

# **Data Mining: Data**

Lecture Notes for Chapter 2

Introduction to Data Mining

by

Tan, Steinbach, Kumar

(modified by Predrag Radivojac, 2018)

# 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

The diagram illustrates a dataset consisting of 10 data objects. Each object is represented by a row in a table. The columns represent attributes: **Tid**, **Refund**, **Marital Status**, **Taxable Income**, and **Cheat**. The **Objects** are grouped by a bracket on the left, and the **Attributes** and **Class** are grouped by brackets at the top right.

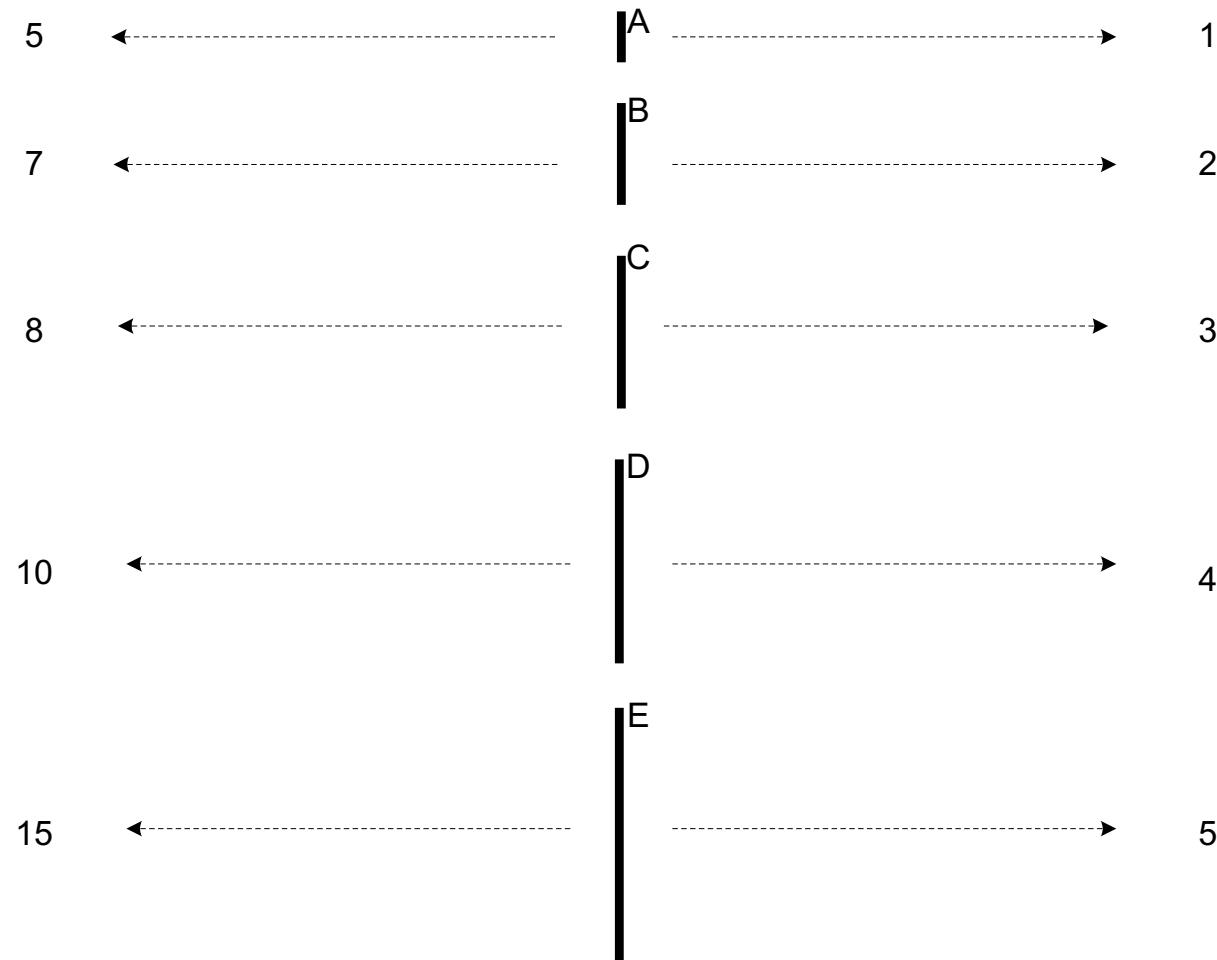
Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

# Attribute Values

- Attribute values are numbers or symbols assigned to an attribute
- Distinction between attributes and attribute values
  - Same attribute can be mapped to different attribute values
    - ◆ Example: height can be measured in feet or meters
  - Different attributes can be mapped to the same set of values
    - ◆ Example: Attribute values for ID and age are integers
    - ◆ But properties of attribute values can be different
      - ID has no limit but age has a maximum and minimum value

# Measurement of Length

- The way you measure an attribute is something that may not match the attributes properties



# Types of Attributes

- There are different types of attributes
  - Nominal
    - ◆ Examples: ID numbers, eye color, zip codes
  - Ordinal
    - ◆ Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height in {tall, medium, short}
  - Interval
    - ◆ Examples: calendar dates, temperatures in Celsius or Fahrenheit.
  - Ratio
    - ◆ Examples: temperature in Kelvin, length, time, counts

# Properties of Attribute Values

- The type of an attribute depends on which of the following properties it possesses:
  - Distinctness:                $= \neq$
  - 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	<p>An order preserving change of values, i.e.,</p> $\text{new\_value} = f(\text{old\_value})$ <p>where <math>f</math> is a monotonic function.</p>	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	$\text{new\_value} = a * \text{old\_value} + b$ <p>where a and b are constants</p>	Thus, the Fahrenheit and Celsius temperature scales differ in terms of where their zero value is and the size of a unit (degree).
Ratio	$\text{new\_value} = a * \text{old\_value}$	Length can be measured in meters or feet.

# Discrete and Continuous Attributes

## ● Discrete Attribute

- Has only a finite or countably infinite set of values
- Examples: zip codes, counts, or the set of words in a collection of documents
- Often represented as integer variables.
- Note: binary attributes are a special case of discrete attributes

## ● Continuous Attribute

- Has real numbers as attribute values
- Examples: 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.

# Types of data sets

## ● Record

- Data Matrix
- Document Data
- Transaction Data

## ● Graph

- World Wide Web
- Molecular Structures

## ● Ordered

- Spatial Data
- Temporal Data
- Sequential Data
- Genetic Sequence Data

# Important Characteristics of Structured Data

- **Dimensionality**

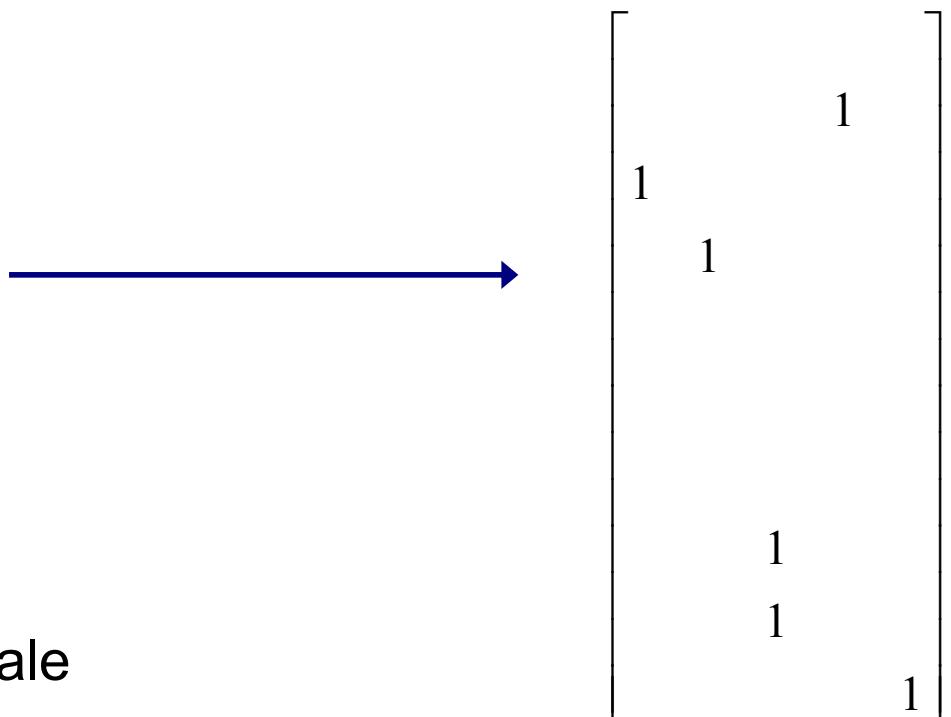
- ◆ Curse of Dimensionality

- **Sparsity**

- ◆ Only presence counts

- **Resolution**

- ◆ Patterns depend on the scale



# Record Data

- Data that consists of a collection of records, each of which consists of a fixed set of attributes

<i>Tid</i>	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

# Data Matrix

- If data objects have the same fixed set of numeric attributes, then the data objects can be thought of as points in a multi-dimensional space, where each dimension represents a distinct attribute
- Such data set can be represented by an  $m$ -by- $n$  matrix, where there are  $m$  rows, one for each object, and  $n$  columns, one for each attribute

Projection of x Load	Projection of y load	Distance	Load	Thickness
10.23	5.27	15.22	2.7	1.2
12.65	6.25	16.22	2.2	1.1

# Document Data

- Each document becomes a 'term' vector,
  - each term is a component (attribute) of the vector,
  - the value of each component is the number of times the corresponding term occurs in the document.

	team	coach	play	ball	score	game	winn	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

# Example

ESPN: The Worldwide Leader

www.espn.com

TOP EVENTS Tennis (M) Final 1 R. Nadal 7<sup>8</sup> 6 6 D. Lajovic 6<sup>6</sup> 2 2

NFL NBA MLB NCAAF Soccer Tennis ...

Quick Links US Open 2017

Arizona for Super Bowl - NFL - Yahoo! Sports - Mozilla ...

sports.yahoo.com/nfl/news;\_ylt=AjrnNjsD... Google

Dictionary Yahoo! Informatics Pedja MapQuest

VIS | Middle East | UN says Ga...

Golf Tennis Soccer MMA Boxing All Sports Shop

STATS TEAMS PLAYERS ODDS VIDEO RUMORS BLOG

## headed to Arizona for Super Bowl

last hurdle: a New York team of road warriors hoping for a

," said Tom Brady, the Patriots' dimple-chinned,

Brady and the Patriots (18-0) will try to match the 1972 Miami Dolphins as the only teams to complete an undefeated season when they face Eli Manning and the Giants on Feb. 3 in the Super Bowl at Glendale, Ariz.

"I think you enter the season and you're hoping to put together a bunch of great wins and you realize there's

Final 3 R. Federer 4 6 6 1 6 F. Tiafoe 6 2 1 6 4

Tennis (W) Final 1 K. Pliskova 6 6 M. Linette 2 1

Final N. Osaka 6 6 6 A. Kerber 3 1



### Tests, turns and concerns in the Kyrie Irving trade

Teams are offering the Cavs a lot for the superstar. Zach Lowe examines why that could work well.



#### Isaiah Thomas on health: 'I am not damaged'

Thomas had a message for those doubting his health after his trade to Cleveland stalled, telling ESPN on Tuesday that his injured hip "won't be a problem in the future" and that he is confident he'll be the same player.

#### vs lowering trade compensation demands.

Adrian Wojnarowski reports that Thomas is no longer asking for an "elite" player from Boston after Thomas' physical.



# Recommendations Data

## ● Sparse matrix

- each row is a person
- each column is a movie (book, disease, ...)
- each number is a rating

The diagram illustrates a sparse matrix structure. A horizontal red double-headed arrow labeled "Movies" spans the width of the matrix, indicating the columns. A vertical red double-headed arrow labeled "Persons" spans the height of the matrix, indicating the rows. The matrix itself is a grid of gray cells. The columns are labeled at the top with movie titles: Spiderman, Ocean's 11, Matrix, Titanic, JFK, Star wars, Creed, and Rocky. The rows are labeled on the left with person names: Person 1, Person 2, Person 3, and Person 4. Numerical values are placed in specific cells to represent ratings: Person 1 rated Spiderman a 3 and Titanic a 4; Person 2 rated Rocky a 5; Person 3 rated Star wars a 4 and Rocky a 5; and Person 4 rated Ocean's 11 a 1 and Titanic a 3.

	Spiderman	Ocean's 11	Matrix	Titanic	JFK	Star wars	Creed	Rocky
Person 1		3		4				
Person 2								5
Person 3						4		5
Person 4	1		3				2	

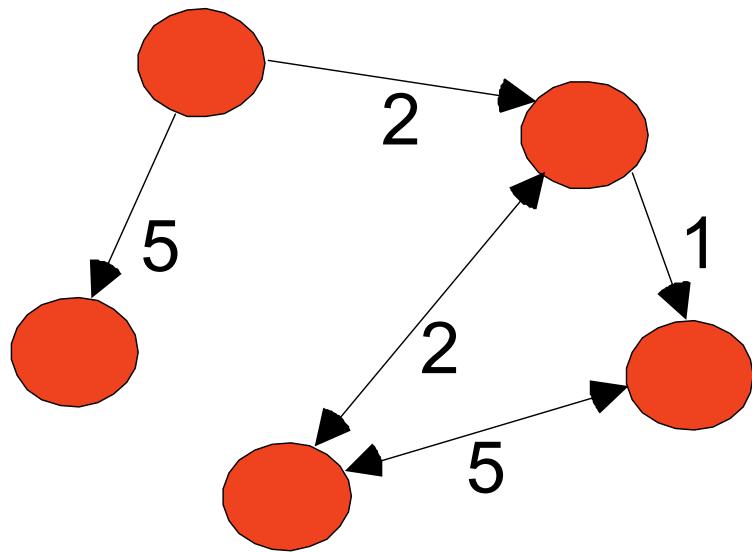
# Transaction Data

- A special type of record data, where
  - each record (transaction) involves a set of items.
  - For example, consider a grocery store. The set of products purchased by a customer during one shopping trip constitute a transaction, while the individual products that were purchased are the items.

<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

# Graph Data

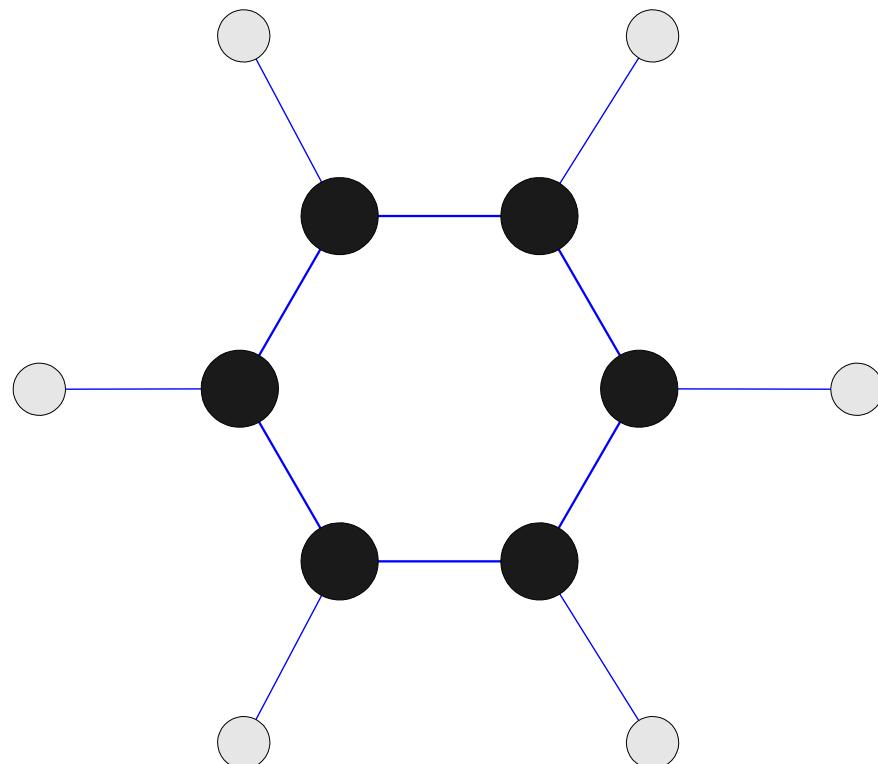
- Examples: Generic graph and HTML Links



```
<a href="papers/papers.html#bbbb">  
Data Mining </a>  
<li>  
<a href="papers/papers.html#aaaa">  
Graph Partitioning </a>  
<li>  
<a href="papers/papers.html#aaaa">  
Parallel Solution of Sparse Linear System of Equations </a>  
<li>  
<a href="papers/papers.html#ffff">  
N-Body Computation and Dense Linear System Solvers
```

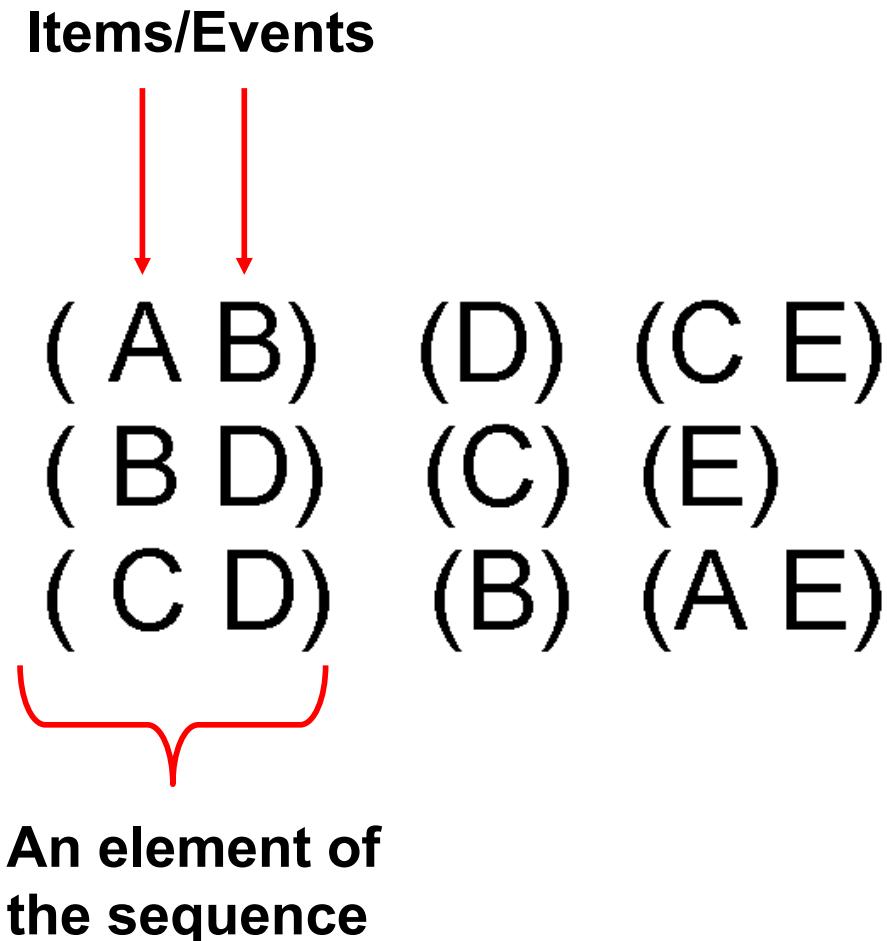
# Chemical Data

- Benzene Molecule: C<sub>6</sub>H<sub>6</sub>



# Ordered Data

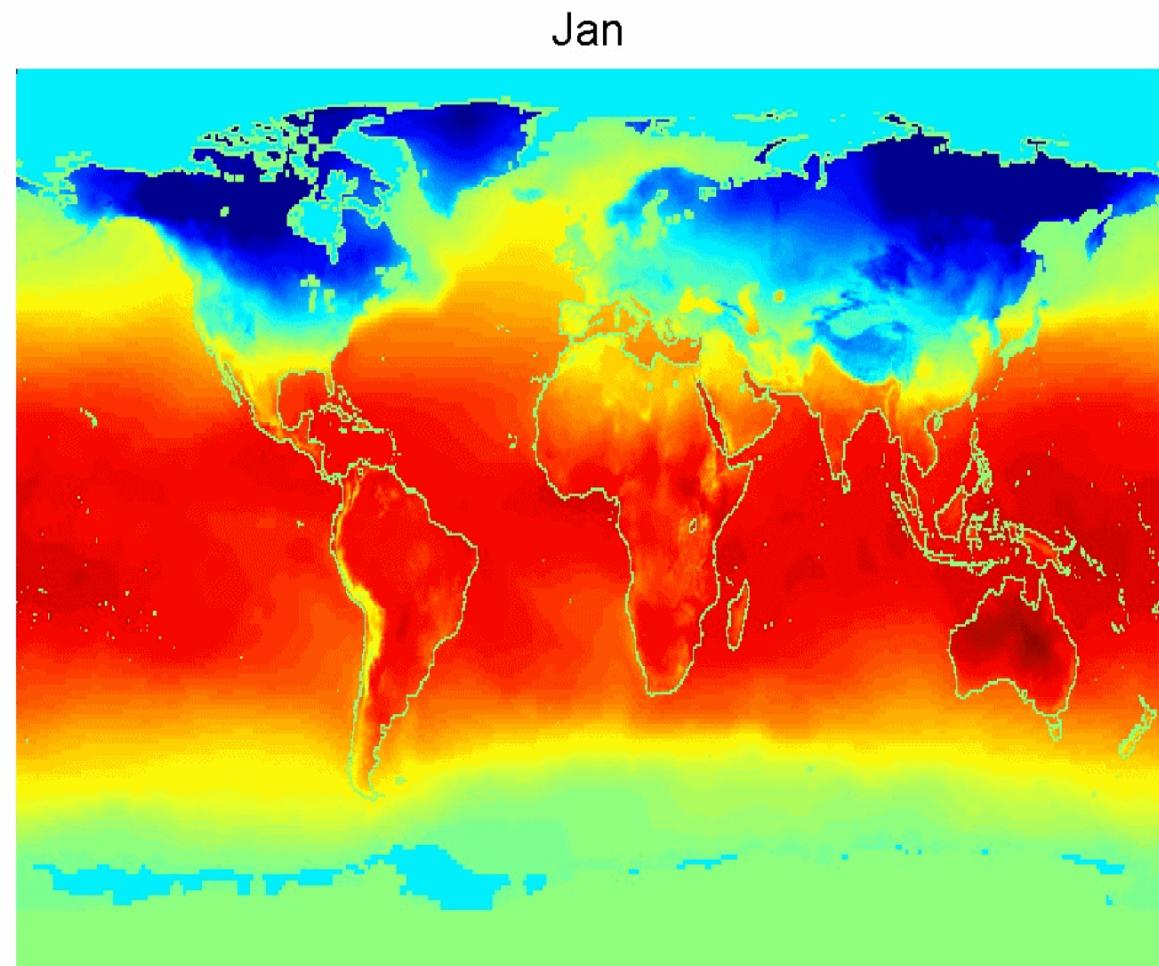
- Sequences of transactions



# Ordered Data

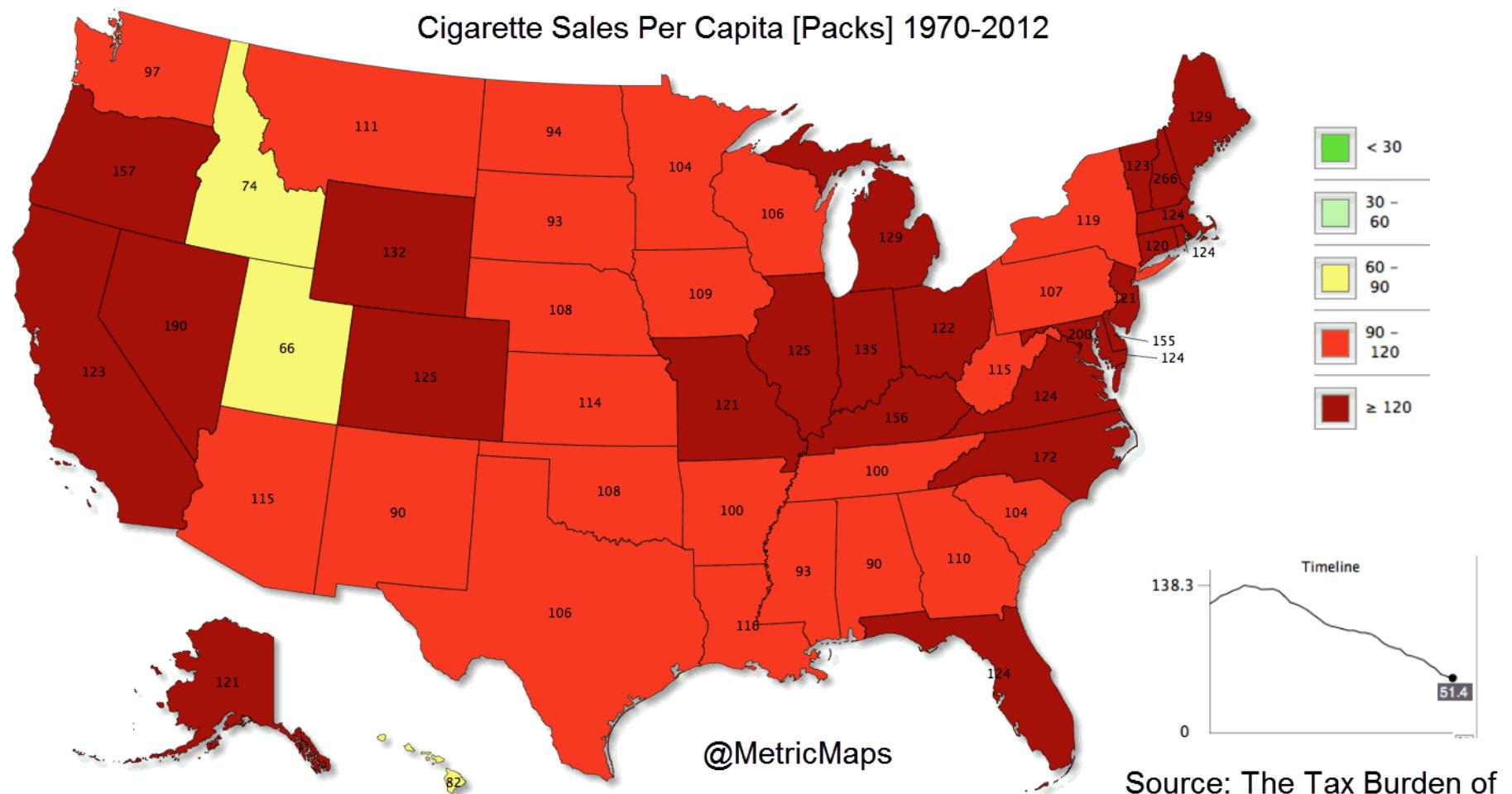
- Spatio-Temporal Data

Average monthly temperature of land and ocean



# Ordered Data

## ● Spatio-Temporal Data

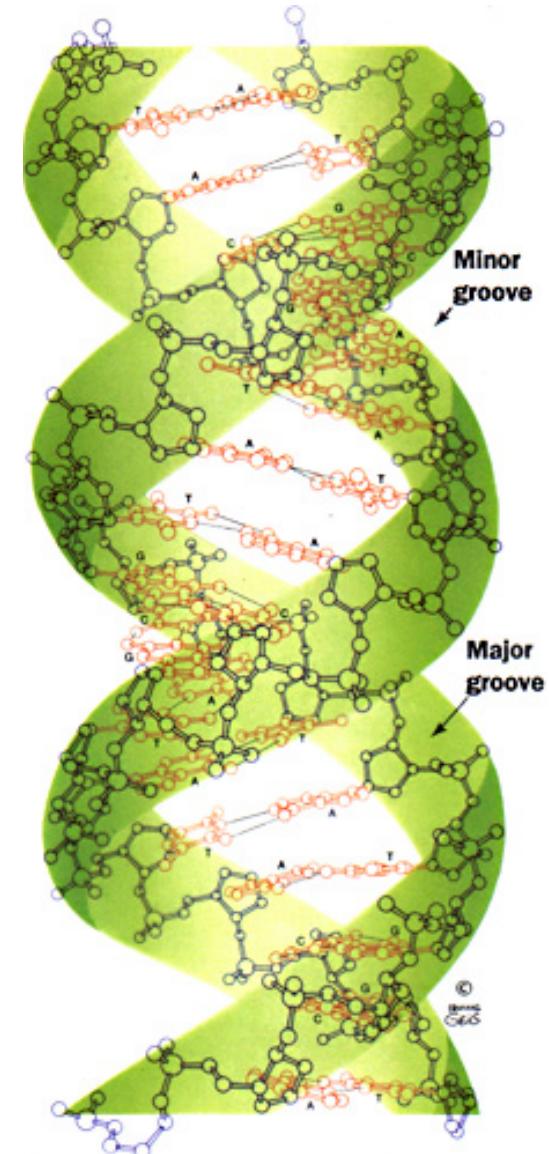


1970

# Ordered/Sequence Data

- Genomic sequence data

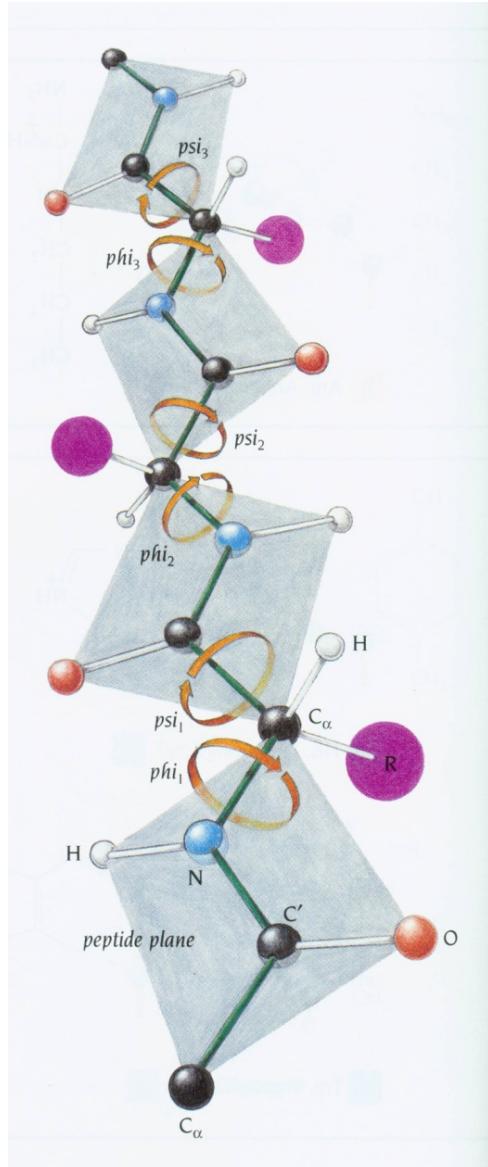
1:	...GGTTCCGGCCTTCAGCCCCCGGCC...	0
2:	...GGTTCCGGCGTTTCAGCCCCGCGGCC...	1
3:	...GGTTCCGGCCTTCAGCCCCCGGCC...	0
4:	...GGTTCCGGCCTTCAGCCCCGCGGCC...	0
5:	...GGTTCCGGCCTTCAGCCCCCTCGGCC...	0
6:	...GGTTCCGGCCTTCAGCCCCGCGGCC...	0
7:	...GGTTCCGGCCTTCAGCCCCCTCGGCC...	0
8:	...GGTTCCGGCATTCAGCCCCCGGCC...	1
9:	...GGTTCCGGCCTTCAGCCCCGCGGCC...	0



# Machine Learning Repository at UCI

- contains a number of user deposited ML problems
- <ftp://ftp.ics.uci.edu/pub/machine-learning-databases>
- Discussion:
  - Pima Indians diabetes example ([link](#))
  - Boston housing example ([link](#))
  - German credit example ([link](#))

# Reading Custom File Types



- Need standard I/O for this
  - use fopen, fclose, fgetl, fget for text files
  - use fread, fwrite, fseek, ftell for binary files
  
- Examples
  - reading protein sequence data
  - reading and writing binary data

# Similarity and Dissimilarity

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

# Similarity/Dissimilarity for Simple Attributes

$p$  and  $q$  are the attribute values for two data objects.

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$	$s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$
Ordinal	$d = \frac{ p-q }{n-1}$ (values mapped to integers 0 to $n-1$ , where $n$ is the number of values)	$s = 1 - \frac{ p-q }{n-1}$
Interval or Ratio	$d =  p - q $	$s = -d, s = \frac{1}{1+d}$ or $s = 1 - \frac{d - \min_d}{\max_d - \min_d}$

**Table 5.1.** Similarity and dissimilarity for simple attributes

# Euclidean Distance

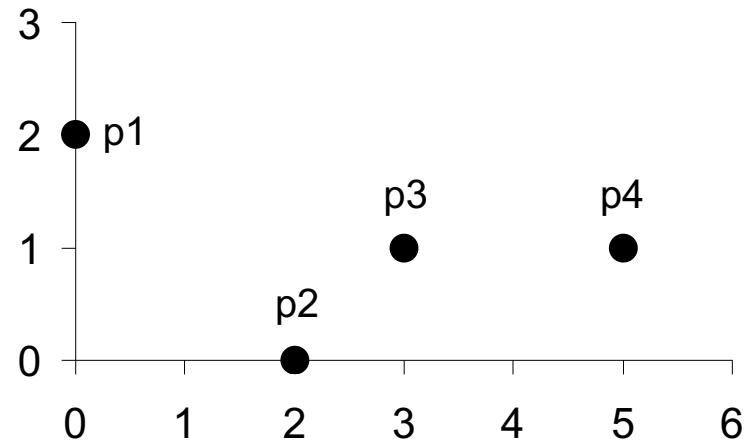
- Euclidean Distance

$$dist = \sqrt{\sum_{k=1}^n (p_k - q_k)^2}$$

Where  $n$  is the number of dimensions (attributes) and  $p_k$  and  $q_k$  are, respectively, the  $k^{\text{th}}$  attributes (components) or data objects  $p$  and  $q$ .

- Standardization is necessary, if scales differ.

# Euclidean Distance



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Distance Matrix

# Minkowski Distance

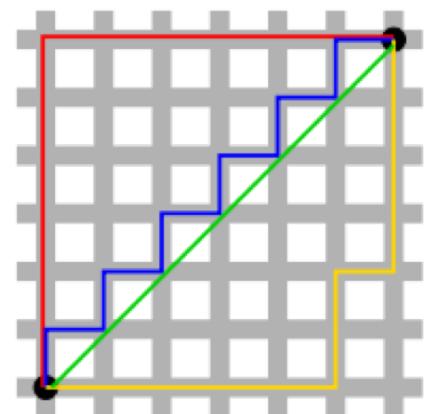
- Minkowski Distance is a generalization of Euclidean Distance

$$dist = \left( \sum_{k=1}^n |p_k - q_k|^r \right)^{\frac{1}{r}}$$

Where  $r$  is a parameter,  $n$  is the number of dimensions (attributes) and  $p_k$  and  $q_k$  are, respectively, the  $k$ th attributes (components) or data objects  $p$  and  $q$ .

# Minkowski Distance: Examples

- $r = 1$ . City block (Manhattan, taxicab,  $L_1$  norm) distance.
  - A common example of this is the Hamming distance, which is just the number of bits that are different between two binary vectors
- $r = 2$ . Euclidean distance
- $r \rightarrow \infty$ . “supremum” ( $L_{\max}$  norm,  $L_\infty$  norm) distance.
  - This is the maximum difference between any component of the vectors
- Do not confuse  $r$  with  $n$ , i.e., all these distances are defined for all numbers of dimensions.



# Minkowski Distance

point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

L1	p1	p2	p3	p4
p1	0	4	4	6
p2	4	0	2	4
p3	4	2	0	2
p4	6	4	2	0

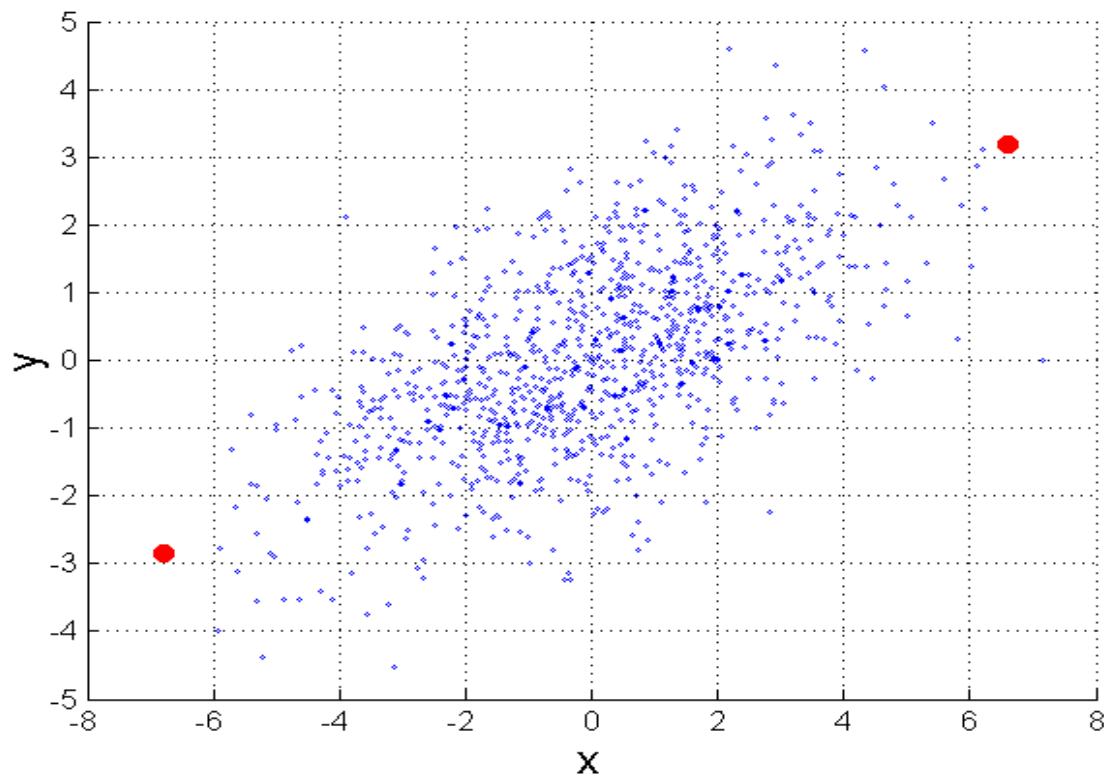
L2	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

L $\infty$	p1	p2	p3	p4
p1	0	2	3	5
p2	2	0	1	3
p3	3	1	0	2
p4	5	3	2	0

Distance Matrix

# Mahalanobis Distance

$$d_M(p, q) = (p - q) \Sigma^{-1} (p - q)^T$$

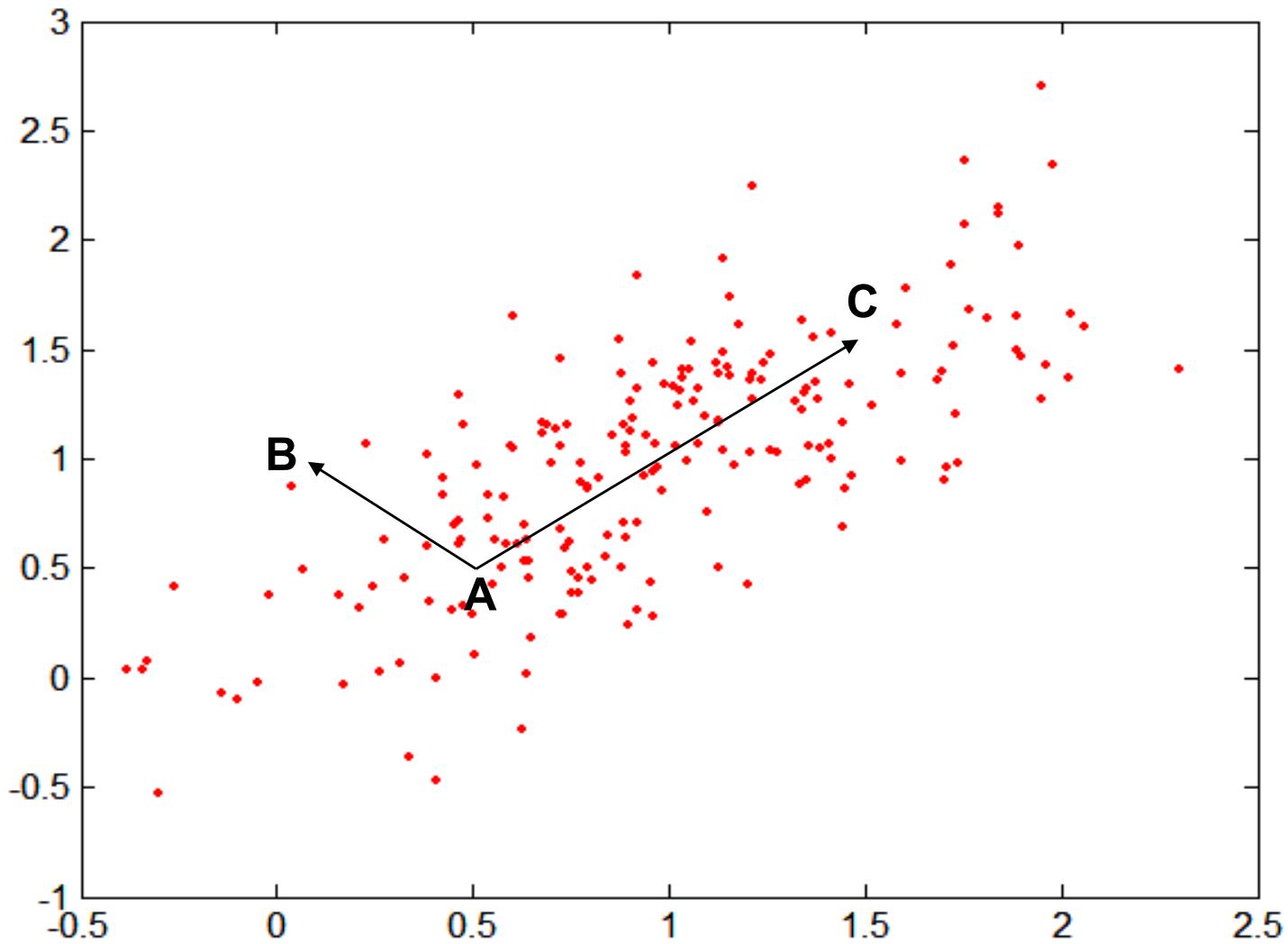


$\Sigma$  is the covariance matrix of the input data  $X$

$$\Sigma_{j,k} = \frac{1}{n-1} \sum_{i=1}^n (X_{ij} - \bar{X}_j)(X_{ik} - \bar{X}_k)$$

For red points, the Euclidean distance is 14.7, Mahalanobis distance is 6.

# Mahalanobis Distance



Covariance Matrix:

$$\Sigma = \begin{bmatrix} 0.3 & 0.2 \\ 0.2 & 0.3 \end{bmatrix}$$

A: (0.5, 0.5)

B: (0, 1)

C: (1.5, 1.5)

$\text{Mahal}(A,B) = 5$

$\text{Mahal}(A,C) = 4$

# Common Properties of a Distance

- Distances, such as the Euclidean distance, have some well known properties.
  1.  $d(p, q) \geq 0$  for all  $p$  and  $q$  and  $d(p, q) = 0$  only if  $p = q$ . (Positive definiteness)
  2.  $d(p, q) = d(q, p)$  for all  $p$  and  $q$ . (Symmetry)
  3.  $d(p, r) \leq d(p, q) + d(q, r)$  for all points  $p$ ,  $q$ , and  $r$ . (Triangle Inequality)

where  $d(p, q)$  is the distance (dissimilarity) between points (data objects),  $p$  and  $q$ .
- A distance that satisfies these properties is a **metric**

# Common Properties of a Similarity

- Similarities, also have some well known properties.
    1.  $s(p, q) = 1$  (or maximum similarity) only if  $p = q$ .
    2.  $s(p, q) = s(q, p)$  for all  $p$  and  $q$ . (Symmetry)
- where  $s(p, q)$  is the similarity between points (data objects),  $p$  and  $q$ .

# Similarity Between Binary Vectors

- Common situation is that objects,  $p$  and  $q$ , have only binary attributes
- Compute similarities using the following quantities

$M_{01}$  = the number of attributes where  $p$  was 0 and  $q$  was 1

$M_{10}$  = the number of attributes where  $p$  was 1 and  $q$  was 0

$M_{00}$  = the number of attributes where  $p$  was 0 and  $q$  was 0

$M_{11}$  = the number of attributes where  $p$  was 1 and  $q$  was 1

- Simple Matching and Jaccard Coefficients

SMC = number of matches / number of attributes

$$= (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00})$$

J = number of 11 matches / number of not-both-zero attributes values

$$= (M_{11}) / (M_{01} + M_{10} + M_{11})$$

# SMC versus Jaccard: Example

$p = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0$

$q = 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1$

$M_{01} = 2$  (the number of attributes where p was 0 and q was 1)

$M_{10} = 1$  (the number of attributes where p was 1 and q was 0)

$M_{00} = 7$  (the number of attributes where p was 0 and q was 0)

$M_{11} = 0$  (the number of attributes where p was 1 and q was 1)

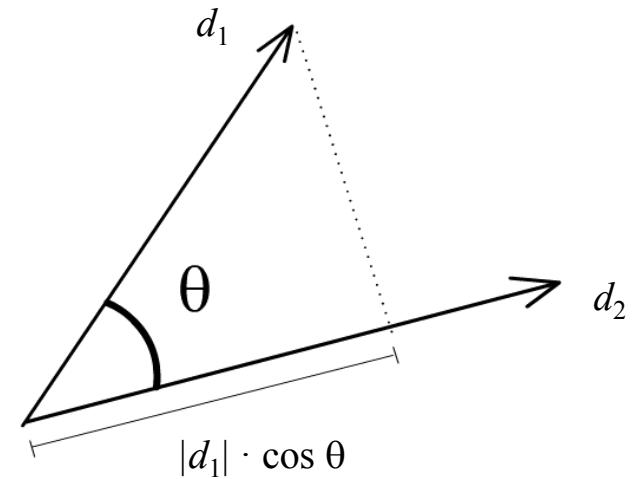
$$\text{SMC} = (M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00}) = (0+7) / (2+1+0+7) = 0.7$$

$$J = (M_{11}) / (M_{01} + M_{10} + M_{11}) = 0 / (2 + 1 + 0) = 0$$

# Cosine Similarity

- If  $d_1$  and  $d_2$  are two document vectors, then

$$\cos(d_1, d_2) = \frac{d_1 \cdot d_2}{\|d_1\| \cdot \|d_2\|}$$



where • indicates vector dot product and  $\| d \|$  is the length of vector  $d$ .

- Example:

$$d_1 = 3 \ 2 \ 0 \ 5 \ 0 \ 0 \ 0 \ 2 \ 0 \ 0$$

$$d_2 = 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 2$$

$$d_1 \cdot d_2 = 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5$$

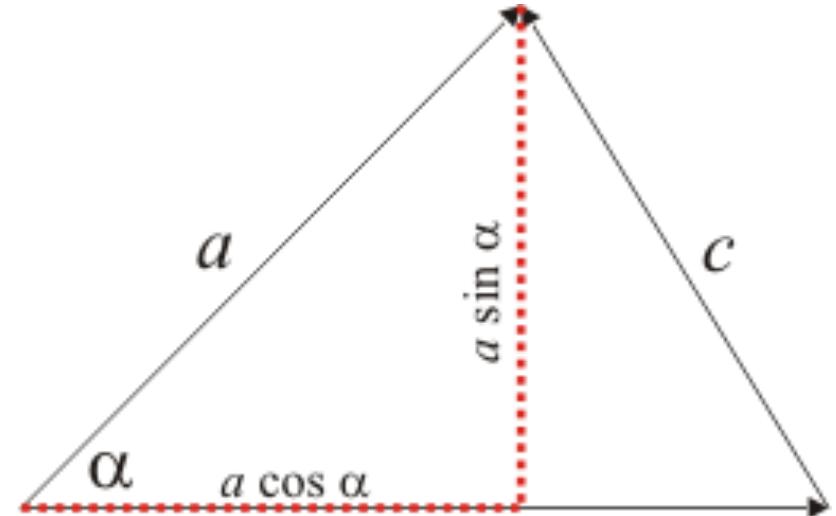
$$\|d_1\| = (3^2 + 2^2 + 0^2 + 5^2 + 0^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2)^{0.5} = (42)^{0.5} = 6.481$$

$$\|d_2\| = (1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1^2 + 0^2 + 2^2)^{0.5} = (6)^{0.5} = 2.245$$

$$\cos(d_1, d_2) = .3150$$

# Proof:

$$\begin{aligned}c^2 &= (b - a \cos \alpha)^2 + (a \sin \alpha)^2 \\&= b^2 - 2ab \cos \alpha + a^2 \cos^2 \alpha + a^2 \sin^2 \alpha \\&= a^2 + b^2 - 2ab \cos \alpha\end{aligned}$$



$$\begin{aligned}c^2 &= \vec{c} \cdot \vec{c} \\&= (\vec{a} - \vec{b}) \cdot (\vec{a} - \vec{b}) \\&= \vec{a} \cdot \vec{a} - 2\vec{a} \cdot \vec{b} + \vec{b} \cdot \vec{b} \\&= a^2 - 2\vec{a} \cdot \vec{b} + b^2\end{aligned}$$

$$\cos \alpha = \frac{\vec{a} \cdot \vec{b}}{ab}$$

By combining previous equations we get:

$$a^2 + b^2 - 2ab \cos \alpha = a^2 - 2\vec{a} \cdot \vec{b} + b^2$$

# Correlation

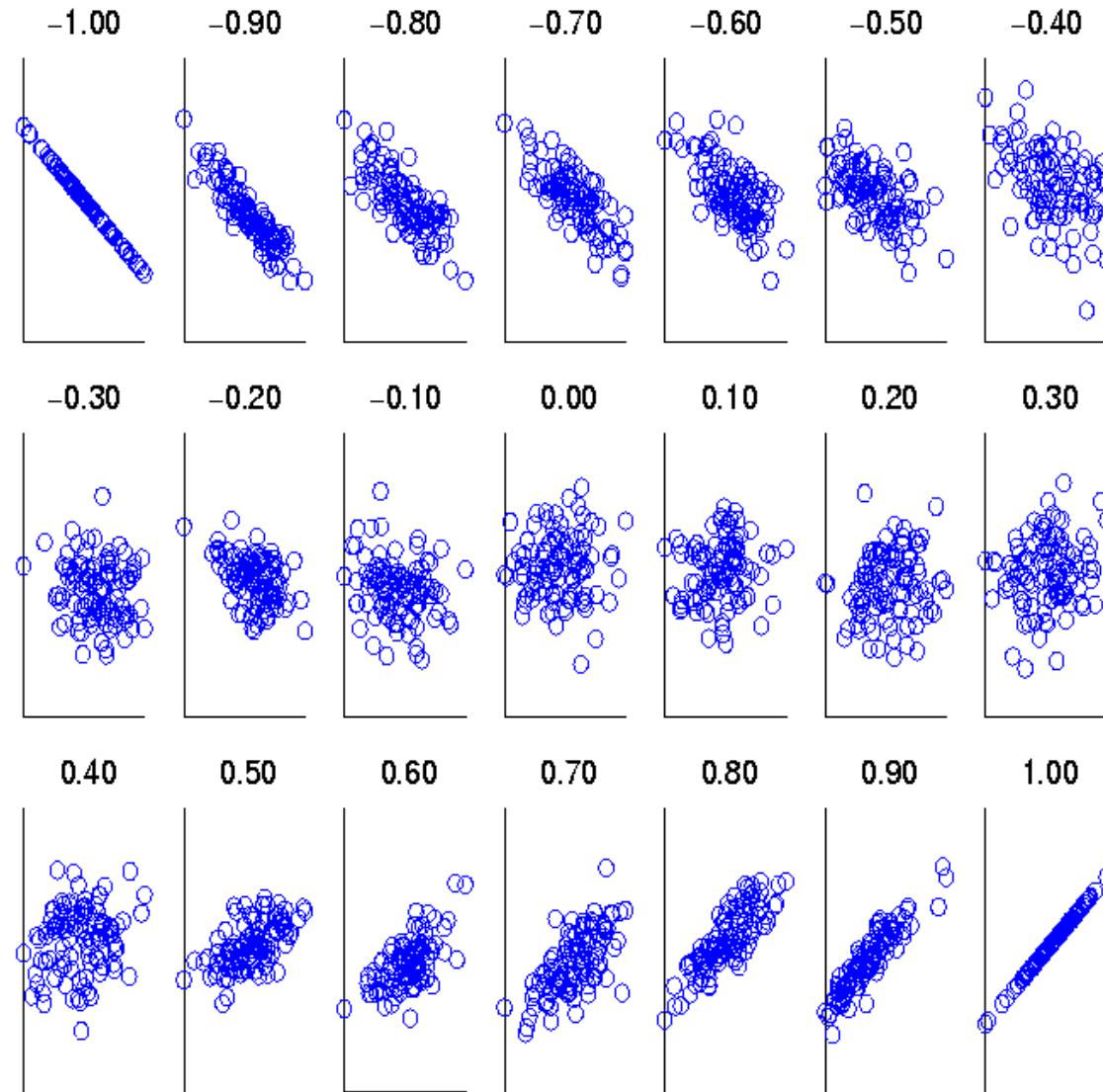
- Correlation measures the linear relationship between objects
- To compute correlation, we standardize data objects,  $p$  and  $q$ , and then take their dot product

$$p'_k = (p_k - \text{mean}(p)) / \text{std}(p)$$

$$q'_k = (q_k - \text{mean}(q)) / \text{std}(q)$$

$$\text{correlation}(p, q) = \frac{p' \bullet q'}{n - 1}$$

# Visually Evaluating Correlation



Scatter plots showing the similarity from –1 to 1.

# General Approach for Combining Similarities

- Sometimes attributes are of many different types, but an overall similarity is needed.
  1. For the  $k^{th}$  attribute, compute a similarity,  $s_k$ , in the range [0, 1].
  2. Define an indicator variable,  $\delta_k$ , for the  $k_{th}$  attribute as follows:
$$\delta_k = \begin{cases} 0 & \text{if the } k^{th} \text{ attribute is a binary asymmetric attribute and both objects have} \\ & \text{a value of 0, or if one of the objects has a missing values for the } k^{th} \text{ attribute} \\ 1 & \text{otherwise} \end{cases}$$
  3. Compute the overall similarity between the two objects using the following formula:

$$similarity(p, q) = \frac{\sum_{k=1}^n \delta_k s_k}{\sum_{k=1}^n \delta_k}$$

# Using Weights to Combine Similarities

- May not want to treat all attributes the same.
  - Use weights  $w_k$  which are between 0 and 1 and sum to 1.

$$\text{similarity}(p, q) = \frac{\sum_{k=1}^n w_k \delta_k s_k}{\sum_{k=1}^n \delta_k}$$

$$\text{distance}(p, q) = \left( \sum_{k=1}^n w_k |p_k - q_k|^r \right)^{1/r}.$$

# Data Quality

- What kinds of data quality problems?
  - How can we detect problems with the data?
  - What can we do about these problems?
- 
- Examples of data quality problems:
    - noise and outliers
    - missing values
    - duplicate data

# Why Is Data Dirty?

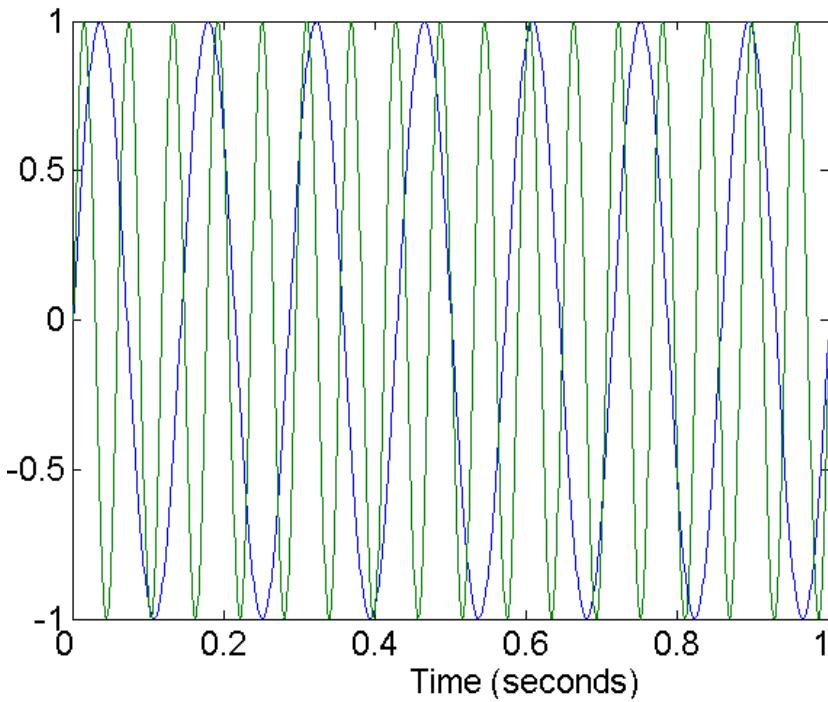
- Incomplete data may come from
  - “Not applicable” data value when collected
  - Different considerations between the time when the data was collected and when it is analyzed.
  - Human/hardware/software problems
- Noisy data (incorrect values) may come from
  - Faulty data collection instruments
  - Human or computer error at data entry
  - Errors in data transmission
- Inconsistent data may come from
  - Different data sources
  - Functional dependency violation (e.g., modify some linked data)
- Duplicate records also need data cleaning

# Why Is Data Preprocessing Important?

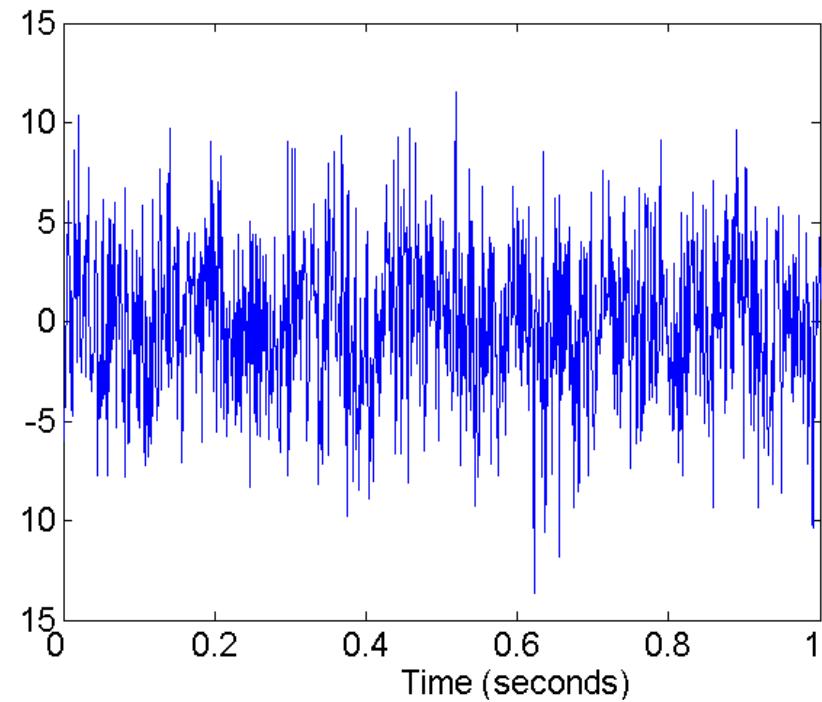
- **No quality data, no quality mining results!**
  - Quality decisions must be based on quality data
    - ◆ e.g., duplicate or missing data may cause incorrect or even misleading statistics.
  - Data warehouse needs consistent integration of quality data
- Data extraction, cleaning, and transformation comprises the majority of the work of building a data warehouse

# Noise

- Noise refers to modification of original values
  - Examples: distortion of a person's voice when talking on a poor phone and “snow” on television screen



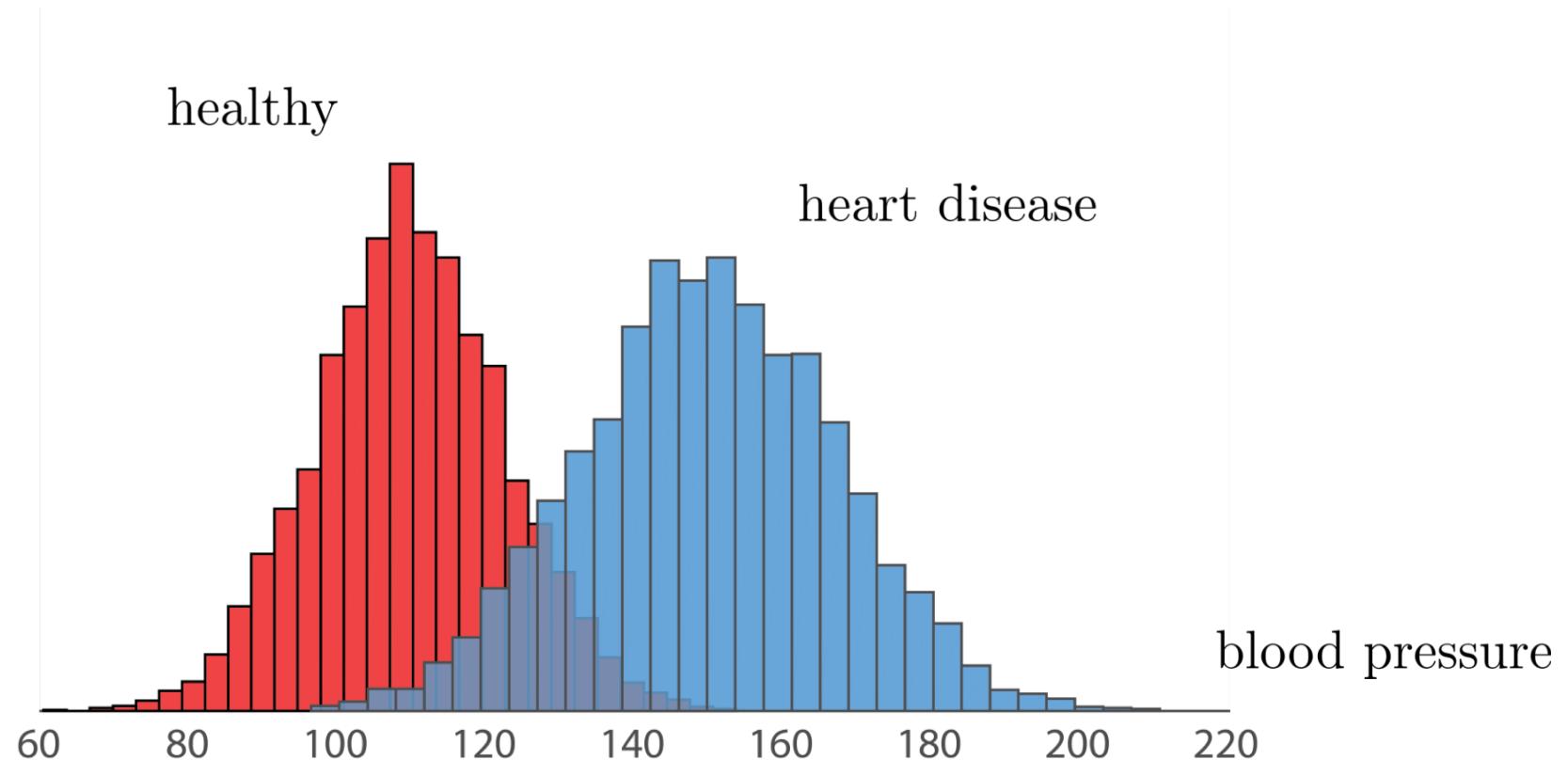
Two Sine Waves



Two Sine Waves + Noise

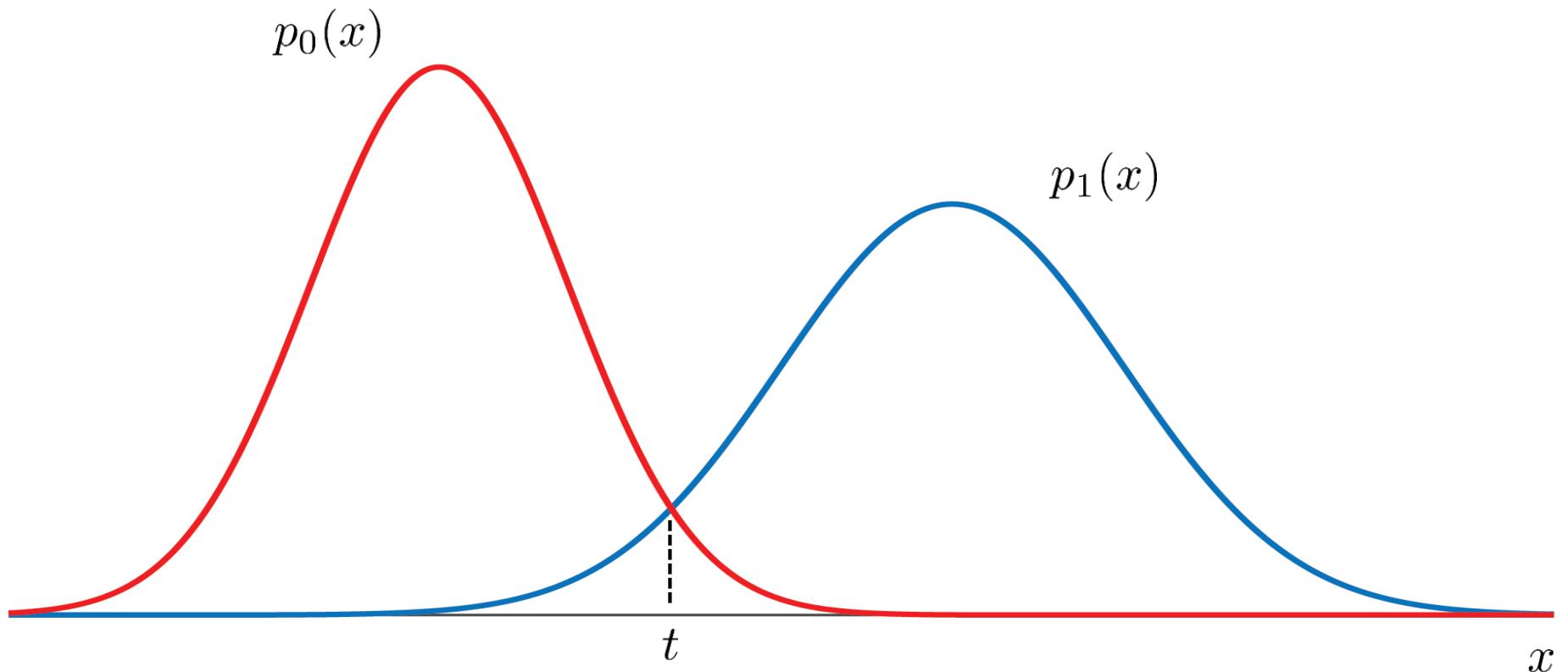
# Noise vs. uncertainty

- Distribution overlap is commonly confused with noise
  - noise implies the true value is modified



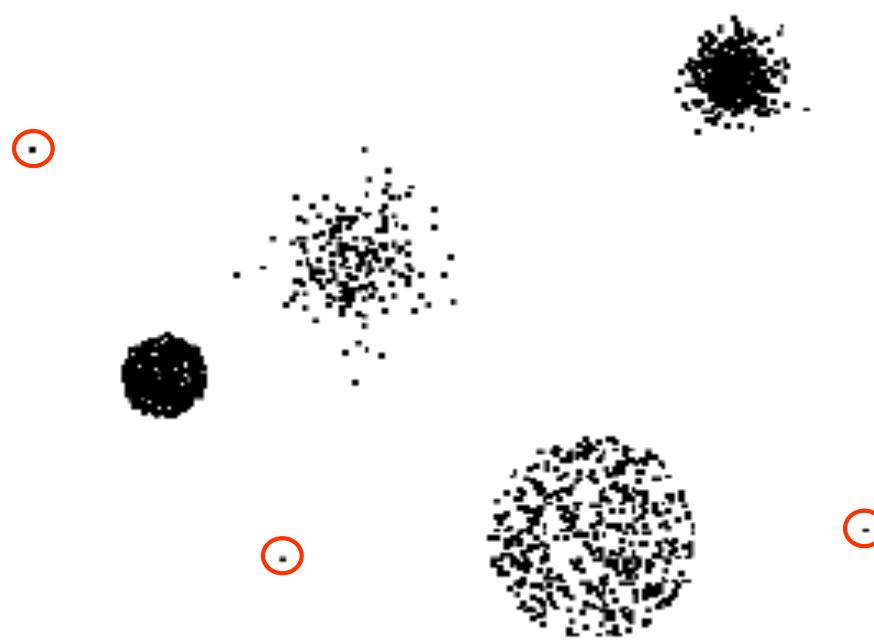
# Noise vs. uncertainty

- Distribution overlap is commonly confused with noise
  - noise implies the true value is modified



# Outliers

- Outliers are data objects with characteristics that are considerably different than most of the other data objects in the data set



# Missing Values

- Reasons for missing values
  - Information is not collected  
(e.g., people decline to give their age and weight)
  - Attributes may not be applicable to all cases  
(e.g., annual income is not applicable to children)

# Missing Values

- Reasons for missing values
  - Information is not collected  
(e.g., people decline to give their age and weight)
  - Attributes may not be applicable to all cases  
(e.g., annual income is not applicable to children)
  
- Handling missing values
  - Eliminate data objects
  - Estimate missing values
  - Ignore the missing value during analysis
  - Replace with all possible values (weighted by their probabilities)

# Duplicate Data

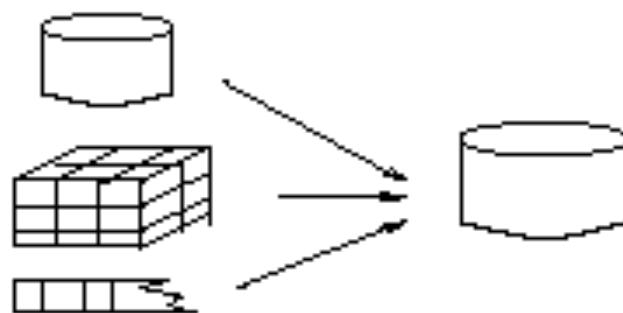
- Data set may include data objects that are duplicates, or almost duplicates of one another
  - Major issue when merging data from heterogeneous sources
- Examples:
  - Same person with multiple email addresses
- Data cleaning
  - Process of dealing with duplicate data issues

# Forms of Data Preprocessing

## Data Cleaning



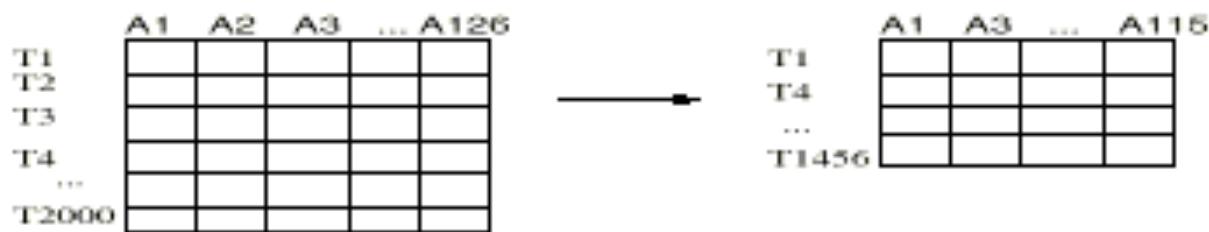
## Data Integration



## Data Transformation

-2, 32, 100, 59, 48 → -0.02, 0.32, 1.00, 0.59, 0.48

## Data Reduction



# Data Preprocessing

- Integration
- Data cleaning
- Aggregation
- Sampling
- Dimensionality reduction
- Feature subset selection
- Feature creation
- Discretization and binarization
- Data transformation

# Data Cleaning

- Importance
  - “Data cleaning is one of the three biggest problems in data warehousing”—Ralph Kimball
  - “Data cleaning is the number one problem in data warehousing”—DCI survey
- Data cleaning tasks
  - Fill in missing values
  - Identify outliers and smooth out noisy data
  - Correct inconsistent data
  - Resolve redundancy caused by data integration

# How to Handle Missing Data?

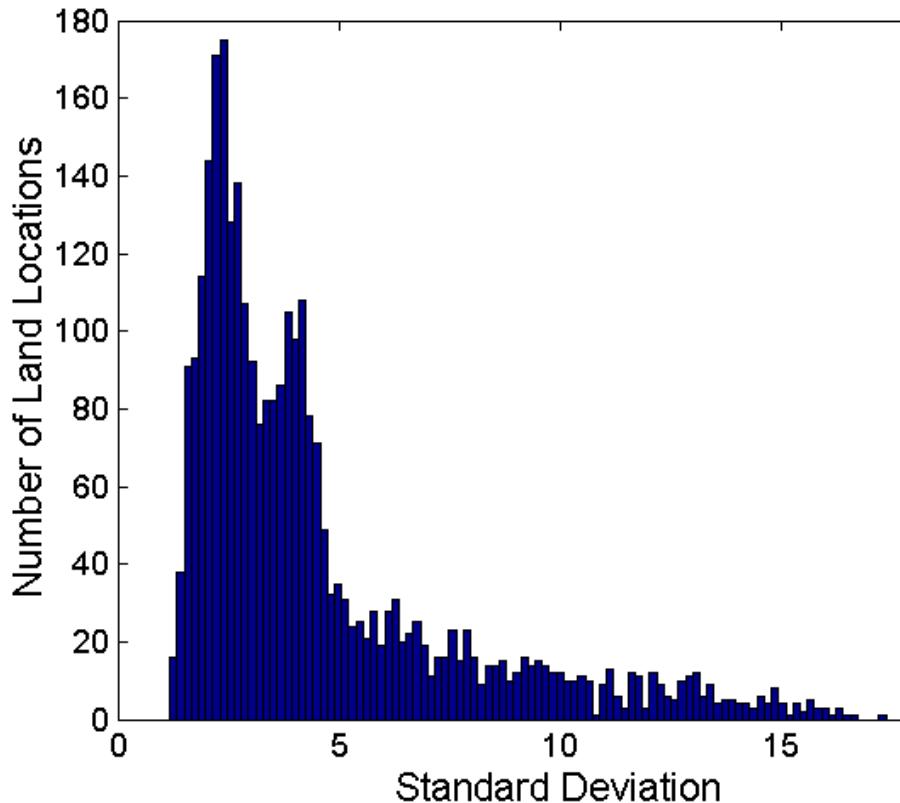
- Ignore the tuple: usually done when class label is missing (assuming the tasks in classification—not effective when the percentage of missing values per attribute varies considerably.)
- Fill in the missing value manually: tedious + infeasible?
- Fill in it automatically with
  - a global constant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

# Aggregation

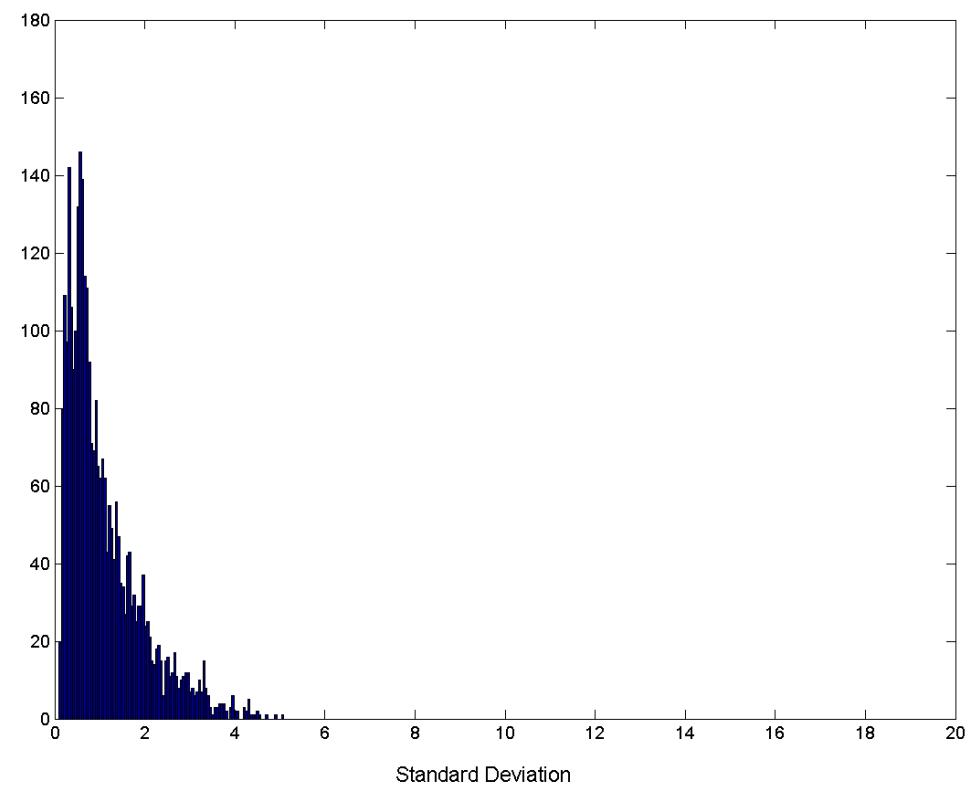
- Combining two or more attributes (or objects) into a single attribute (or object)
- Purpose
  - Data reduction
    - ◆ Reduce the number of attributes or objects
  - Change of scale
    - ◆ Cities aggregated into regions, states, countries, etc
  - More “stable” data
    - ◆ Aggregated data tends to have less variability

# Aggregation

## Variation of Precipitation in Australia



Standard Deviation of Average  
Monthly Precipitation



Standard Deviation of Average  
Yearly Precipitation

# Data Integration

- Data integration:
  - Combines data from multiple sources into a coherent store
- Schema integration: e.g., A.cust-id  $\equiv$  B.cust-#
  - Integrate metadata from different sources
- Entity identification problem:
  - Identify real world entities from multiple data sources, e.g., Bill Clinton = William Clinton
- Detecting and resolving data value conflicts
  - For the same real world entity, attribute values from different sources are different
  - 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*
- Careful integration of the data from multiple sources may help reduce/avoid redundancies and inconsistencies and improve mining speed and quality

# Sampling

- Sampling is the main technique employed for data selection.
  - It is often used for both the preliminary investigation of the data and the final data analysis.
- Statisticians sample because **obtaining** the entire set of data of interest is too expensive or time consuming.
- Sampling is used in data mining because **processing** the entire set of data of interest is too expensive or time consuming.

# Sampling ...

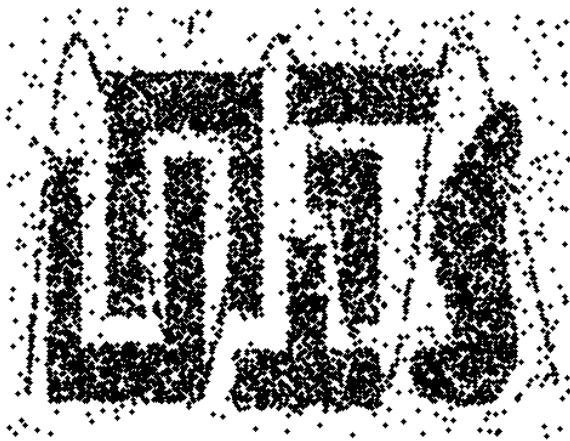
- The key principle for effective sampling is the following:
  - using a sample will work almost as well as using the entire data sets, if the sample is representative
  - a sample is representative if it has approximately the same property (of interest) as the original set of data

# Types of Sampling

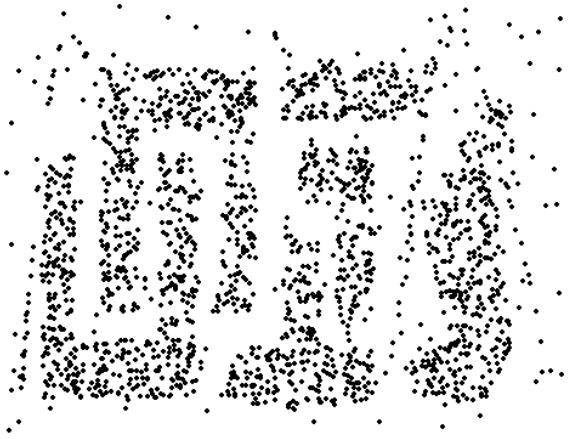
- Simple random sampling
  - There is an equal probability of selecting any particular item
- Stratified sampling
  - Split the data into several partitions; then draw random samples from each partition

---
- Sampling without replacement
  - As each item is selected, it is removed from the population
- Sampling with replacement
  - Objects are not removed from the population as they are selected for the sample. The same object can be picked up more than once.

# Sample Size



8000 points



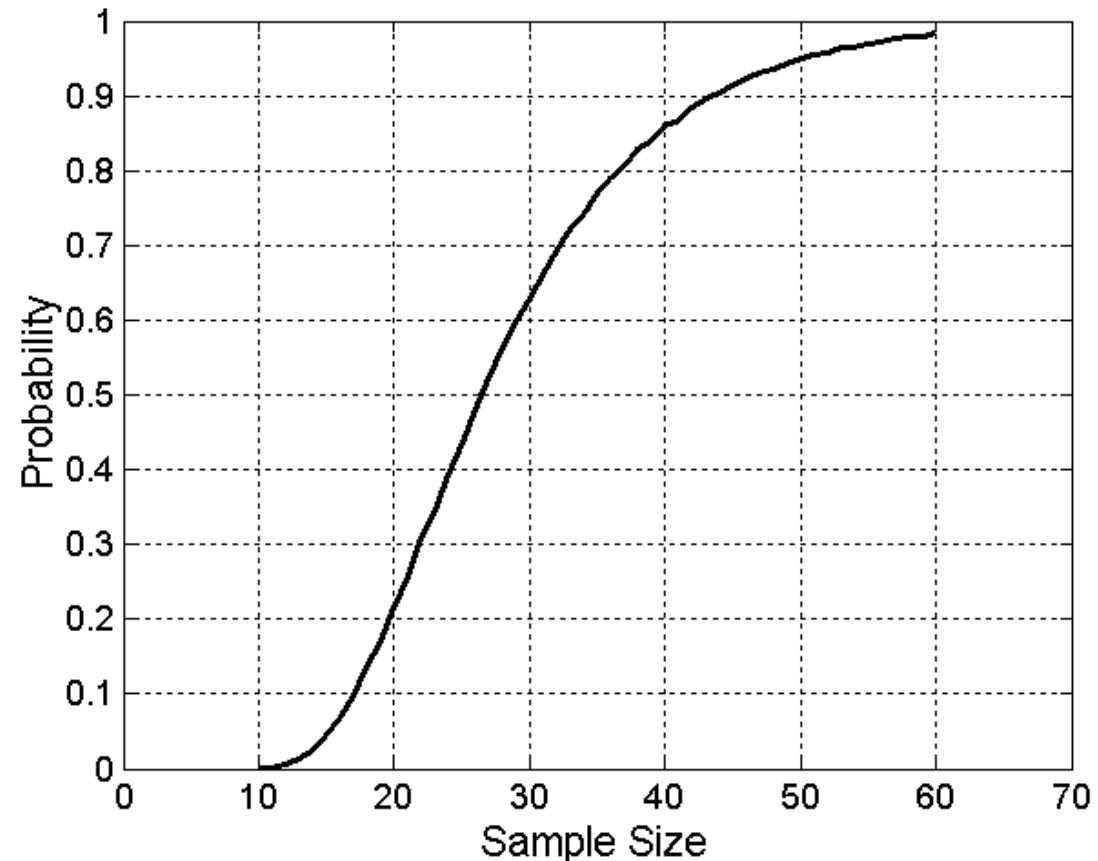
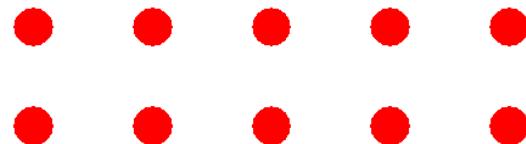
2000 Points



500 Points

# Sample Size

- What sample size is necessary to get at least one object from each of 10 groups.



# Dimensionality Reduction

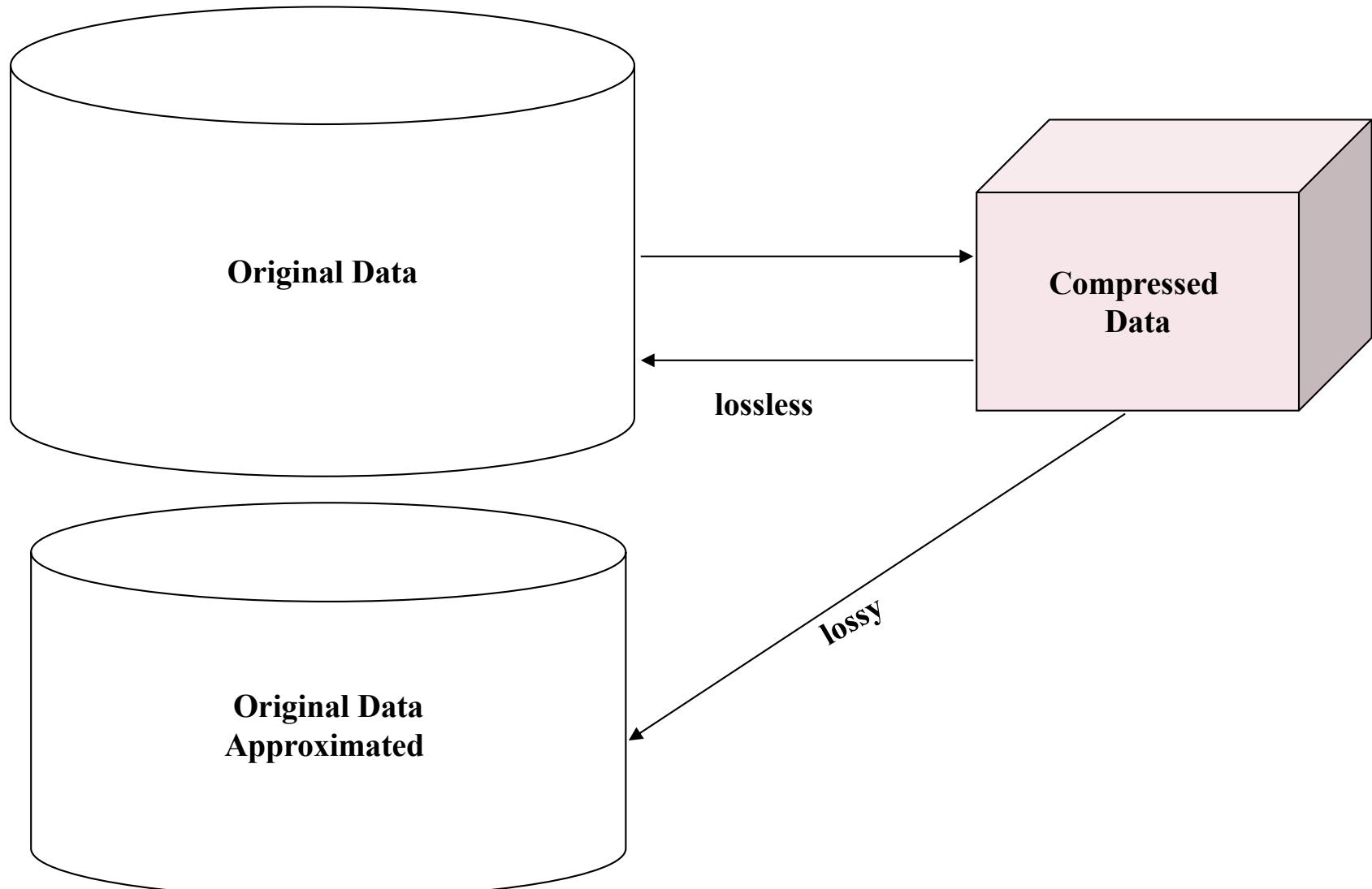
- Purpose:

- avoid curse of dimensionality
- reduce amount of time and memory required by data mining algorithms
- allow data to be more easily visualized
- may help to eliminate irrelevant features or reduce noise
- may help to avoid stability problems

- Techniques

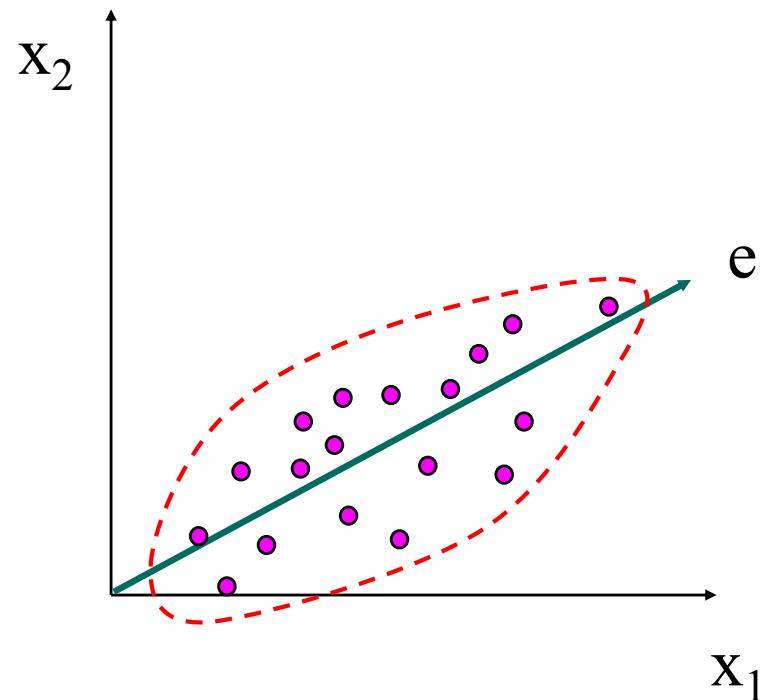
- Principal Component Analysis (PCA)
- Singular Value Decomposition (SVD)
- Others: supervised and non-linear techniques

# Relationship to Data Compression



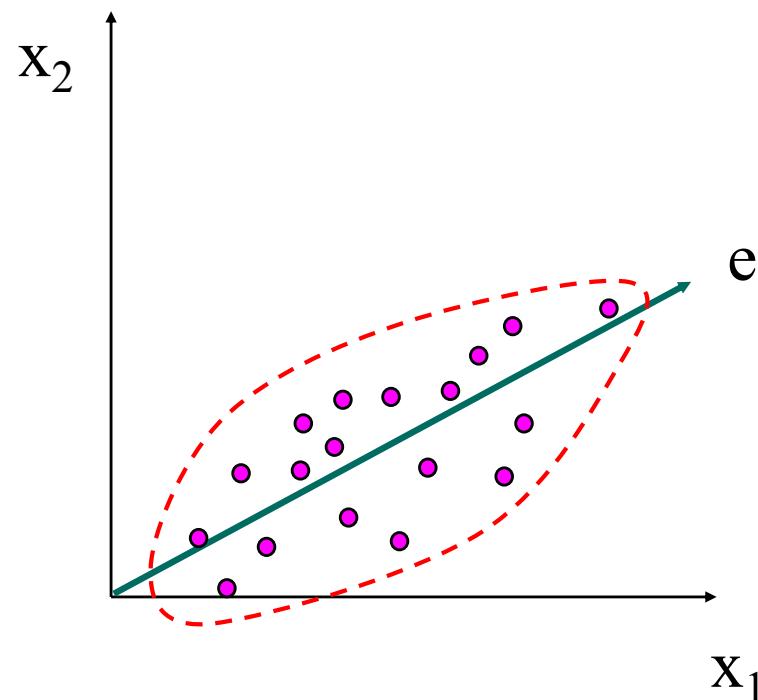
# Dimensionality Reduction: PCA

- Goal is to find a projection that captures the largest amount of variation in data



# Dimensionality Reduction: PCA

- Find the eigenvectors of the covariance matrix
- The eigenvectors define the new space



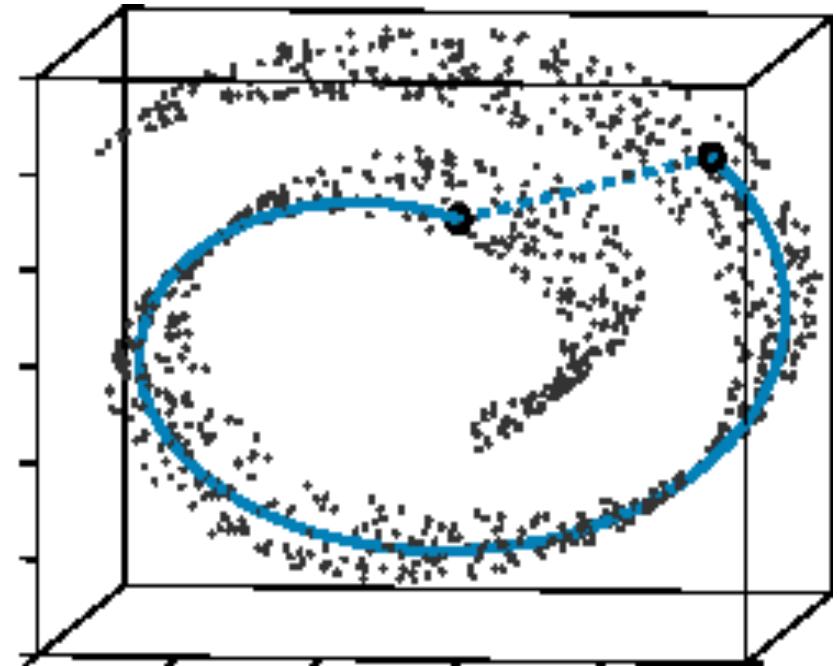
# Dimensionality Reduction: PCA

Dimensions = 206



# Dimensionality Reduction: ISOMAP

By: Tenenbaum, de Silva, Langford (2000)



- Construct a neighbourhood graph
- For each pair of points in the graph, compute the shortest path distances – geodesic distances

# Feature Subset Selection

- Another way to reduce dimensionality of data
- Redundant features
  - duplicate much or all of the information contained in one or more other attributes
  - **example:** purchase price of a product and the amount of sales tax paid
- Irrelevant features
  - contain no information that is useful for the data mining task at hand
  - **example:** students' ID is often irrelevant to the task of predicting students' GPA

# Feature Subset Selection

- **Brute-force approach:**
  - Try all possible feature subsets as input to data mining algorithm
- **Embedded approaches:**
  - Feature selection occurs naturally as part of the data mining algorithm
- **Filter approaches (usually one pass through data):**
  - Features are selected before data mining algorithm is run
- **Wrapper approaches (usually many passes through data):**
  - Use the data mining algorithm as a black box to find best subset of attributes

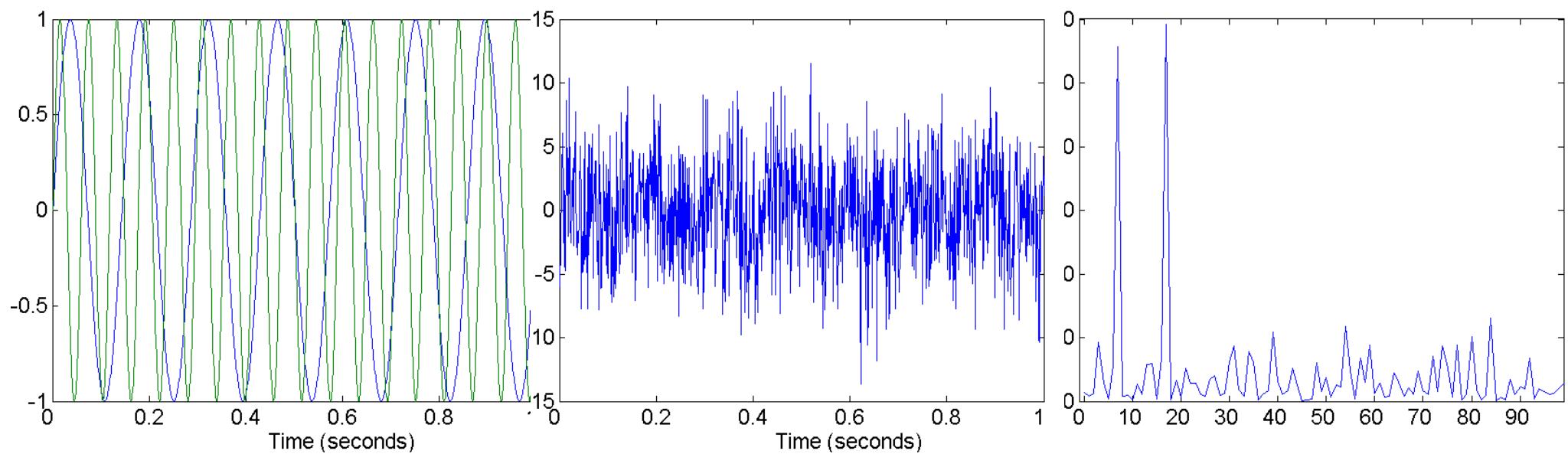
	The	Game	Play	Football	Baseball	Brady	Deflate	Gate
Document 1	12	2	3	14		4	4	6
Document 2	18	5	5		3			5
Document 3	24					4		5
Document 4	56	15					24	

# Feature Creation

- Create new attributes that can capture the important information in a data set much more efficiently than the original attributes
- Three general methodologies:
  - Feature Extraction
    - ◆ domain-specific
  - Mapping Data to New Space
  - Feature Construction
    - ◆ combining features

# Mapping Data to a New Space

- Fourier transform
- Wavelet transform

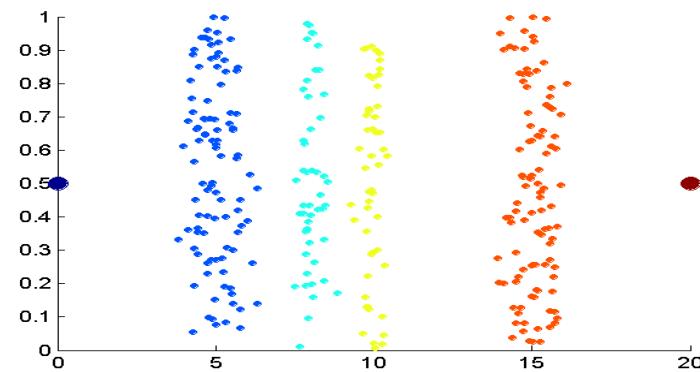


Two Sine Waves

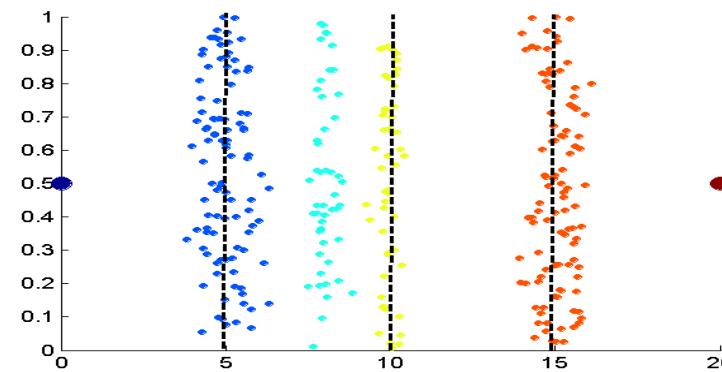
Two Sine Waves + Noise

Frequency

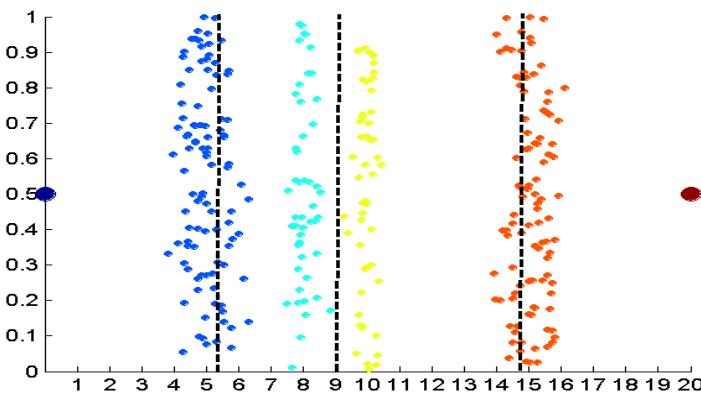
# Discretization Without Using Class Labels



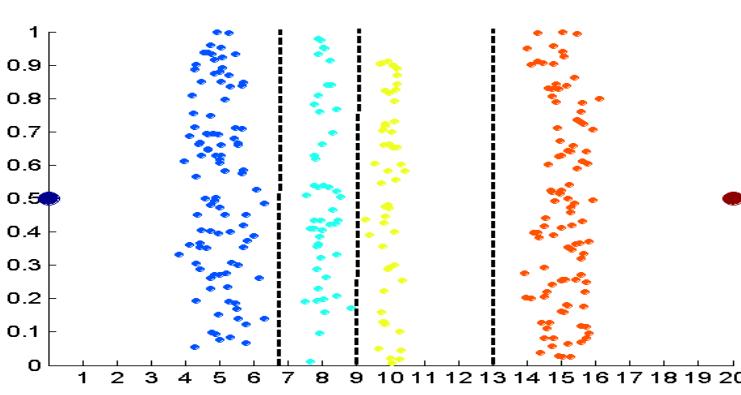
Data



Equal interval width



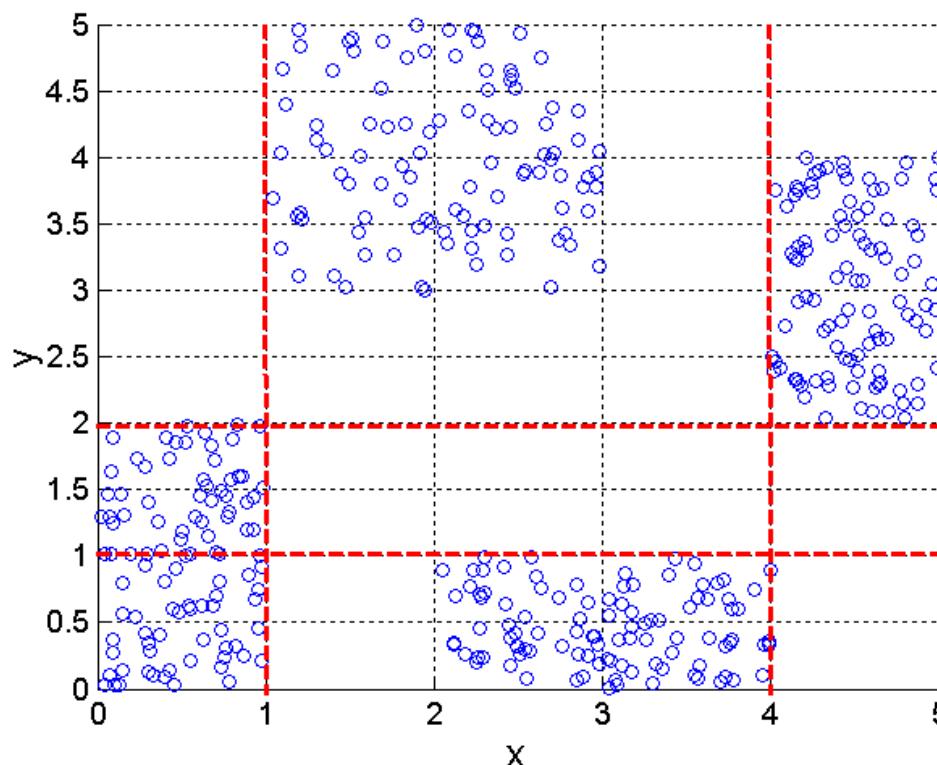
Equal frequency



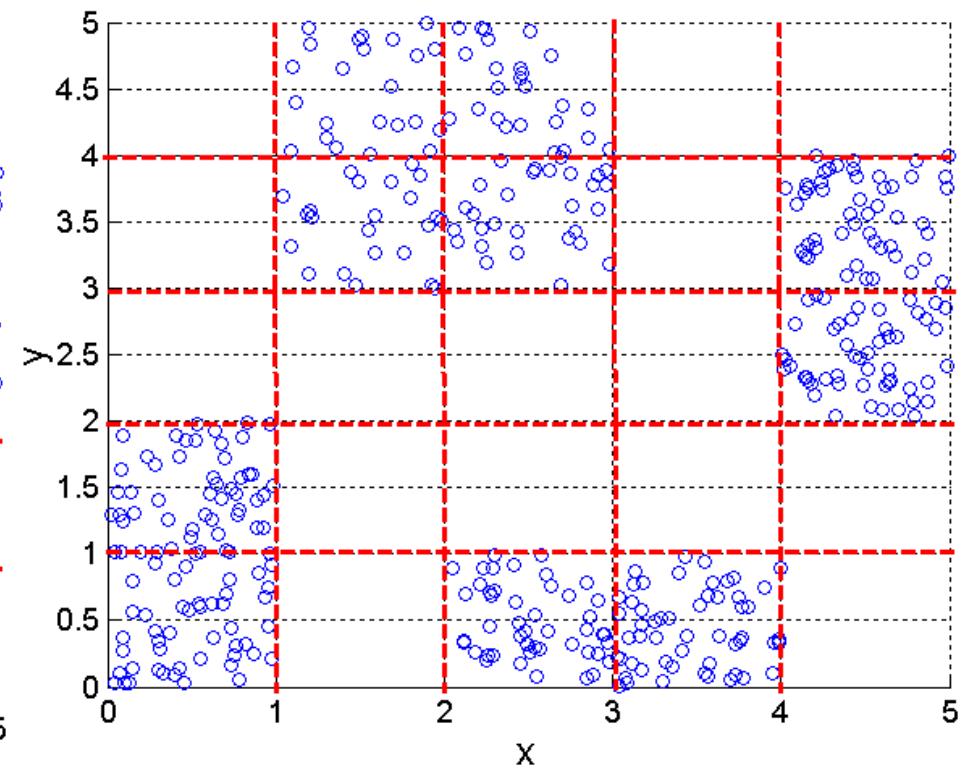
K-means

# Discretization Using Class Labels

- Entropy based approach



3 categories for both x and y



5 categories for both x and y

# How to Handle Noisy Data?

- Binning
  - first sort data and partition into (equal-frequency) bins
  - then one can **smooth by bin means, smooth by bin median, smooth by bin boundaries**, etc.
- Regression
  - smooth by fitting the data into regression functions
- Clustering
  - detect and remove outliers
- Combined computer and human inspection
  - detect suspicious values and check by human (e.g., deal with possible outliers)

# Simple Discretization Methods: Binning

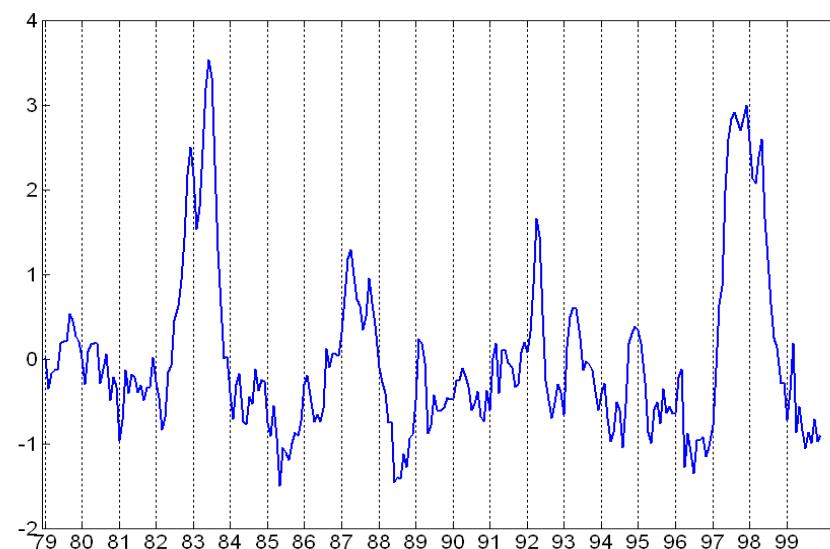
- **Equal-width** (distance) partitioning
  - Divides the range into  $N$  intervals of equal size: uniform grid
  - if  $A$  and  $B$  are the lowest and highest values of the attribute, the width of intervals will be:  $W = (B - A)/N$ .
  - The most straightforward, but outliers may dominate presentation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

# Binning Methods for Data Smoothing

- 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

# Attribute Transformation

- A function that maps the entire set of values of a given attribute to a new set of replacement values such that each old value can be identified with one of the new values
  - Simple functions:  $x^k$ ,  $\log(x)$ ,  $e^x$ ,  $|x|$
  - Standardization and normalization



# Data Transformation: Normalization

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

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

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

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

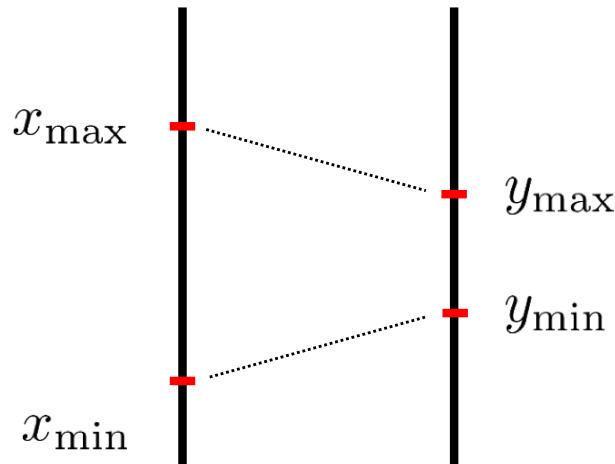
$$v' = \frac{v - \mu_A}{\sigma_A}$$

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

- Normalization by decimal scaling

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

# Deriving the Min-Max Normalization



**Find linear transform:**

$$y = ax + b$$

$$x_{\min} \rightarrow y_{\min}$$

$$x_{\max} \rightarrow y_{\max}$$

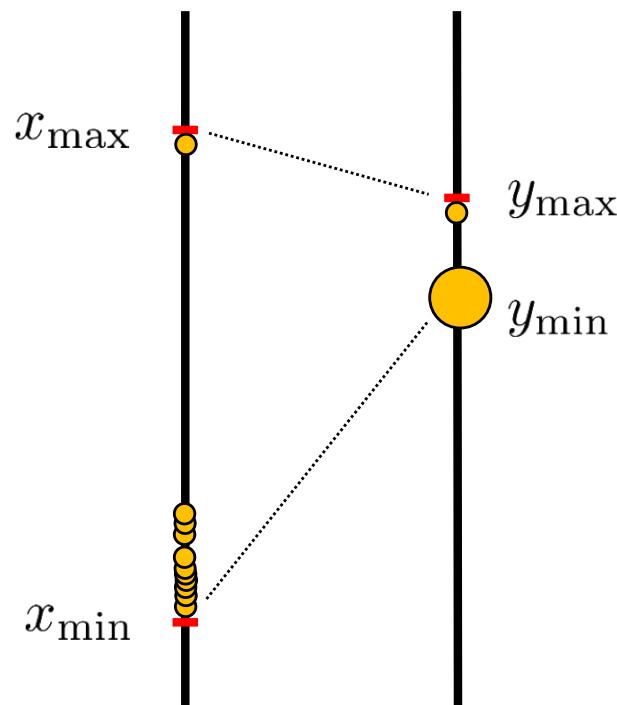
Equation of a line given two points;

i.e.,  $(x_1, y_1)$  and  $(x_2, y_2)$

$$y - y_1 = \frac{y_2 - y_1}{x_2 - x_1} (x - x_1)$$

You take it from here...

# Min-max Normalization: Problems?

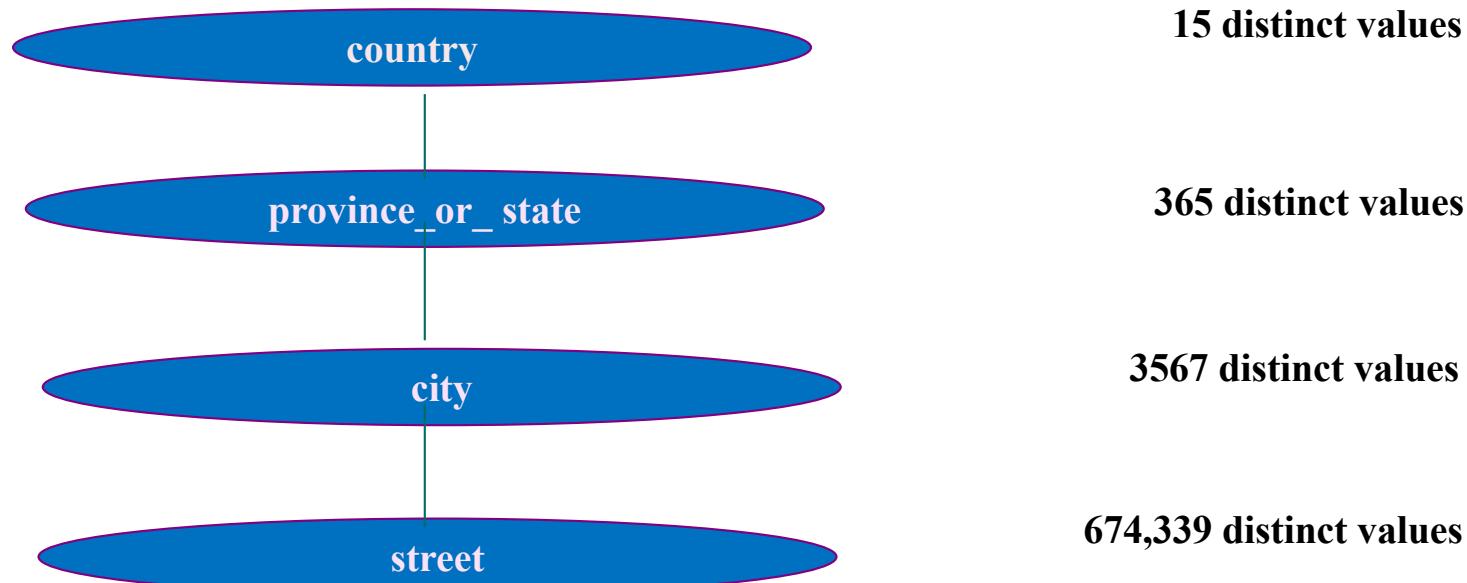


# Concept Hierarchy Generation for Categorical Data

- 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
  - {Bloomington, Indianapolis, South Bend} < Indiana
- 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: {street, city, state, country}

# Automatic Concept Hierarchy Generation

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



# Summary

- Data preparation or preprocessing is a big issue for both data warehousing and data mining
- Descriptive data summarization is need for quality data preprocessing
- Data preparation includes
  - Data cleaning and data integration
  - Data reduction and feature selection
  - Discretization
- A lot a methods have been developed but data preprocessing still an active area of research