

# Data Mining: Concepts and Techniques

---

(3<sup>rd</sup> ed.)

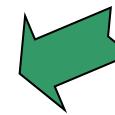
## — Chapter 3 —

Jiawei Han, Micheline Kamber, and Jian Pei  
University of Illinois at Urbana-Champaign &  
Simon Fraser University

©2011 Han, Kamber & Pei. All rights reserved.

# Chapter 3: Data Preprocessing

---

- n Data Preprocessing: An Overview 
- n Data Quality
- n Major Tasks in Data Preprocessing
- n Data Cleaning
- n Data Integration
- n Data Reduction
- n Data Transformation and Data Discretization
- n Summary

# Data Quality: Why Preprocess the Data ?

---

- n Measures for data quality: A multidimensional view
  - n Accuracy: correct or wrong, accurate or not
  - n Completeness: not recorded, unavailable, ...
  - n Consistency: some modified but some not, dangling, ...
  - n Timeliness: timely update?
  - n Believability: how trustable the data are correct?
  - n Interpretability: how easily the data can be understood?

# Major Tasks in Data Preprocessing

---

- n **Data cleaning**

- n Fill in missing values, smooth noisy data, identify or remove outliers, and resolve inconsistencies

- n **Data integration**

- n Integration of multiple databases, data cubes, or files

- n **Data reduction**

- n Dimensionality reduction
  - n Numerosity reduction
  - n Data compression

- n **Data transformation and data discretization**

- n Normalization
  - n Concept hierarchy generation

# Chapter 3: Data Preprocessing

---

- n Data Preprocessing: An Overview
  - n Data Quality
  - n Major Tasks in Data Preprocessing
- n Data Cleaning
- n Data Integration
- n Data Reduction
- n Data Transformation and Data Discretization
- n Summary



# Data Cleaning

---

- n Data in the Real World Is Dirty: Lots of potentially incorrect data, e.g., instrument faulty, human or computer error, transmission error
  - n incomplete: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
    - n e.g., *Occupation*=“ ” (missing data)
  - n noisy: containing noise, errors, or outliers
    - n e.g., *Salary*=“−10” (an error)
  - n inconsistent: containing discrepancies in codes or names, e.g.,
    - n *Age*=“42”, *Birthday*=“03/07/2010”
    - n Was rating “1, 2, 3”, now rating “A, B, C”
    - n discrepancy between duplicate records
  - n Intentional (e.g., *disguised missing* data)
    - n Jan. 1 as everyone’s birthday?

# Incomplete (Missing) Data

---

- n Data is not always available
  - n E.g., many tuples have no recorded value for several attributes, such as customer income in sales data
- n Missing data may be due to
  - n equipment malfunction
  - n inconsistent with other recorded data and thus deleted
  - n data not entered due to misunderstanding
  - n certain data may not be considered important at the time of entry
  - n not register history or changes of the data
- n Missing data may need to be inferred

# How to Handle Missing Data ?

---

- n Ignore the tuple: usually done when class label is missing (when doing classification)—not effective when the % of missing values per attribute varies considerably
- n Fill in the missing value manually: tedious + infeasible?
- n Fill in it automatically with
  - n a global constant : e.g., “unknown”, a new class?!
  - n the attribute mean
  - n the attribute mean for all samples belonging to the same class: smarter
  - n the most probable value: inference-based such as Bayesian formula or decision tree

# Noisy Data

---

- n **Noise**: random error or variance in a measured variable
- n **Incorrect attribute values** may be due to
  - n faulty data collection instruments
  - n data entry problems
  - n data transmission problems
  - n technology limitation
  - n inconsistency in naming convention
- n **Other data problems** which require data cleaning
  - n duplicate records
  - n incomplete data
  - n inconsistent data

# How to Handle Noisy Data ?

---

- n Binning

- n first sort data and partition into (equal-frequency) bins
  - n then one can smooth by bin means, smooth by bin median, smooth by bin boundaries, etc.

- n Regression

- n smooth by fitting the data into regression functions

- n Clustering

- n detect and remove outliers

- n Combined computer and human inspection

- n detect suspicious values and check by human (e.g., deal with possible outliers)

# Data Cleaning as a Process

---

- n Data discrepancy detection
  - n Use metadata (e.g., domain, range, dependency, distribution)
  - n Check field overloading
  - n Check uniqueness rule, consecutive rule and null rule
  - n Use commercial tools
    - n Data scrubbing: use simple domain knowledge (e.g., postal code, spell-check) to detect errors and make corrections
    - n Data auditing: by analyzing data to discover rules and relationship to detect violators (e.g., correlation and clustering to find outliers)
- n Data migration and integration
  - n Data migration tools: allow transformations to be specified
  - n ETL (Extraction/Transformation>Loading) tools: allow users to specify transformations through a graphical user interface
- n Integration of the two processes
  - n Iterative and interactive (e.g., Potter's Wheels)

# Chapter 3: Data Preprocessing

---

- n Data Preprocessing: An Overview
  - n Data Quality
  - n Major Tasks in Data Preprocessing
- n Data Cleaning
- n Data Integration
- n Data Reduction
- n Data Transformation and Data Discretization
- n Summary



# Data Integration

---

- n **Data integration:**

- n Combines data from multiple sources into a coherent store

- n Schema integration: e.g., A.cust-id     B.cust-#

- n Integrate metadata from different sources

- n Entity identification problem:

- n Identify real world entities from multiple data sources, e.g., Bill Clinton = William Clinton

- n Detecting and resolving data value conflicts

- n For the same real world entity, attribute values from different sources are different
  - n Possible reasons: different representations, different scales, e.g., metric vs. British units

# Handling Redundancy in Data Integration

---

- Redundant data occur often when integration of multiple databases
  - *Object identification:* The same attribute or object may have different names in different databases
  - *Derivable data:* One attribute may be a “derived” attribute in another table, e.g., annual revenue
- Redundant attributes may be able to be detected by *correlation analysis* and *covariance analysis*
- Careful integration of the data from multiple sources may help reduce/avoid redundancies and inconsistencies and improve mining speed and quality

# Correlation Analysis (Nominal Data)

## n **X<sup>2</sup> (chi-square) test**

$$\chi^2 = \sum \frac{(Observed - Expected)^2}{Expected}$$

- n The larger the X<sup>2</sup> value, the more likely the variables are related
- n The cells that contribute the most to the X<sup>2</sup> value are those whose actual count is very different from the expected count
- n Correlation does not imply causality
  - n # of hospitals and # of car-theft in a city are correlated
  - n Both are causally linked to the third variable: population

# Chi-Square Calculation: An Example

	Play chess	Not play chess	Sum (row)
Like science fiction	250(90)	200(360)	450
Not like science fiction	50(210)	1000(840)	1050
Sum(col.)	300	1200	1500

- n  $\chi^2$  (chi-square) calculation (numbers in parenthesis are expected counts calculated based on the data distribution in the two categories)

$$\chi^2 = \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(200-360)^2}{360} + \frac{(1000-840)^2}{840} = 507.93$$

- n It shows that like\_science\_fiction and play\_chess are correlated in the group

# Correlation Analysis (Numeric Data)

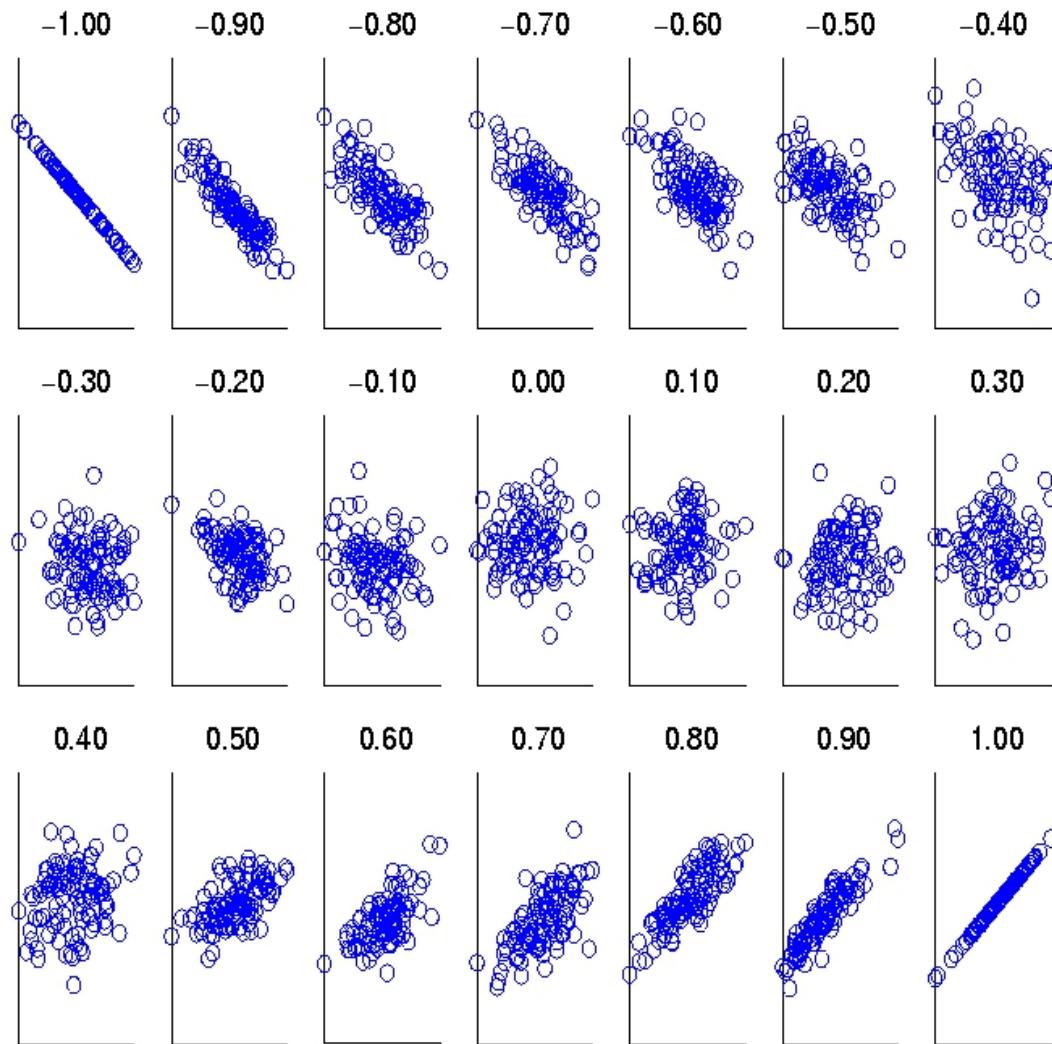
- n Correlation coefficient (also called Pearson's product moment coefficient)

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{(n-1)\sigma_A \sigma_B} = \frac{\sum_{i=1}^n (a_i b_i) - n \bar{A} \bar{B}}{(n-1)\sigma_A \sigma_B}$$

where  $n$  is the number of tuples,  $\bar{A}$  and  $\bar{B}$  are the respective means of A and B,  $\sigma_A$  and  $\sigma_B$  are the respective standard deviation of A and B, and  $\Sigma(a_i b_i)$  is the sum of the AB cross-product.

- n If  $r_{A,B} > 0$ , A and B are positively correlated (A's values increase as B's). The higher, the stronger correlation.
- n  $r_{A,B} = 0$ : independent;  $r_{AB} < 0$ : negatively correlated

# Visually Evaluating Correlation



**Scatter plots  
showing the  
similarity from  
-1 to 1.**

# Correlation (viewed as linear relationship)

---

- n Correlation measures the linear relationship between objects
- n To compute correlation, we standardize data objects, A and B, and then take their dot product

$$a'_k = (a_k - \text{mean}(A)) / \text{std}(A)$$

$$b'_k = (b_k - \text{mean}(B)) / \text{std}(B)$$

$$\text{correlation}(A, B) = A' \bullet B'$$

# Covariance (Numeric Data)

- Covariance is similar to correlation

$$Cov(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

Correlation coefficient:

$$r_{A,B} = \frac{Cov(A, B)}{\sigma_A \sigma_B}$$

where n is the number of tuples,  $\bar{A}$  and  $\bar{B}$  are the respective mean or **expected values** of A and B,  $\sigma_A$  and  $\sigma_B$  are the respective standard deviation of A and B.

- **Positive covariance:** If  $Cov_{A,B} > 0$ , then A and B both tend to be larger than their expected values.
- **Negative covariance:** If  $Cov_{A,B} < 0$  then if A is larger than its expected value, B is likely to be smaller than its expected value.
- **Independence:**  $Cov_{A,B} = 0$  but the converse is not true:
  - Some pairs of random variables may have a covariance of 0 but are not independent. Only under some additional assumptions (e.g., the data follow multivariate normal distributions) does a covariance of 0 imply independence.

# Co-Variance: An Example

$$Cov(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n}$$

- n It can be simplified in computation as

$$Cov(A, B) = E(A \cdot B) - \bar{A}\bar{B}$$

- n Suppose two stocks A and B have the following values in one week:  
(2, 5), (3, 8), (5, 10), (4, 11), (6, 14).
- n Question: If the stocks are affected by the same industry trends, will their prices rise or fall together?

n  $E(A) = (2 + 3 + 5 + 4 + 6)/ 5 = 20/5 = 4$

n  $E(B) = (5 + 8 + 10 + 11 + 14) /5 = 48/5 = 9.6$

n  $Cov(A,B) = (2 \times 5 + 3 \times 8 + 5 \times 10 + 4 \times 11 + 6 \times 14) / 5 - 4 \times 9.6 = 4$

- n Thus, A and B rise together since  $Cov(A, B) > 0$ .

# Chapter 3: Data Preprocessing

---

- n Data Preprocessing: An Overview
  - n Data Quality
  - n Major Tasks in Data Preprocessing
- n Data Cleaning
- n Data Integration
- n Data Reduction 
- n Data Transformation and Data Discretization
- n Summary

# Data Reduction Strategies

- n **Data reduction:** Obtain a reduced representation of the data set that is much smaller in volume but yet produces the same (or almost the same) analytical results
- n Why data reduction? — A database/data warehouse may store terabytes of data. Complex data analysis may take a very long time to run on the complete data set.
- n Data reduction strategies
  - n Dimensionality reduction, e.g., remove unimportant attributes
    - n Wavelet transforms
    - n Principal Components Analysis (PCA)
    - n Feature subset selection, feature creation
  - n Numerosity reduction (some simply call it: Data Reduction)
    - n Regression and Log-Linear Models
    - n Histograms, clustering, sampling
    - n Data cube aggregation
  - n Data compression

# Data Reduction 1: Dimensionality Reduction

---

## n **Curse of dimensionality**

- n When dimensionality increases, data becomes increasingly sparse
- n Density and distance between points, which is critical to clustering, outlier analysis, becomes less meaningful
- n The possible combinations of subspaces will grow exponentially

## n **Dimensionality reduction**

- n Avoid the curse of dimensionality
- n Help eliminate irrelevant features and reduce noise
- n Reduce time and space required in data mining
- n Allow easier visualization

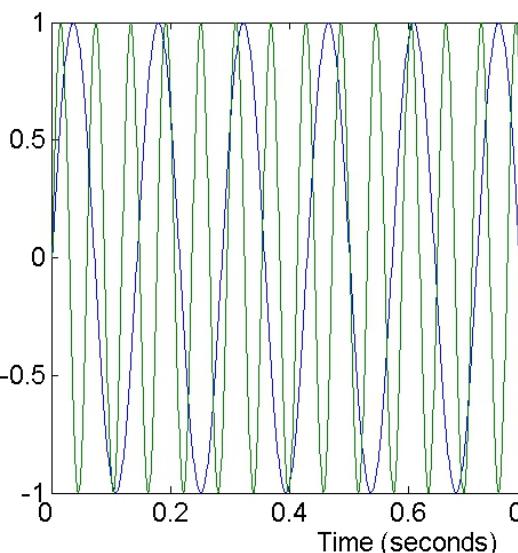
## n **Dimensionality reduction techniques**

- n Wavelet transforms
- n Principal Component Analysis
- n Supervised and nonlinear techniques (e.g., feature selection)

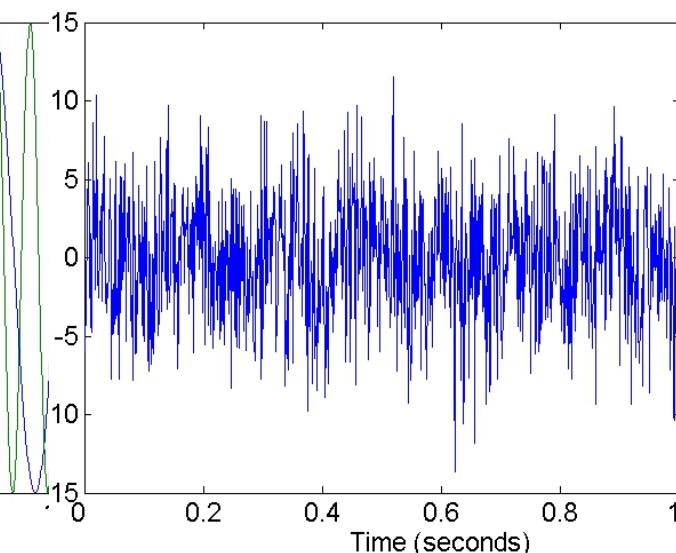
# Mapping Data to a New Space

n Fourier transform

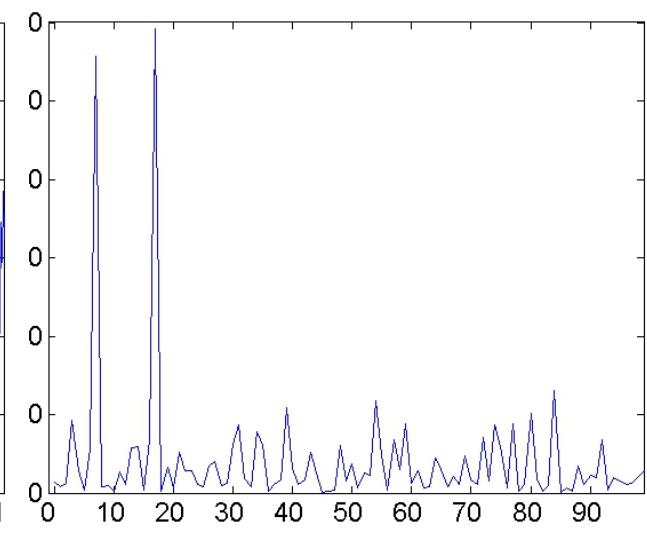
n Wavelet transform



Two Sine Waves



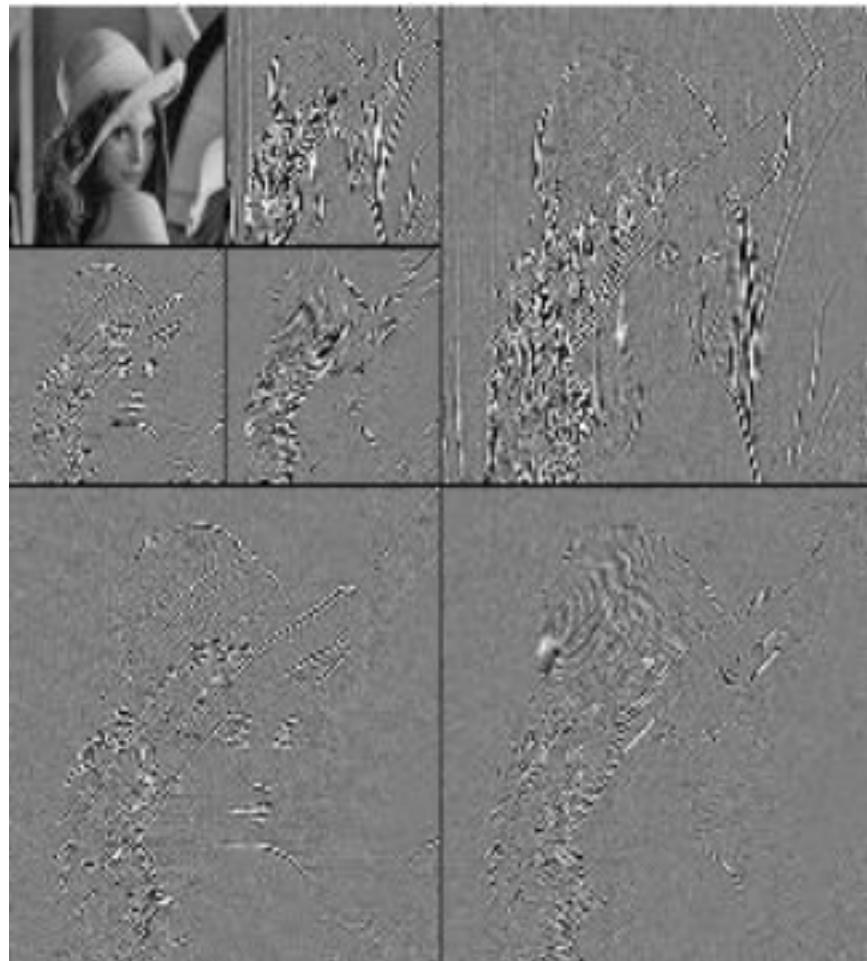
Two Sine Waves + Noise



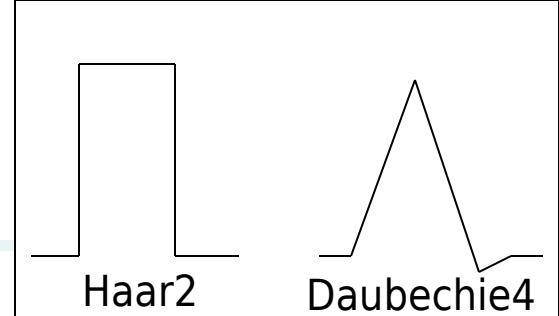
Frequency

# What Is Wavelet Transform ?

- n Decomposes a signal into different frequency subbands
  - n Applicable to n-dimensional signals
- n Data are transformed to preserve relative distance between objects at different levels of resolution
- n Allow natural clusters to become more distinguishable
- n Used for image compression



# Wavelet Transformation



- n Discrete wavelet transform (DWT) for linear signal processing, multi-resolution analysis
- n Compressed approximation: store only a small fraction of the strongest of the wavelet coefficients
- n Similar to discrete Fourier transform (DFT), but better lossy compression, localized in space
- n Method:
  - n Length,  $L$ , must be an integer power of 2 (padding with 0's, when necessary)
  - n Each transform has 2 functions: smoothing, difference
  - n Applies to pairs of data, resulting in two set of data of length  $L/2$
  - n Applies two functions recursively, until reaches the desired length

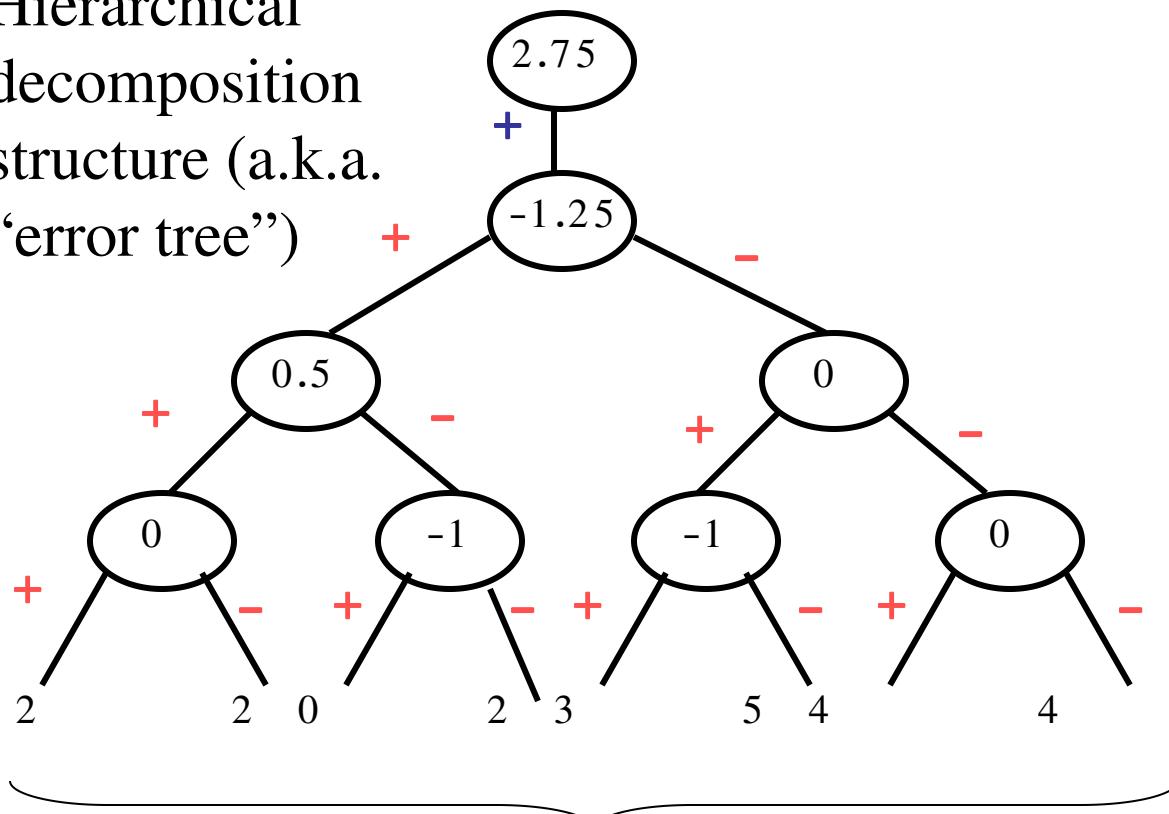
# Wavelet Decomposition

- n Wavelets: A math tool for space-efficient hierarchical decomposition of functions
- n  $S = [2, 2, 0, 2, 3, 5, 4, 4]$  can be transformed to  $\hat{S} = [2^3/4, -1^1/4, 1/2, 0, 0, -1, -1, 0]$
- n Compression: many small detail coefficients can be replaced by 0's, and only the significant coefficients are retained

Resolution	Averages	Detail Coefficients
8	$[2, 2, 0, 2, 3, 5, 4, 4]$	
4	$[2, 1, 4, 4]$	$[0, -1, -1, 0]$
2	$[1\frac{1}{2}, 4]$	$[\frac{1}{2}, 0]$
1	$[2\frac{3}{4}]$	$[-1\frac{1}{4}]$

# Haar Wavelet Coefficients

Hierarchical decomposition structure (a.k.a.  
“error tree”)



Original frequency distribution

Coefficient “Supports”

2.75



-1.25



0.5



0



0



-1



-1



0



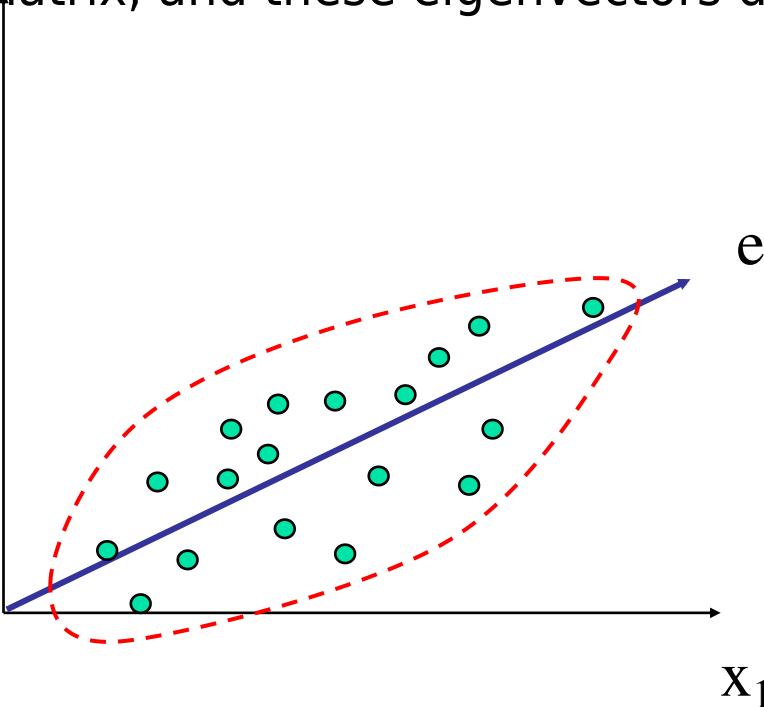
# Why Wavelet Transform ?

---

- n Use hat-shape filters
  - n Emphasize region where points cluster
  - n Suppress weaker information in their boundaries
- n Effective removal of outliers
  - n Insensitive to noise, insensitive to input order
- n Multi-resolution
  - n Detect arbitrary shaped clusters at different scales
- n Efficient
  - n Complexity  $O(N)$
- n Only applicable to low dimensional data

# Principal Component Analysis (PCA)

- Find a projection that captures the largest amount of variation in data
- The original data are projected onto a much smaller space, resulting in dimensionality reduction. We find the eigenvectors of the covariance matrix, and these eigenvectors define the new space



# Principal Component Analysis (Steps)

---

- Given  $N$  data vectors from  $n$ -dimensions, find  $k \leq n$  orthogonal vectors (*principal components*) that can be best used to represent data
  - Normalize input data: Each attribute falls within the same range
  - Compute  $k$  orthonormal (unit) vectors, i.e., *principal components*
  - Each input data (vector) is a linear combination of the  $k$  principal component vectors
  - The principal components are sorted in order of decreasing “significance” or strength
  - Since the components are sorted, the size of the data can be reduced by eliminating the *weak components*, i.e., those with low variance (i.e., using the strongest principal components, it is possible to reconstruct a good approximation of the original data)

# Attribute Subset Selection

---

- ▀ Another way to reduce dimensionality of data
- ▀ Redundant attributes
  - ▀ Duplicate much or all of the information contained in one or more other attributes
    - ▀ E.g., purchase price of a product and the amount of sales tax paid
- ▀ Irrelevant attributes
  - ▀ Contain no information that is useful for the data mining task at hand
    - ▀ E.g., students' ID is often irrelevant to the task of predicting students' GPA

# Heuristic Search in Attribute Selection

---

- n There are  $2^d$  possible attribute combinations of  $d$  attributes
- n Typical heuristic attribute selection methods:
  - n Best single attribute under the attribute independence assumption: choose by significance tests
  - n Best step-wise feature selection:
    - n The best single-attribute is picked first
    - n Then next best attribute condition to the first, ...
  - n Step-wise attribute elimination:
    - n Repeatedly eliminate the worst attribute
  - n Best combined attribute selection and elimination
  - n Optimal branch and bound:
    - n Use attribute elimination and backtracking

# Attribute Creation (Feature Generation)

---

- n Create new attributes (features) that can capture the important information in a data set more effectively than the original ones
- n Three general methodologies
  - n Attribute extraction
    - n Domain-specific
    - n Mapping data to new space (see: data reduction)
      - n E.g., Fourier transformation, wavelet transformation, manifold approaches (not covered)
  - n Attribute construction
    - n Combining features (see: discriminative frequent patterns in Chapter 7)
    - n Data discretization

# Data Reduction 2: Numerosity Reduction

---

- n Reduce data volume by choosing alternative, *smaller forms* of data representation
- n **Parametric methods** (e.g., regression)
  - n Assume the data fits some model, estimate model parameters, store only the parameters, and discard the data (except possible outliers)
  - n Ex.: Log-linear models—obtain value at a point in  $m$ -D space as the product on appropriate marginal subspaces
- n **Non-parametric** methods
  - n Do not assume models
  - n Major families: histograms, clustering, sampling, ...

# Parametric Data Reduction: Regression and Log-Linear Models

---

## n **Linear regression**

- n Data modeled to fit a straight line
- n Often uses the least-square method to fit the line

## n **Multiple regression**

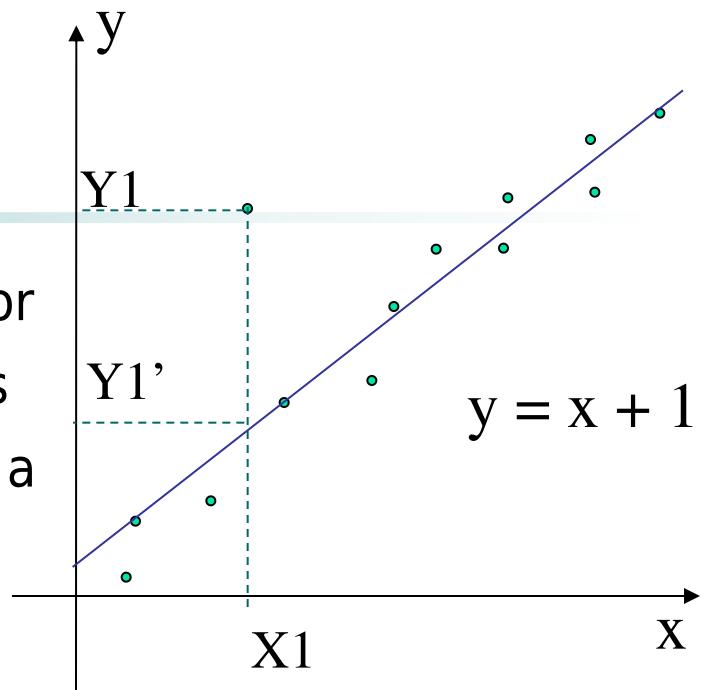
- n Allows a response variable  $Y$  to be modeled as a linear function of multidimensional feature vector

## n **Log-linear model**

- n Approximates discrete multidimensional probability distributions

# Regression Analysis

- n Regression analysis: A collective name for techniques for the modeling and analysis of numerical data consisting of values of a **dependent variable** (also called **response variable** or *measurement*) and of one or more *independent variables* (aka. **explanatory variables** or **predictors**)
- n The parameters are estimated so as to give a "**best fit**" of the data
- n Most commonly the best fit is evaluated by using the **least squares method**, but other criteria have also been used



- n Used for prediction (including forecasting of time-series data), inference, hypothesis testing, and modeling of causal relationships

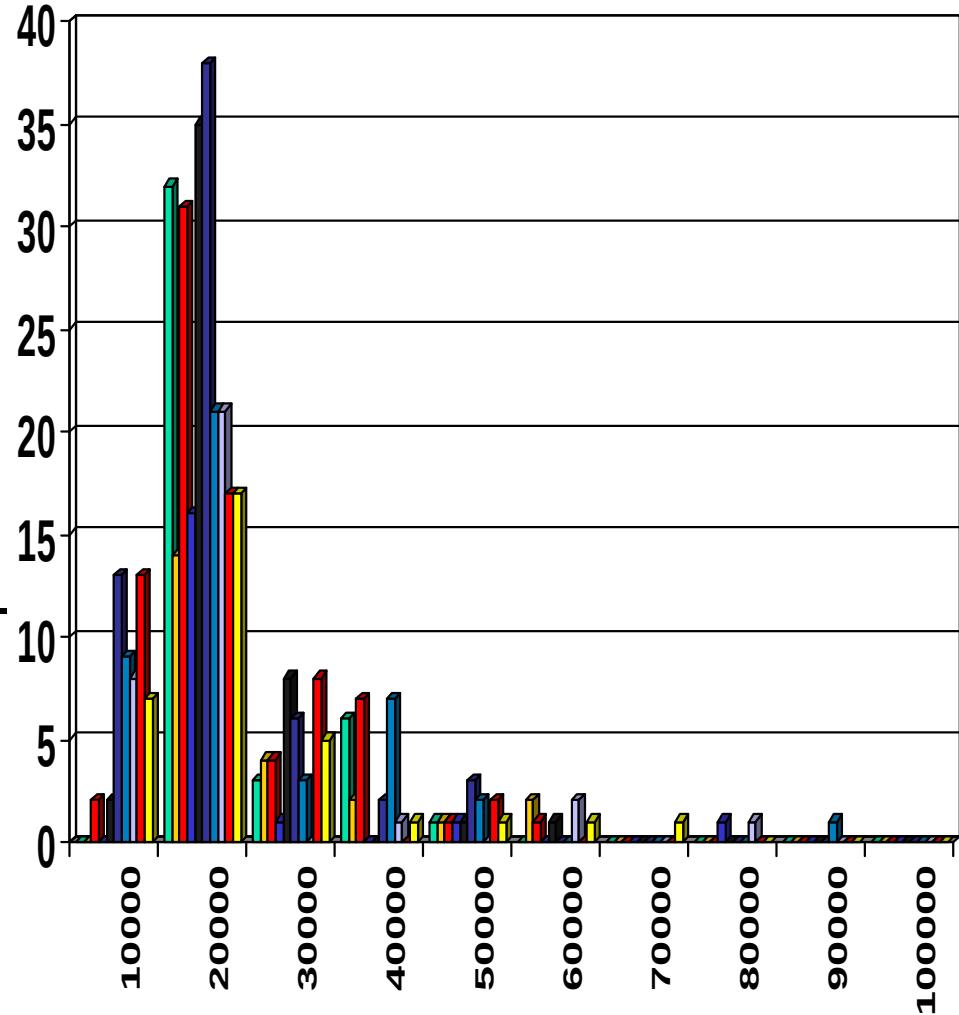
# Regress Analysis and Log-Linear Models

---

- n Linear regression:  $Y = w X + b$ 
  - n Two regression coefficients,  $w$  and  $b$ , specify the line and are to be estimated by using the data at hand
  - n Using the least squares criterion to the known values of  $Y_1, Y_2, \dots, X_1, X_2, \dots$
- n Multiple regression:  $Y = b_0 + b_1 X_1 + b_2 X_2$ 
  - n Many nonlinear functions can be transformed into the above
- n Log-linear models:
  - n Approximate discrete multidimensional probability distributions
  - n Estimate the probability of each point (tuple) in a multi-dimensional space for a set of discretized attributes, based on a smaller subset of dimensional combinations
  - n Useful for dimensionality reduction and data smoothing

# Histogram Analysis

- n Divide data into buckets and store average (sum) for each bucket
- n Partitioning rules:
  - n Equal-width: equal bucket range
  - n Equal-frequency (or equal-depth)



# Clustering

---

- n Partition data set into clusters based on similarity, and store cluster representation (e.g., centroid and diameter) only
- n Can be very effective if data is clustered but not if data is “smeared”
- n Can have hierarchical clustering and be stored in multi-dimensional index tree structures
- n There are many choices of clustering definitions and clustering algorithms
- n Cluster analysis will be studied in depth in Chapter 10

# Sampling

---

- n Sampling: obtaining a small sample  $s$  to represent the whole data set  $N$
- n Allow a mining algorithm to run in complexity that is potentially sub-linear to the size of the data
- n Key principle: Choose a **representative** subset of the data
  - n Simple random sampling may have very poor performance in the presence of skew
  - n Develop adaptive sampling methods, e.g., stratified sampling:
- n Note: Sampling may not reduce database I/Os (page at a time)

# Types of Sampling

---

- n **Simple random sampling**

- n There is an equal probability of selecting any particular item

- n **Sampling without replacement**

- n Once an object is selected, it is removed from the population

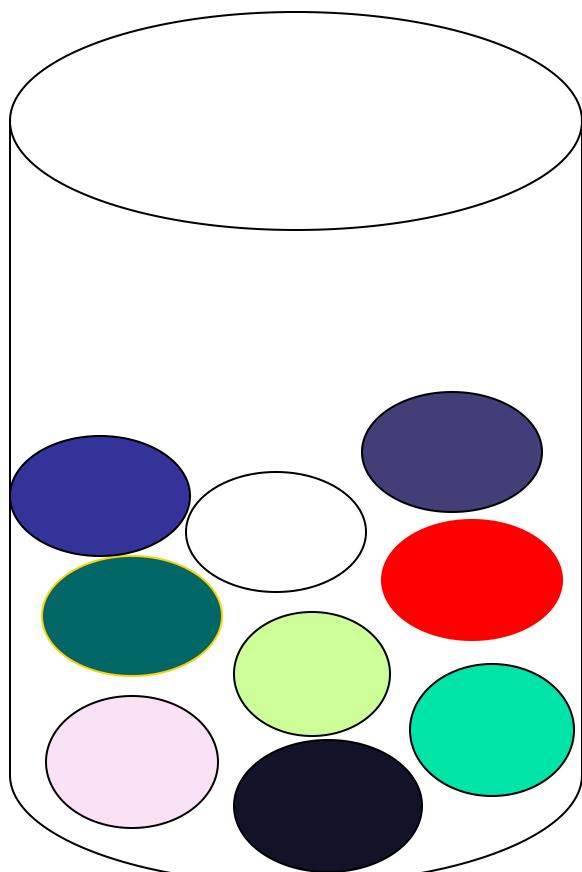
- n **Sampling with replacement**

- n A selected object is not removed from the population

- n **Stratified sampling:**

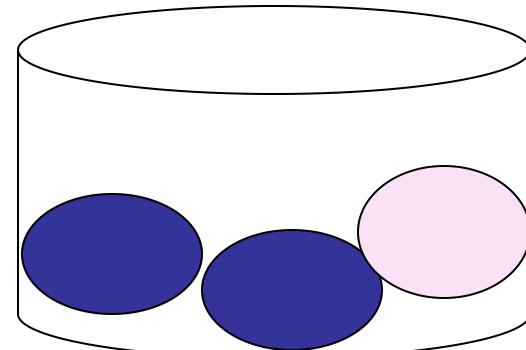
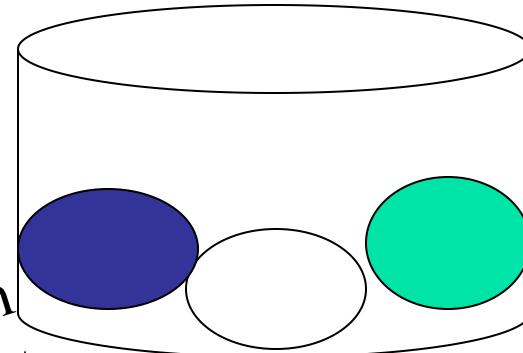
- n Partition the data set, and draw samples from each partition (proportionally, i.e., approximately the same percentage of the data)
  - n Used in conjunction with skewed data

# Sampling: With or without Replacement



*SRSWOR*  
(simple random  
sample without  
replacement)

*SRSWR*

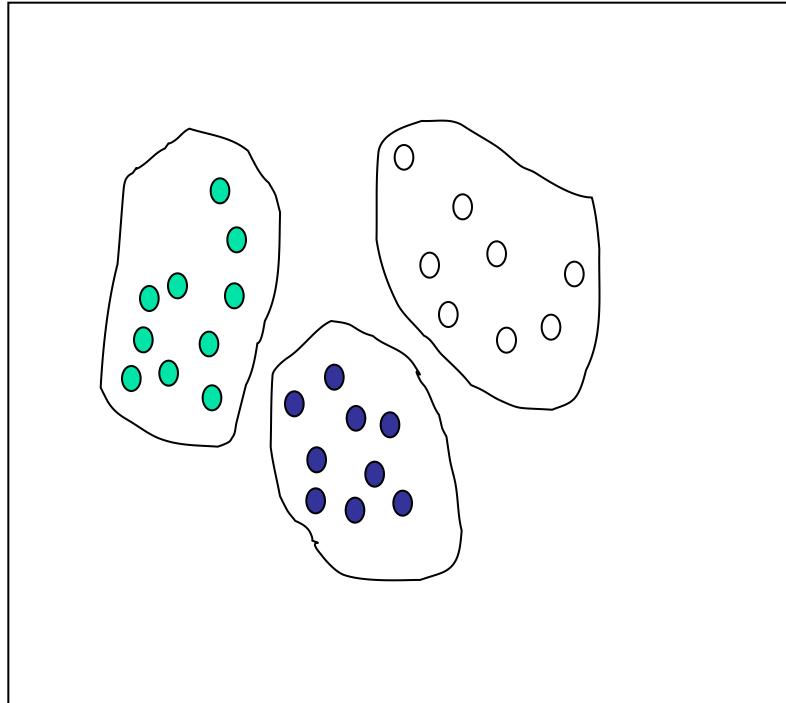


Raw Data

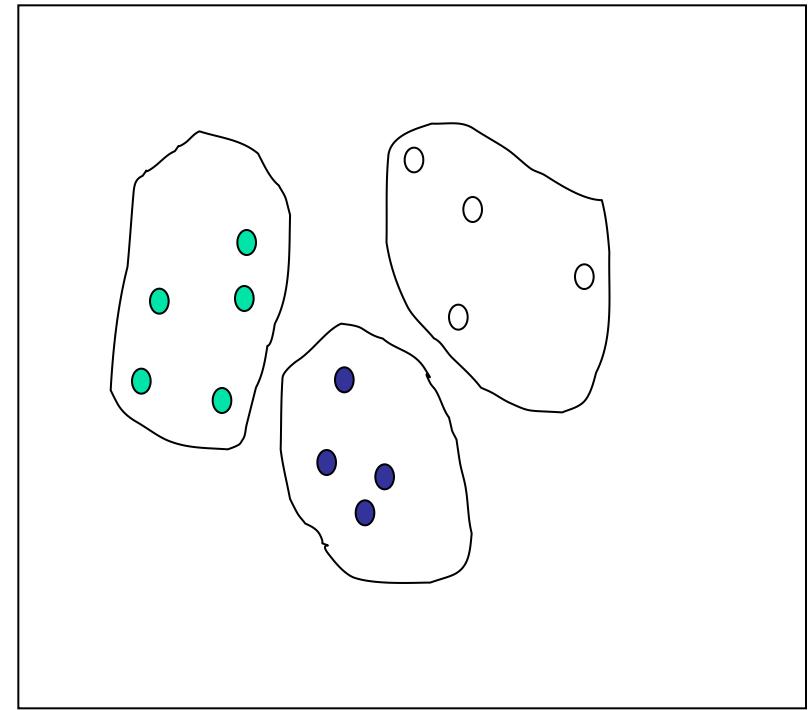
# Sampling: Cluster or Stratified Sampling

---

Raw Data



Cluster/Stratified Sample



# Data Cube Aggregation

---

- n The lowest level of a data cube (base cuboid)
  - n The aggregated data for an **individual entity of interest**
    - n E.g., a customer in a phone calling data warehouse
- n Multiple levels of aggregation in data cubes
  - n Further reduce the size of data to deal with
- n Reference appropriate levels
  - n Use the smallest representation which is enough to solve the task
- n Queries regarding aggregated information should be answered using data cube, when possible

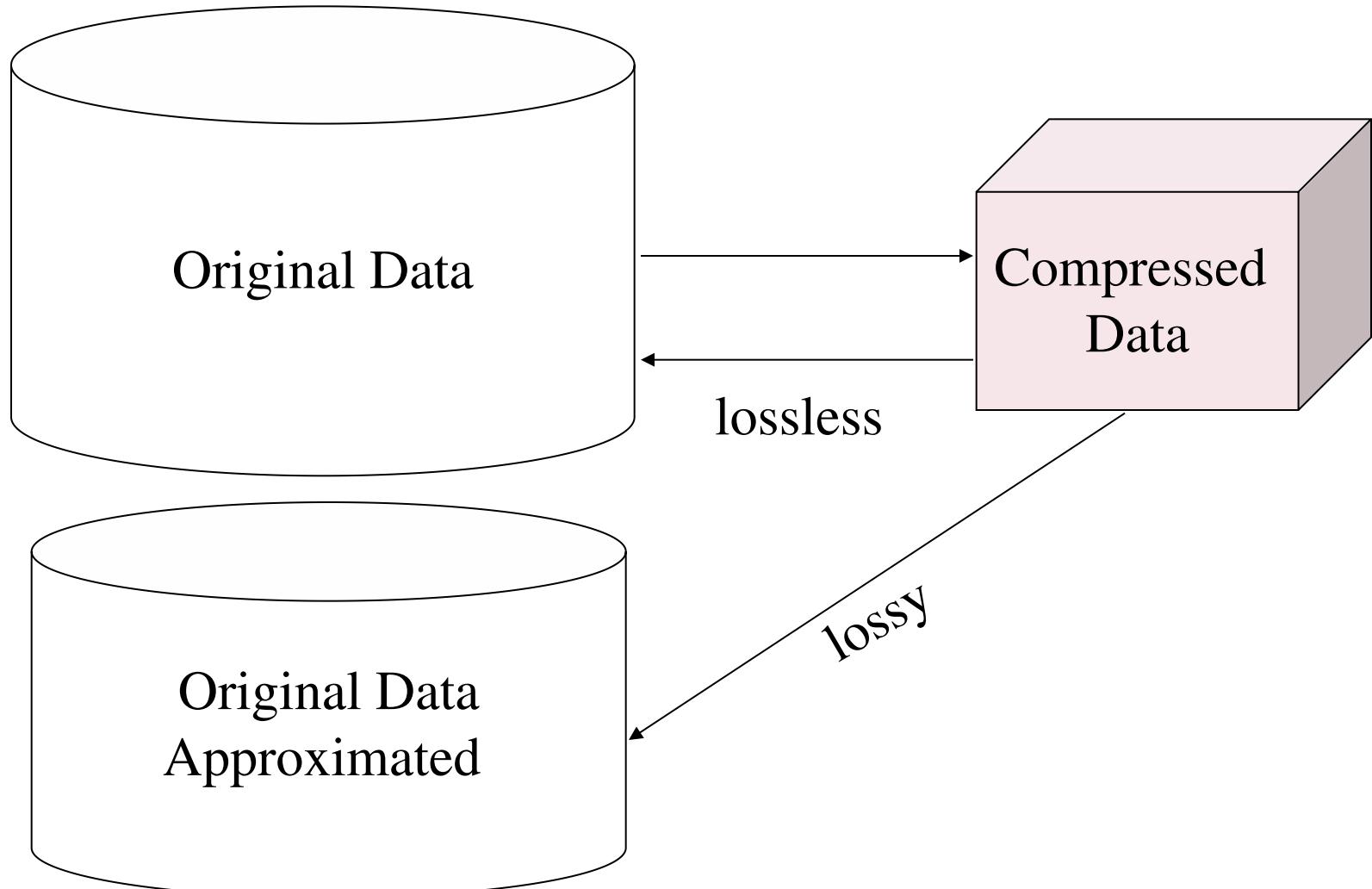
# Data Reduction 3: Data Compression

---

- n String compression
  - n There are extensive theories and well-tuned algorithms
  - n Typically lossless, but only limited manipulation is possible without expansion
- n Audio/video compression
  - n Typically lossy compression, with progressive refinement
  - n Sometimes small fragments of signal can be reconstructed without reconstructing the whole
- n Time sequence is not audio
  - n Typically short and vary slowly with time
- n Dimensionality and numerosity reduction may also be considered as forms of data compression

# Data Compression

---



# Chapter 3: Data Preprocessing

---

- n Data Preprocessing: An Overview
  - n Data Quality
  - n Major Tasks in Data Preprocessing
- n Data Cleaning
- n Data Integration
- n Data Reduction
- n Data Transformation and Data Discretization
- n Summary



# Data Transformation

- n A function that maps the entire set of values of a given attribute to a new set of replacement values s.t. each old value can be identified with one of the new values
- n Methods
  - n Smoothing: Remove noise from data
  - n Attribute/feature construction
    - n New attributes constructed from the given ones
  - n Aggregation: Summarization, data cube construction
  - n Normalization: Scaled to fall within a smaller, specified range
    - n min-max normalization
    - n z-score normalization
    - n normalization by decimal scaling
  - n Discretization: Concept hierarchy climbing

# Normalization

- n **Min-max normalization:** to  $[new\_min_A, new\_max_A]$

$$v' = \frac{v - min_A}{max_A - min_A} (new\_max_A - new\_min_A) + new\_min_A$$

- n Ex. Let income range \$12,000 to \$98,000 normalized to [0.0, 1.0]. Then \$73,000 is mapped to  $\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$

- n **Z-score normalization** ( $\mu$ : mean,  $\sigma$ : standard deviation):

$$v' = \frac{v - \mu_A}{\sigma_A}$$

- n Ex. Let  $\mu = 54,000$ ,  $\sigma = 16,000$ . Then  $\frac{73,600 - 54,000}{16,000} = 1.225$

- n **Normalization by decimal scaling**

$$v' = \frac{v}{10^j} \quad \text{Where } j \text{ is the smallest integer such that } \text{Max}(|v'|) < 1$$

# Discretization

- n Three types of attributes
  - n Nominal—values from an unordered set, e.g., color, profession
  - n Ordinal—values from an ordered set, e.g., military or academic rank
  - n Numeric—real numbers, e.g., integer or real numbers
- n Discretization: Divide the range of a continuous attribute into intervals
  - n Interval labels can then be used to replace actual data values
  - n Reduce data size by discretization
  - n Supervised vs. unsupervised
  - n Split (top-down) vs. merge (bottom-up)
  - n Discretization can be performed recursively on an attribute
  - n Prepare for further analysis, e.g., classification

# Data Discretization Methods

---

- „ Typical methods: All the methods can be applied recursively
  - „ Binning
    - „ Top-down split, unsupervised
  - „ Histogram analysis
    - „ Top-down split, unsupervised
  - „ Clustering analysis (unsupervised, top-down split or bottom-up merge)
  - „ Decision-tree analysis (supervised, top-down split)
  - „ Correlation (e.g.,  $\chi^2$ ) analysis (unsupervised, bottom-up merge)

# Simple Discretization: Binning

---

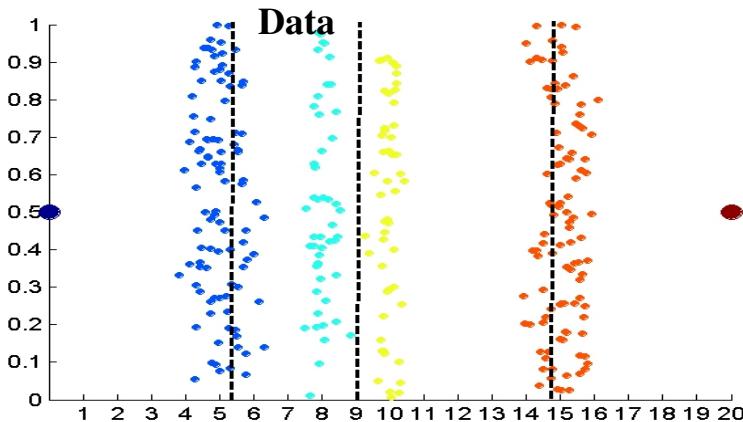
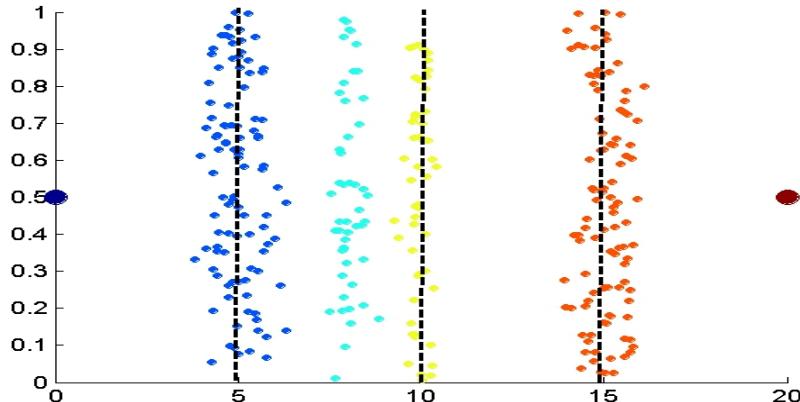
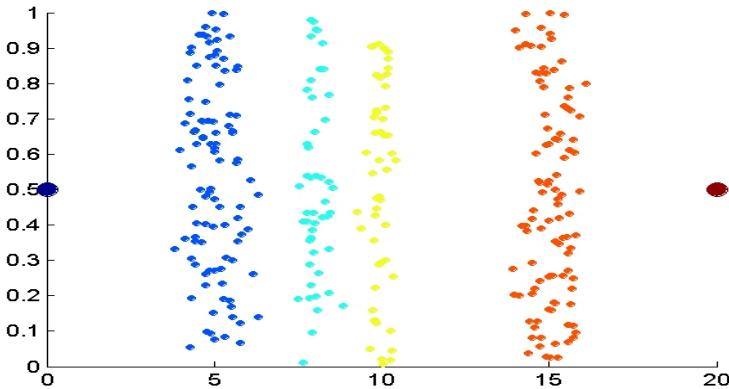
- n Equal-width (distance) partitioning
  - n Divides the range into  $N$  intervals of equal size: uniform grid
  - n if  $A$  and  $B$  are the lowest and highest values of the attribute, the width of intervals will be:  $W = (B - A)/N$ .
  - n The most straightforward, but outliers may dominate presentation
  - n Skewed data is not handled well
- n Equal-depth (frequency) partitioning
  - n Divides the range into  $N$  intervals, each containing approximately same number of samples
  - n Good data scaling
  - n Managing categorical attributes can be tricky

# Binning Methods for Data Smoothing

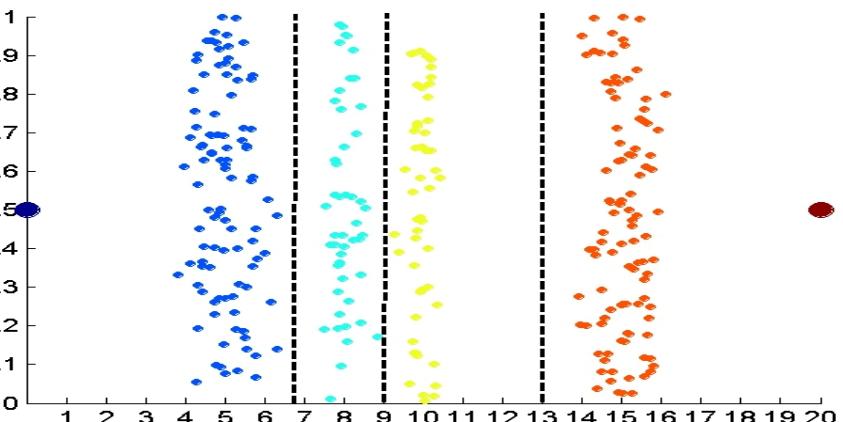
---

- q Sorted data for price (in dollars): 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34
- \* Partition into equal-frequency (**equi-depth**) bins:
  - Bin 1: 4, 8, 9, 15
  - Bin 2: 21, 21, 24, 25
  - Bin 3: 26, 28, 29, 34
- \* Smoothing by **bin means**:
  - Bin 1: 9, 9, 9, 9
  - Bin 2: 23, 23, 23, 23
  - Bin 3: 29, 29, 29, 29
- \* Smoothing by **bin boundaries**:
  - Bin 1: 4, 4, 4, 15
  - Bin 2: 21, 21, 25, 25
  - Bin 3: 26, 26, 26, 34

# Discretization Without Using Class Labels (Binning vs. Clustering)



Equal frequency (binning)



K-means clustering leads to better results

# Discretization by Classification & Correlation Analysis

---

- Classification (e.g., decision tree analysis)
  - Supervised: Given class labels, e.g., cancerous vs. benign
  - Using *entropy* to determine split point (discretization point)
  - Top-down, recursive split
  - Details to be covered in Chapter 7
- Correlation analysis (e.g., Chi-merge:  $\chi^2$ -based discretization)
  - Supervised: use class information
  - Bottom-up merge: find the best neighboring intervals (those having similar distributions of classes, i.e., low  $\chi^2$  values) to merge
  - Merge performed recursively, until a predefined stopping

# Concept Hierarchy Generation

---

- **Concept hierarchy** organizes concepts (i.e., attribute values) hierarchically and is usually associated with each dimension in a data warehouse
- Concept hierarchies facilitate drilling and rolling in data warehouses to view data in multiple granularity
- Concept hierarchy formation: Recursively reduce the data by collecting and replacing low level concepts (such as numeric values for *age*) by higher level concepts (such as *youth*, *adult*, or *senior*)
- Concept hierarchies can be explicitly specified by domain experts and/or data warehouse designers
- Concept hierarchy can be automatically formed for both numeric and nominal data. For numeric data, use discretization methods shown.

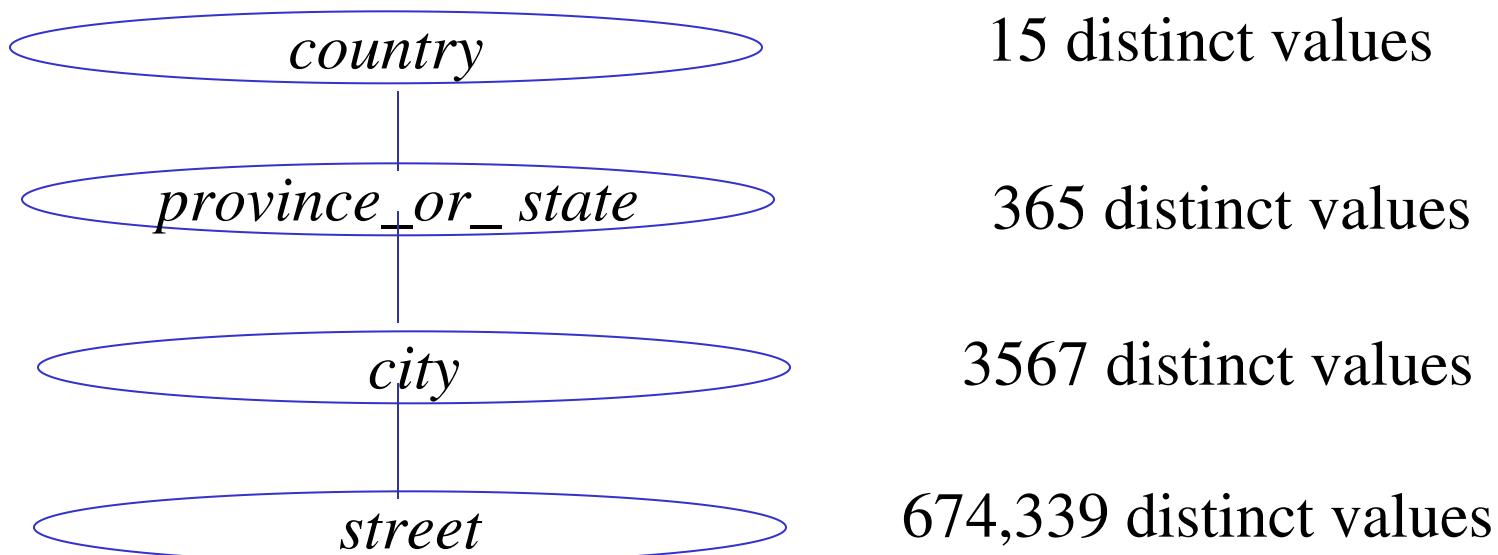
# Concept Hierarchy Generation for Nominal Data

---

- Specification of a partial/total ordering of attributes explicitly at the schema level by users or experts
  - *street < city < state < country*
- Specification of a hierarchy for a set of values by explicit data grouping
  - $\{\text{Urbana, Champaign, Chicago}\} < \text{Illinois}$
- Specification of only a partial set of attributes
  - E.g., only *street < city*, not others
- Automatic generation of hierarchies (or attribute levels) by the analysis of the number of distinct values
  - E.g., for a set of attributes:  $\{\text{street, city, state, country}\}$

# Automatic Concept Hierarchy Generation

- n Some hierarchies can be automatically generated based on the analysis of the number of distinct values per attribute in the data set
  - n The attribute with the most distinct values is placed at the lowest level of the hierarchy
  - n Exceptions, e.g., weekday, month, quarter, year



# Chapter 3: Data Preprocessing

---

- n Data Preprocessing: An Overview
    - n Data Quality
    - n Major Tasks in Data Preprocessing
  - n Data Cleaning
  - n Data Integration
  - n Data Reduction
  - n Data Transformation and Data Discretization
  - n Summary
- 

# Summary

---

- n **Data quality**: accuracy, completeness, consistency, timeliness, believability, interpretability
- n **Data cleaning**: e.g. missing/noisy values, outliers
- n **Data integration** from multiple sources:
  - n Entity identification problem
  - n Remove redundancies
  - n Detect inconsistencies
- n **Data reduction**
  - n Dimensionality reduction
  - n Numerosity reduction
  - n Data compression
- n **Data transformation and data discretization**
  - n Normalization
  - n Concept hierarchy generation

# References

- n D. P. Ballou and G. K. Tayi. Enhancing data quality in data warehouse environments. Comm. of ACM, 42:73-78, 1999
- n A. Bruce, D. Donoho, and H.-Y. Gao. Wavelet analysis. *IEEE Spectrum*, Oct 1996
- n T. Dasu and T. Johnson. *Exploratory Data Mining and Data Cleaning*. John Wiley, 2003
- n J. Devore and R. Peck. *Statistics: The Exploration and Analysis of Data*. Duxbury Press, 1997.
- n H. Galhardas, D. Florescu, D. Shasha, E. Simon, and C.-A. Saita. Declarative data cleaning: Language, model, and algorithms. *VLDB'01*
- n M. Hua and J. Pei. Cleaning disguised missing data: A heuristic approach. *KDD'07*
- n H. V. Jagadish, et al., *Special Issue on Data Reduction Techniques*. Bulletin of the Technical Committee on Data Engineering, 20(4), Dec. 1997
- n H. Liu and H. Motoda (eds.). *Feature Extraction, Construction, and Selection: A Data Mining Perspective*. Kluwer Academic, 1998
- n J. E. Olson. *Data Quality: The Accuracy Dimension*. Morgan Kaufmann, 2003
- n D. Pyle. *Data Preparation for Data Mining*. Morgan Kaufmann, 1999
- n V. Raman and J. Hellerstein. *Potters Wheel: An Interactive Framework for Data Cleaning and Transformation*, VLDB'2001
- n T. Redman. *Data Quality: The Field Guide*. Digital Press (Elsevier), 2001
- n R. Wang, V. Storey, and C. Firth. A framework for analysis of data quality research. *IEEE Trans. Knowledge and Data Engineering*, 7:623-640, 1995