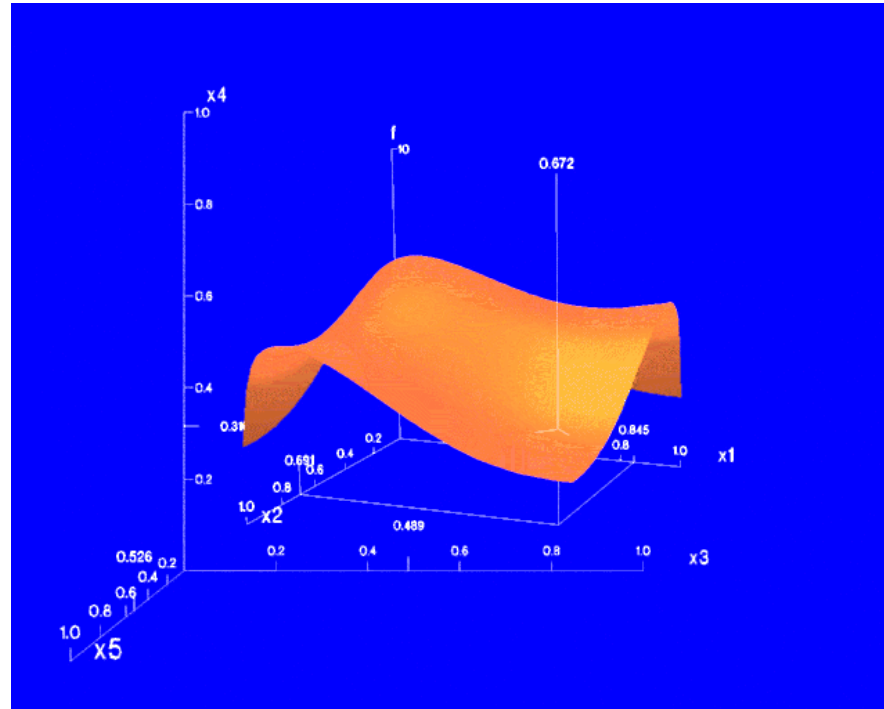
48

- Assign the function and two most important parameters to innermost world
- Fix all other parameters at constant values - draw other (1 or 2 or 3-dimensional worlds choosing these as the axes)
- Software that uses this paradigm
- N-vision: Dynamic interaction through data glove and stereo displays, including rotation, scaling (inner) and translation (inner/outer)
- Auto Visual: Static interaction using queries



# Worlds-within-Worlds (2)

49

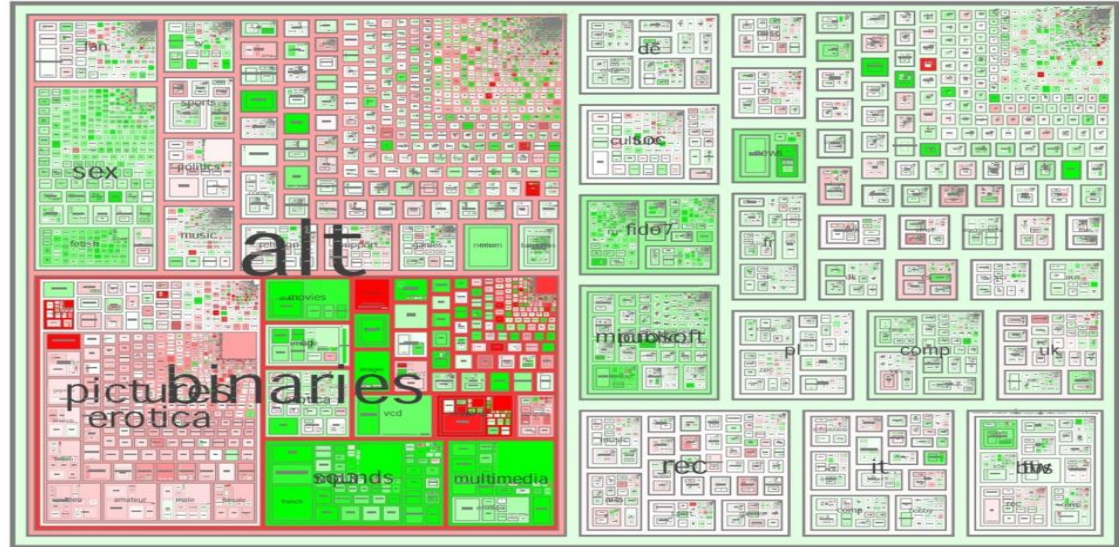


# Tree-Map

50

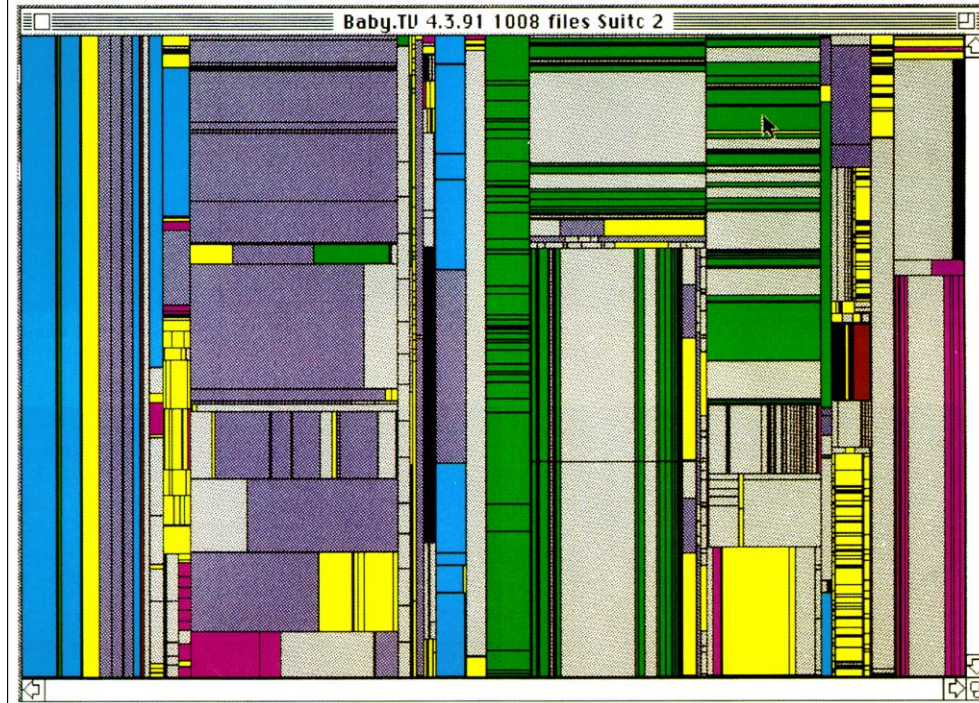
- Screen-filling method which uses a hierarchical partitioning of the screen into regions depending on the attribute values
- The x- and y-dimension of the screen are partitioned alternately according to the attribute values (classes)

MSR Netscan Image



# Tree-Map of a File System (Schneiderman)

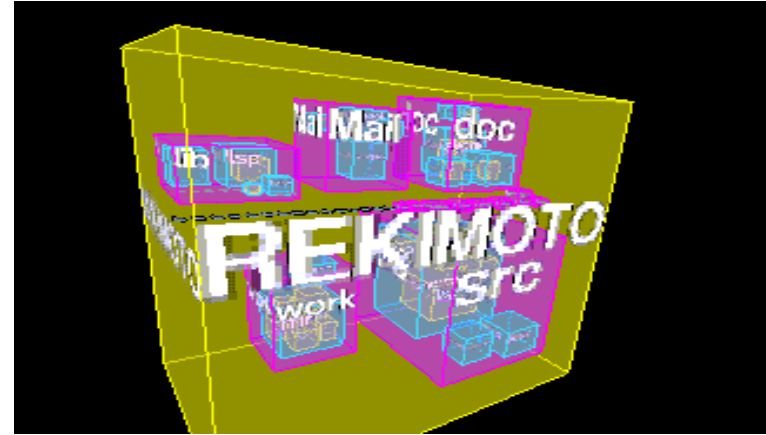
51



# InfoCube

52

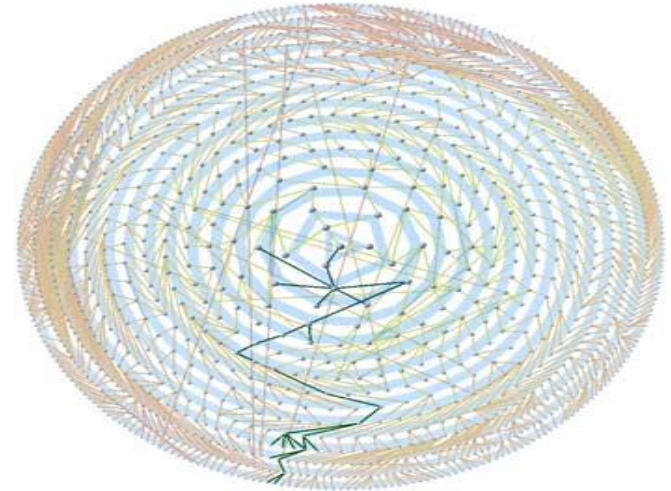
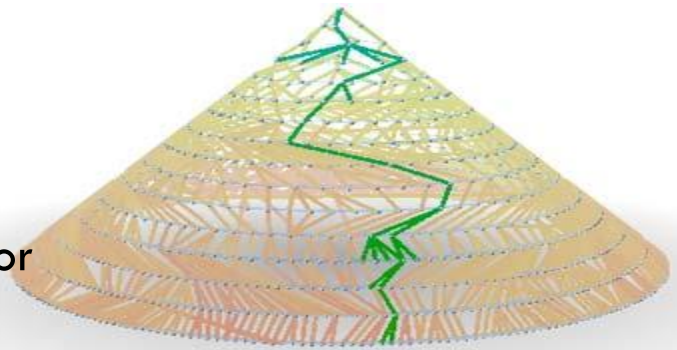
- A 3-D visualization technique where hierarchical information is displayed as nested semi-transparent cubes
- The outermost cubes correspond to the top-level data, while the subnodes or the lower level data are represented as smaller cubes inside the outermost cubes, and so on



# Three-D Cone Trees

53

- 3D cone tree visualization technique works well for up to a thousand nodes or so
- First, build a 2D circle tree that arranges its nodes in concentric circles centered on the root node
- Cannot avoid overlaps when projected to 2D
- G. Robertson, J. Mackinlay, S. Card. “Cone Trees: Animated 3D Visualizations of Hierarchical Information”, ACM SIGCHI'91
- Graph from Nadeau Software Consulting website: Visualize a social network dataset that models the way an infection spreads from one person to the next



# Visualizing Complex Data and Relations

- Visualizing non-numerical data: text and social networks
- Tag cloud: visualizing user-generated tags
  - The importance of tag is represented by font size/color
  - Besides text data, there are also methods to visualize relationships, such as visualizing social networks



Newsmap: Google News Stories

# Summary

55

- 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

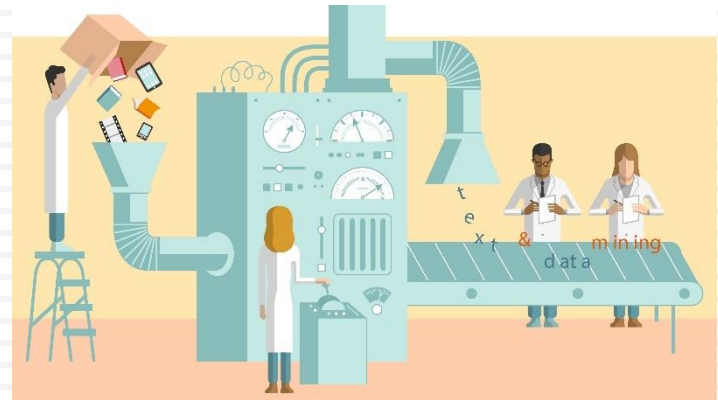
## 6.4: Data Mining

6.1: Data Objects and Attribute Types

6.2: Basic Statistical Descriptions of Data

6.3: Data Visualization

**6.4: Data Mining**





# Learning Objectives

57

- Define data mining
- Understand knowledge discovery process
- List data mining functions
- Explain classification, cluster analysis
- Describe applications of data mining

# Trends leading to Big Data

58

- The Explosive Growth of Data: from terabytes to petabytes
  - ▣ Data collection and data availability
    - Automated data collection tools, database systems, Web, computerized society
  - ▣ Major sources of abundant data
    - Business: Web, e-commerce, transactions, stocks, ...
    - Science: Remote sensing, bioinformatics, scientific simulation, ...
    - Society and everyone: news, digital cameras, YouTube
- We are drowning in data, but starving for knowledge!
- “Necessity is the mother of invention”—**Data mining**—Automated analysis of massive data sets

# Evolution of Database Technology

59

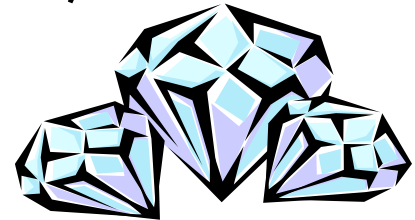
- 1960s:
  - ▣ Data collection, database creation, IMS and network DBMS
- 1970s:
  - ▣ Relational data model, relational DBMS implementation
- 1980s:
  - ▣ RDBMS, advanced data models (extended-relational, OO, deductive, etc.)
  - ▣ Application-oriented DBMS (spatial, scientific, engineering, etc.)
- 1990s:
  - ▣ Data mining, data warehousing, multimedia databases, and Web databases
- 2000s
  - ▣ Stream data management and mining
  - ▣ Data mining and its applications
  - ▣ Web technology (XML, data integration) and global information systems

# Definition of Data Mining



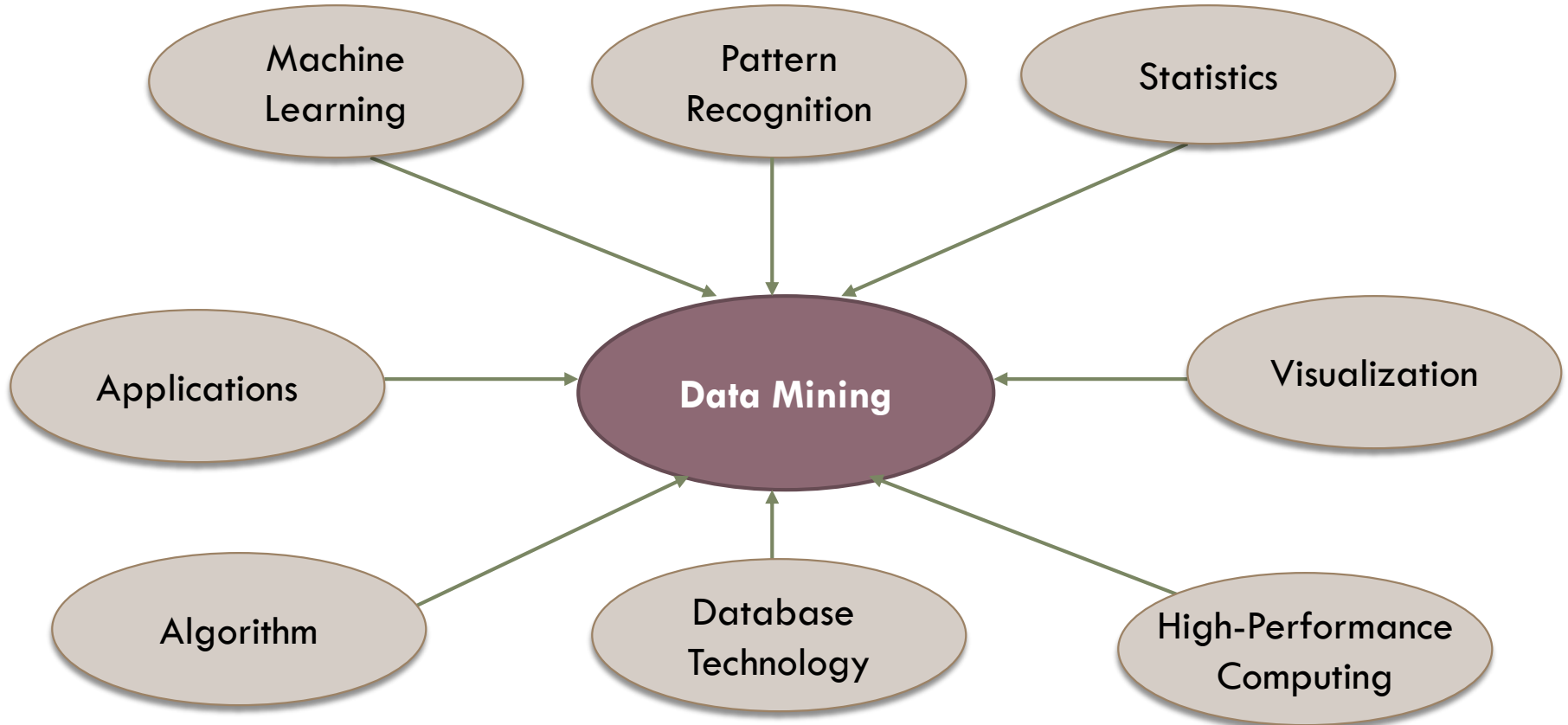
60

- Data mining (knowledge discovery from data)
  - ▣ Extraction of interesting (non-trivial, implicit, previously unknown and potentially useful) patterns or knowledge from huge amount of data
  
- Alternative names
  - ▣ Knowledge discovery (mining) in databases (KDD), knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence, etc.



# Data Mining: Confluence of Multiple Disciplines (1)

61



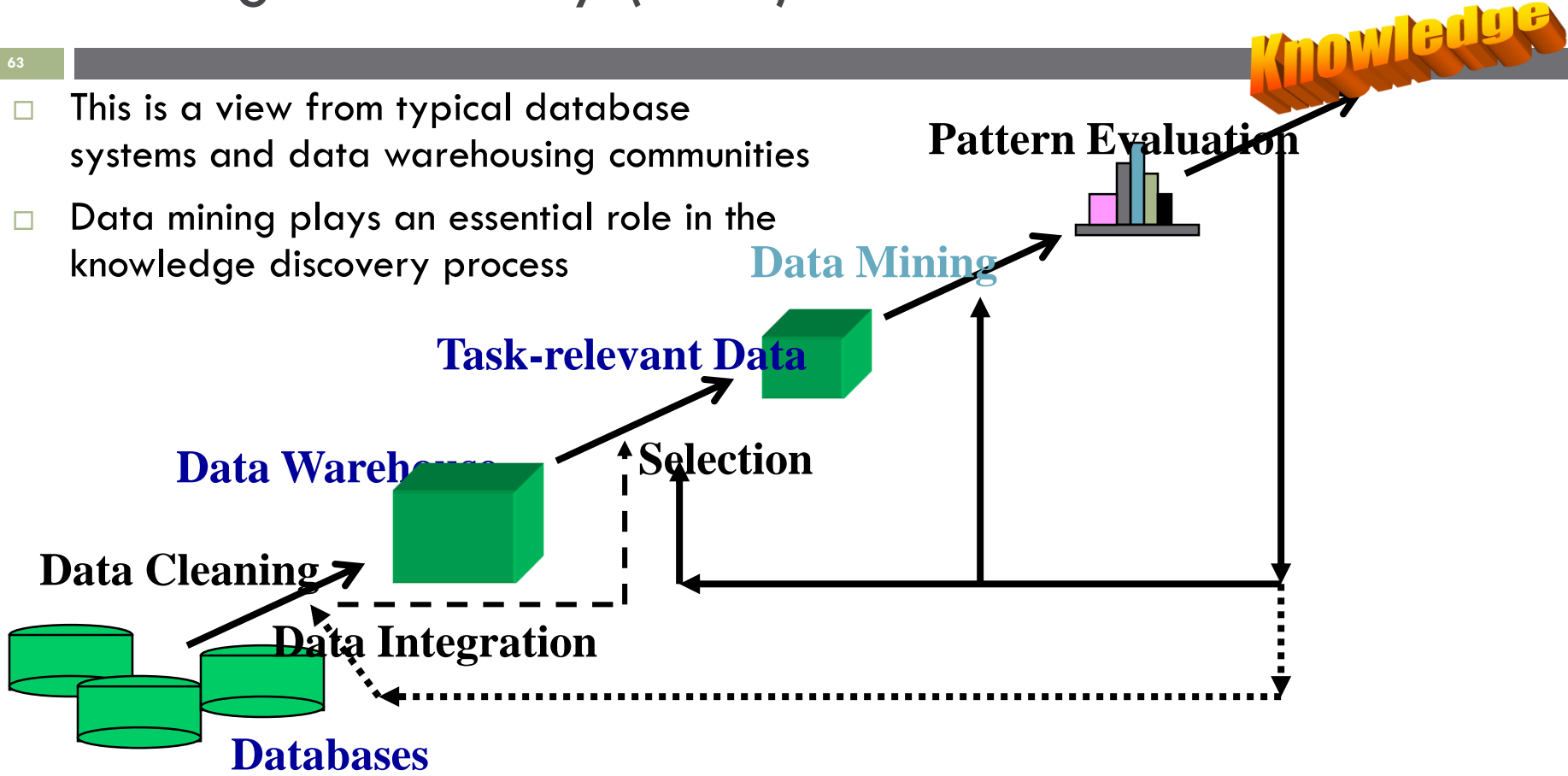
# Data Mining: Confluence of Multiple Disciplines (2)

62

- Tremendous amount of data
  - ▣ Algorithms must be highly scalable to handle such as terabytes of data
- High-dimensionality of data
  - ▣ Micro-array may have tens of thousands of dimensions
- High complexity of data
  - ▣ Data streams and sensor data
  - ▣ Time-series data, temporal data, sequence data
  - ▣ Structure data, graphs, social networks and multi-linked data
  - ▣ Heterogeneous databases and legacy databases
  - ▣ Spatial, spatiotemporal, multimedia, text and Web data
  - ▣ Software programs, scientific simulations
- New and sophisticated applications

# Knowledge Discovery (KDD) Process

- This is a view from typical database systems and data warehousing communities
- Data mining plays an essential role in the knowledge discovery process



# Example: A Web Mining Framework

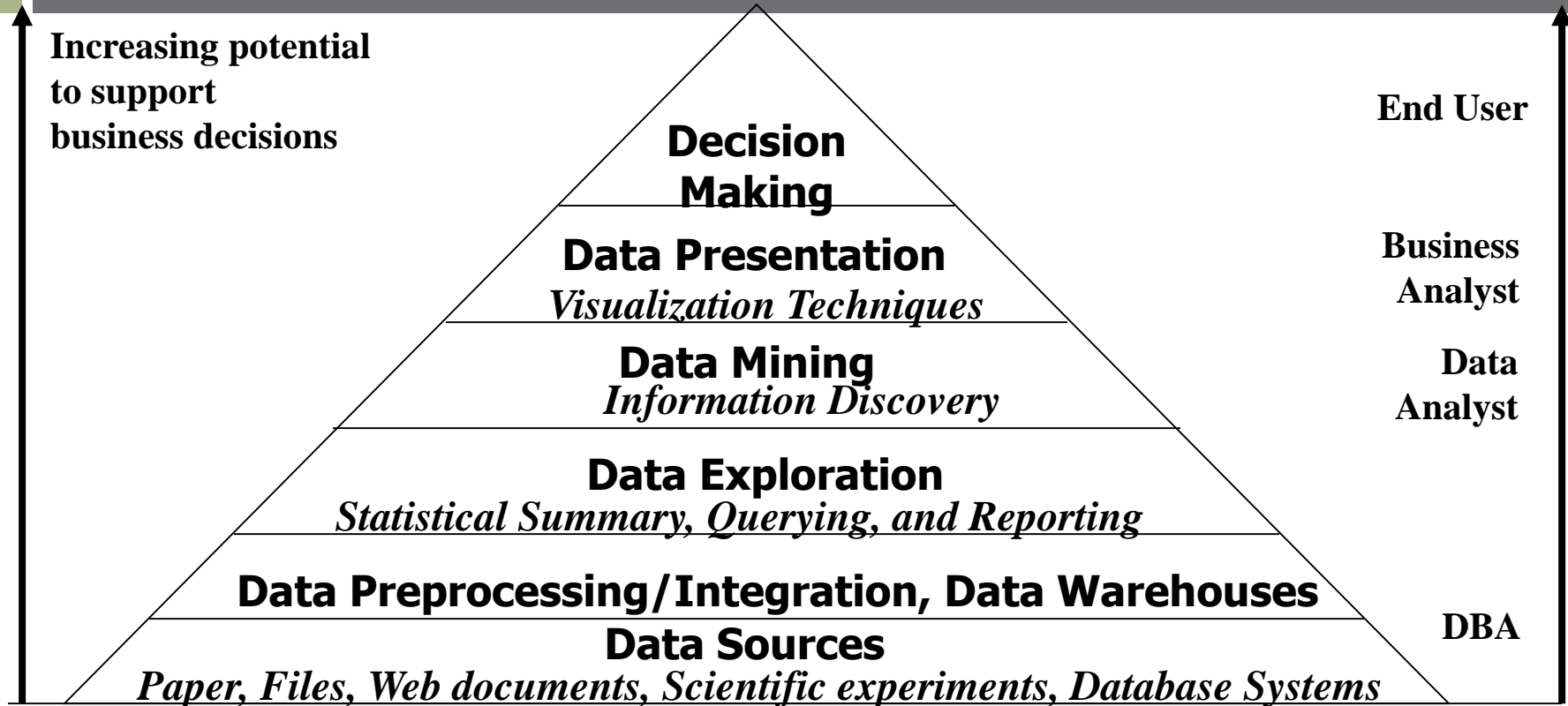
64

- Web mining usually involves
  - Data cleaning
  - Data integration from multiple sources
  - Warehousing the data
  - Data cube construction
  - Data selection for data mining
  - Data mining
  - Presentation of the mining results
  - Patterns and knowledge to be used or stored in a knowledge-base

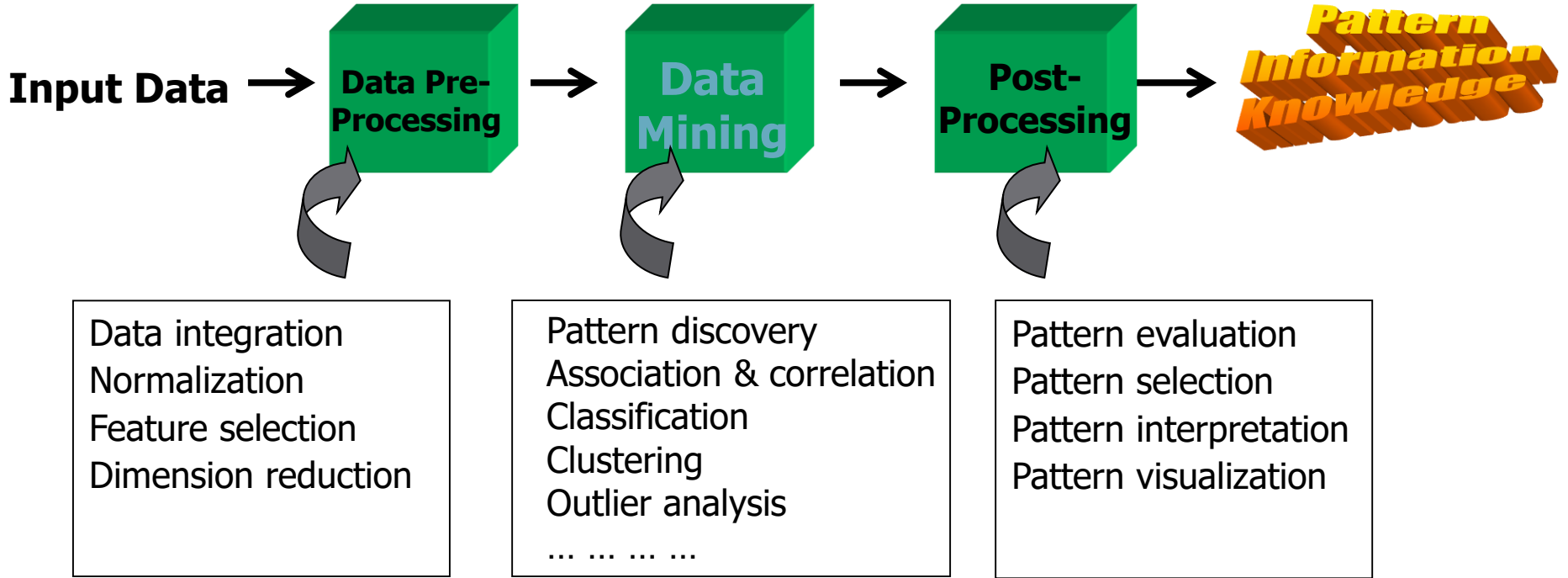


# Data Mining in Business Intelligence

65



# KDD Process: A Typical View from ML and Statistics



□ This is a view from typical machine learning and statistics communities

# Multi-Dimensional View of Data Mining (1)

67

- **Data to be mined**
  - ▣ Database data (extended-relational, object-oriented, heterogeneous, legacy), data warehouse, transactional data, stream, spatiotemporal, time-series, sequence, text and web, multi-media, graphs & social and information networks
- **Knowledge to be mined (or: Data mining functions)**
  - ▣ Characterization, discrimination, association, classification, clustering, trend/deviation, outlier analysis, etc.
  - ▣ Descriptive vs. predictive data mining
  - ▣ Multiple/integrated functions and mining at multiple levels

# Multi-Dimensional View of Data Mining (2)

68

## □ **Techniques utilized**

- ▣ Data-intensive, data warehouse (OLAP), machine learning, statistics, pattern recognition, visualization, high-performance, etc.

## □ **Applications adapted**

- ▣ Retail, telecommunication, banking, fraud analysis, bio-data mining, stock market analysis, text mining, Web mining, etc.

# Data Mining: On What Kinds of Data?

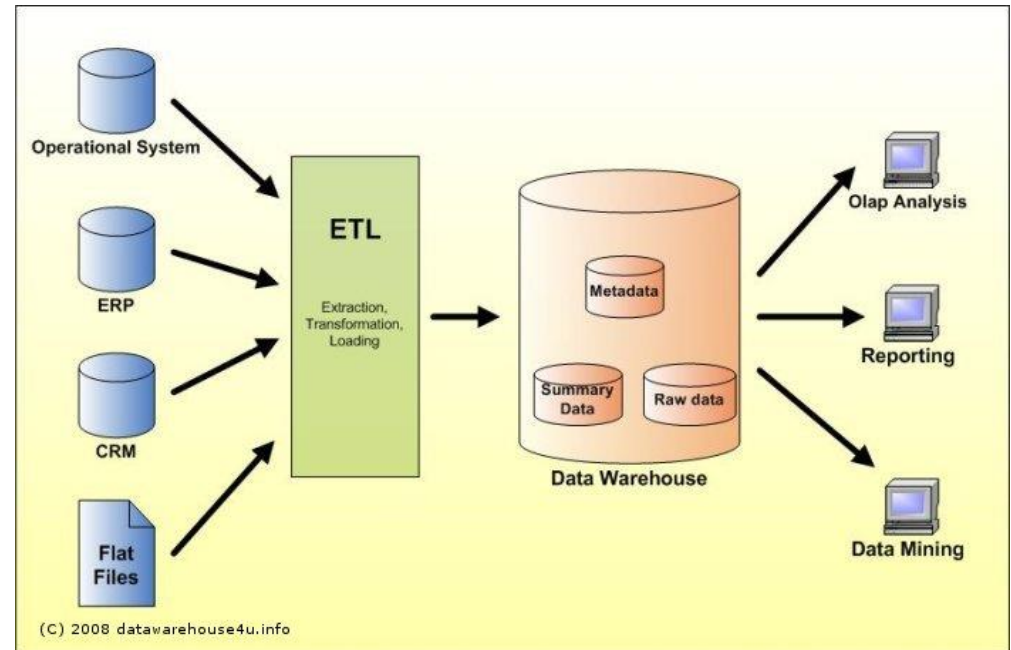
69

- Database-oriented data sets and applications
  - Relational database, data warehouse, transactional database
- Advanced data sets and advanced applications
  - Data streams and sensor data
  - Time-series data, temporal data, sequence data (incl. bio-sequences)
  - Structure data, graphs, social networks and multi-linked data
  - Object-relational databases
  - Heterogeneous databases and legacy databases
  - Spatial data and spatiotemporal data
  - Multimedia database
  - Text databases
  - The World-Wide Web

# Data Mining Function: (1) Generalization

70

- Information integration and data warehouse construction
  - ▣ Data cleaning, transformation, integration, and multidimensional data model





## Data Mining Function: (2) Association and Correlation Analysis

72

- Frequent patterns (or frequent itemsets)
  - What items are frequently purchased together at your grocery store?
- Association, correlation vs. causality
  - A typical association rule
    - Diaper  $\rightarrow$  Beer [0.5%, 75%] (support, confidence)
  - Are strongly associated items also strongly correlated?
- How to mine such patterns and rules efficiently in large datasets?
- How to use such patterns for classification, clustering, and other applications?



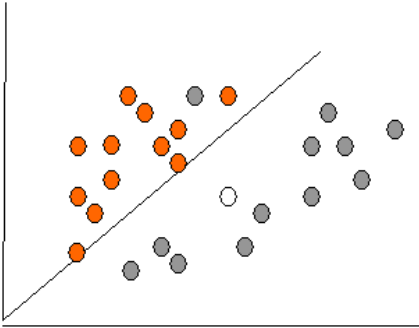
# Data Mining Function: (3) Classification

73

- Classification and label prediction
  - Construct models (functions) based on some training examples
  - Describe and distinguish classes or concepts for future prediction
    - E.g., classify countries based on (climate), or classify cars based on (gas mileage)
  - Predict some unknown class labels
- Typical methods
  - Decision trees, naïve Bayesian classification, support vector machines, neural networks, rule-based classification, pattern-based classification, logistic regression, ...
- Typical applications:
  - Credit card fraud detection, direct marketing, classifying stars, diseases, web-pages, ...

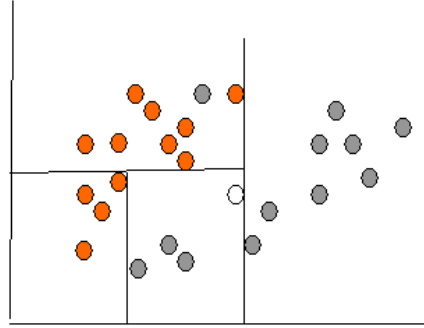
# Classification: Examples

74



## b) Linear Regression

- $w_0 + w_1 x + w_2 y \geq 0$
- Regression computes  $w_i$  from data to minimize squared error to 'fit' the data
- Not flexible enough



## b) Decision tree

- if  $(X > 5)$  then grey
- else if  $(Y > 3)$  then orange
- else if  $(X > 2)$  then grey
- else orange



## b) Neural network

- Can select more complex regions
- Can be more accurate
- Also can overfit the data – find patterns in random noise

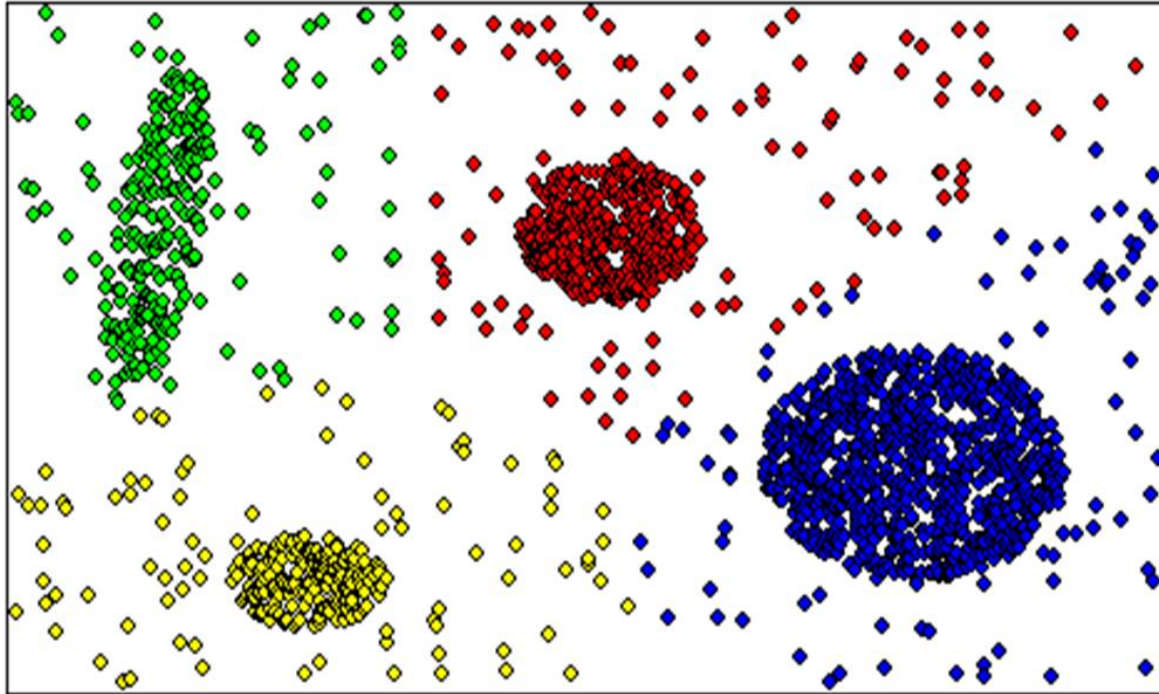
# Data Mining Function: (4) Cluster Analysis

75

- Unsupervised learning (i.e., Class label is unknown)
- Group data to form new categories (i.e., clusters), e.g., cluster houses to find distribution patterns
- Principle: Maximizing intra-class similarity & minimizing interclass similarity
- Many methods and applications

# Cluster Analysis: Example

76

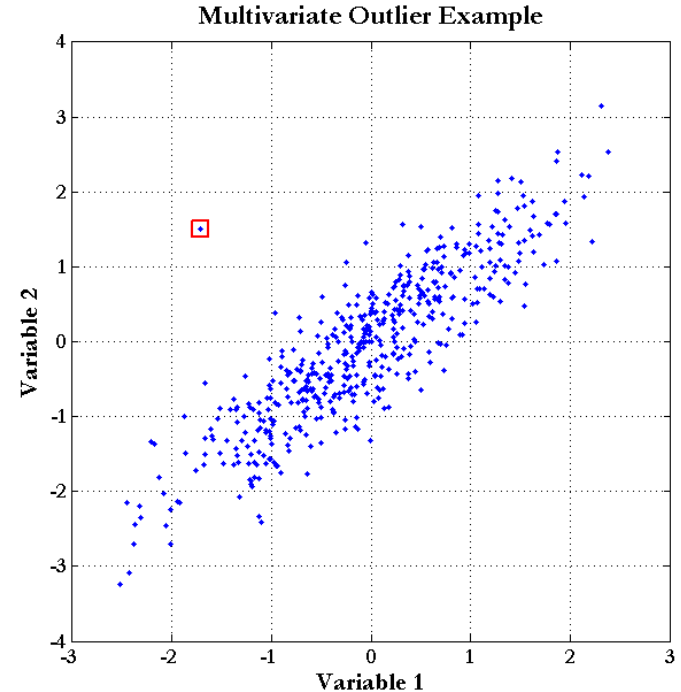


# Data Mining Function: (5) Outlier Analysis

77

## □ Outlier analysis

- Outlier: A data object that does not comply with the general behavior of the data
- Noise or exception? — One person's garbage could be another person's treasure
- Methods: by-product of clustering or regression analysis, ...
- Useful in fraud detection, rare events analysis



## Time and Ordering: Sequential Pattern, Trend and Evolution Analysis

78

- Sequence, trend and evolution analysis
  - Trend, time-series, and deviation analysis: e.g., regression and value prediction
  - Sequential pattern mining
    - e.g., first buy digital camera, then buy large SD memory cards
  - Periodicity analysis
  - Motifs and biological sequence analysis
    - Approximate and consecutive motifs
  - Similarity-based analysis
- Mining data streams
  - Ordered, time-varying, potentially infinite, data streams

# Structure and Network Analysis

79

- Graph mining
  - ▣ Finding frequent subgraphs (e.g., chemical compounds), trees (XML), substructures (web fragments)
- Information network analysis
  - ▣ Social networks: actors (objects, nodes) and relationships (edges)
    - e.g., author networks in CS, terrorist networks
  - ▣ Multiple heterogeneous networks
    - A person could be multiple information networks: friends, family, classmates, ...
  - ▣ Links carry a lot of semantic information: Link mining
- Web mining
  - ▣ Web is a big information network: from PageRank to Google
  - ▣ Analysis of Web information networks
    - Web community discovery, opinion mining, usage mining, ...

# Applications of Data Mining

- Web page analysis: from web page classification, clustering to PageRank & HITS algorithms
- Collaborative analysis & recommender systems
- Basket data analysis to targeted marketing
- Biological and medical data analysis: classification, cluster analysis (microarray data analysis), biological sequence analysis, biological network analysis
- Data mining and software engineering (e.g., IEEE Computer, Aug. 2009 issue)
- From major dedicated data mining systems/tools (e.g., SAS, MS SQL-Server Analysis Manager, Oracle Data Mining Tools) to invisible data mining



# Summary

- Data mining: Discovering interesting patterns and knowledge from massive amount of data
- A natural evolution of database technology, in great demand, with wide applications
- A KDD process includes data cleaning, data integration, data selection, transformation, data mining, pattern evaluation, and knowledge presentation
- Mining can be performed in a variety of data
- Data mining functionalities: characterization, discrimination, association, classification, clustering, outlier and trend analysis, etc.