



# Data Mining: Data (Chapter 2)

郑子彬 副教授  
中山大学 数据科学与计算机学院  
[zhzibin@mail.sysu.edu.cn](mailto:zhzibin@mail.sysu.edu.cn)  
2016年



# Project: Crime Classification



- 组队要求：3个小组成员，如有特殊情况（如不够人）可向TA申请
- 组队信息3月12日前发送到邮箱：[dm2016sysu@sina.com](mailto:dm2016sysu@sina.com)
- TA统计好组队信息后将邮件通知各小组组号，Kaggle注册账号请使用DM\_组号（如DM\_001），评分将根据leaderboard上指定的账号名的分数排名

# Example



You receive an email from a medical researcher concerning a project that you are eager to work on.

Hi,

I've attached the data file.

Each line contains the information for a single patient and consists of five fields.

We want to predict the last field using the other fields.

Thanks and see you in a couple of days.

# Example



The first few rows of the file are as follows:

012	232	33.5	0	10.7
020	121	16.9	2	210.1
027	165	24.0	0	427.6
:				

**Nothing looks strange.  
You put your doubts aside  
and start the analysis.**

**Two days later you arrive for the meeting, and before the meeting, you strike up a conversation with a statistician who is working on the project.**

# Example



**Statistician:** So, you got the data for all the patients?

**Data Miner:** Yes. I haven't had much time for analysis, but I do have a few interesting results.

**Statistician:** Amazing. There were so many data issues with this set of patients that I couldn't do much.

**Data Miner:** Oh? I didn't hear about any possible problems.

**Statistician:** But surely you heard about what happened to field 4? It's supposed to be measured on a scale from 1 to 10, with 0 indicating a missing value, but because of a data entry error, all 10's were changed into 0's.

**Data Miner:** Interesting. Were there any other problems?

**Statistician:** Yes, fields 2 and 3 are basically the same, but I assume that you probably noticed that.

**Data Miner:** Yes, but these fields were only weak predictors of field 5.

012	232	33.5	0	10.7
020	121	16.9	2	210.1
027	165	24.0	0	427.6
⋮				

# Example



**Statistician:** Anyway, given all those problems, I'm surprised you were able to accomplish anything.

**Data Miner:** True, but my results are really quite good. Field 1 is a very strong predictor of field 5. I'm surprised that this wasn't noticed before.

**Statistician:** What? Field 1 is just an identification number.

**Data Miner:** Nonetheless, my results speak for themselves.

**Statistician:** Oh, no! I just remembered. We assigned ID numbers after we sorted the records based on field 5. There is a strong connection, but it's meaningless. Sorry.

**Lesson: Get to know your data!**

012	232	33.5	0	10.7
020	121	16.9	2	210.1
027	165	24.0	0	427.6

# What is Data?



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

## Attributes

## Objects

<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

# Attribute Values



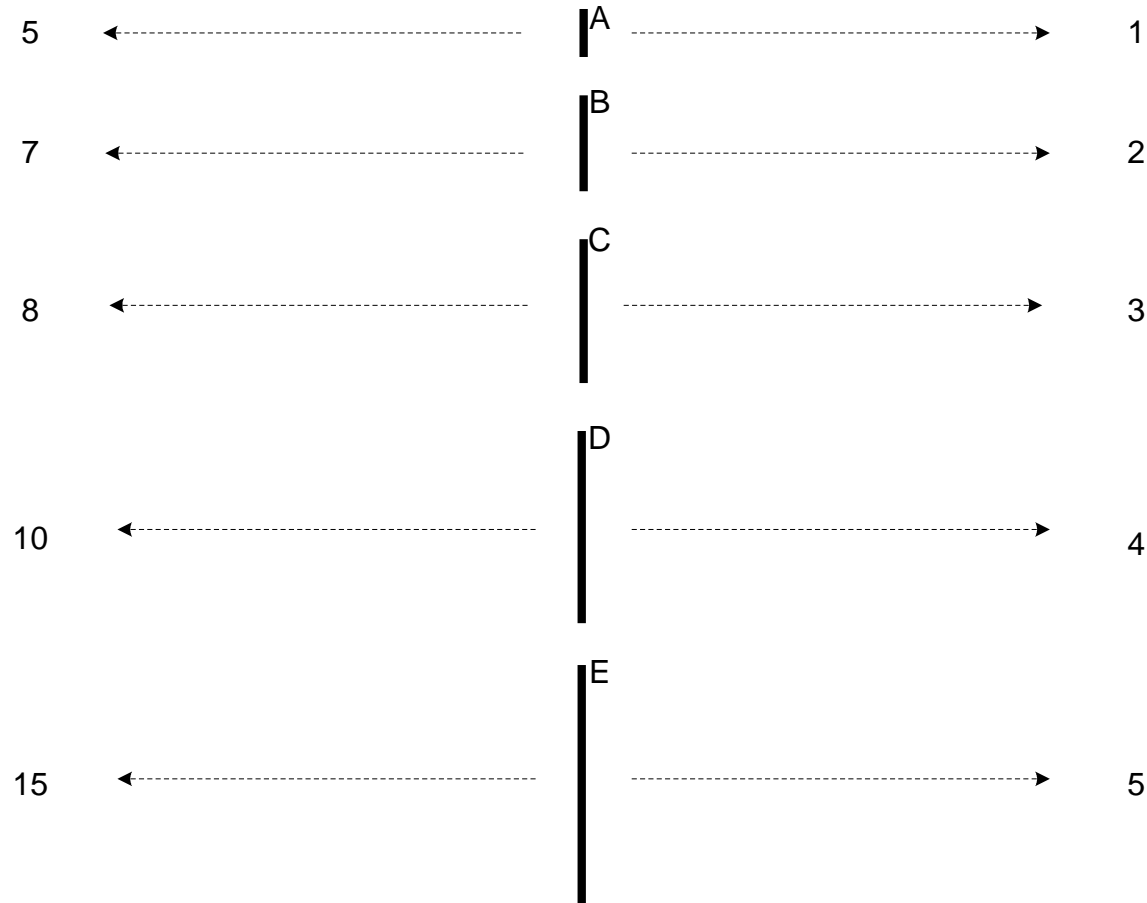
- Attribute values are numbers or symbols assigned to an attribute
- Distinction between attributes and attribute values
  - Same attribute can be mapped to different attribute values
    - ◆ Example: height can be measured in feet or meters
  - Different attributes can be mapped to the same set of values
    - ◆ 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 somewhat 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: { <i>male</i> , <i>female</i> }	mode, entropy, contingency correlation, $\chi^2$ test
Ordinal	The values of an ordinal attribute provide enough information to order objects. ( $<$ , $>$ )	hardness of minerals, { <i>good</i> , <i>better</i> , <i>best</i> }, 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 one-to-one mapping	If all employee ID numbers were reassigned, would it make any difference?
Ordinal	An order preserving change of values, i.e., $new\_value = f(old\_value)$ where $f$ is a monotonic function.	An attribute encompassing the notion of good, better best can be represented equally well by the values {1, 2, 3} or by {0.5, 1, 10}.
Interval	$new\_value = a * old\_value + b$ where a and b are constants	Thus, the Fahrenheit and Celsius temperature scales differ in terms of where their zero value is and the size of a unit (degree).
Ratio	$new\_value = a * old\_value$	Length can be measured in meters or feet.

# Discrete 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.

# Important Characteristics of Structured Data



- **Dimensionality**
  - ◆ **Curse of Dimensionality**
  
- **Sparsity**
  - ◆ **Only presence counts**
  
- **Resolution**
  - ◆ **Patterns depend on the scale**

# 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



# 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	player	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

# Transaction Data



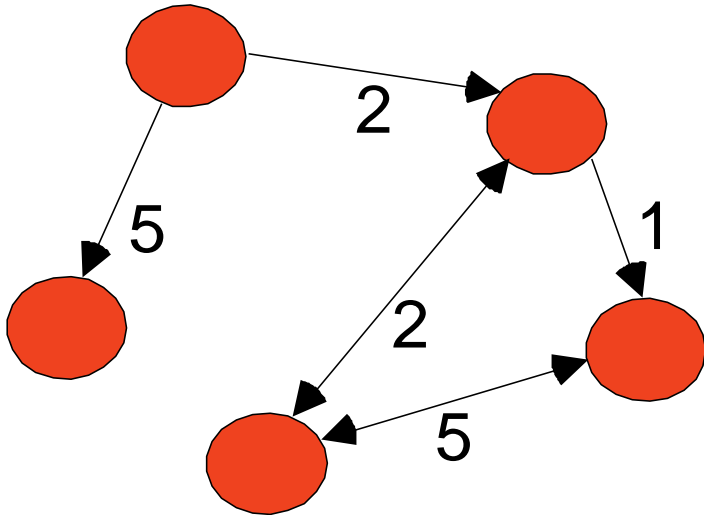
- A special type of record data, where
  - each record (transaction) involves a set of items.
  - For example, consider a grocery store. The set of products purchased by a customer during one shopping trip constitute a transaction, while the individual products that were purchased are the items.

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

# Graph Data

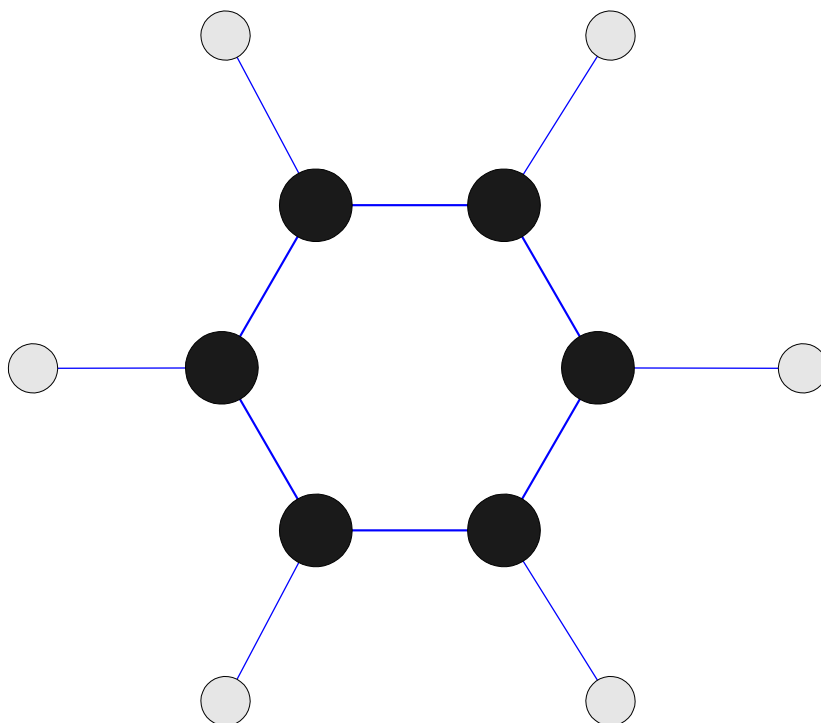


- Examples: Generic graph and HTML Links



```
<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
```

- Benzene Molecule (苯分子) :  $\text{C}_6\text{H}_6$



# Ordered Data: Sequential Data



- Sequential Data

Time	Customer	Items Purchased
t1	C1	A, B
t2	C3	A, C
t2	C1	C, D
t3	C2	A, D
t4	C2	E
t5	C1	A, E

Customer	Time and Items Purchased
C1	(t1: A,B) (t2:C,D) (t5:A,E)
C2	(t3: A, D) (t4: E)
C3	(t2: A, C)

# Ordered Data: Sequence Data



- Genomic sequence data

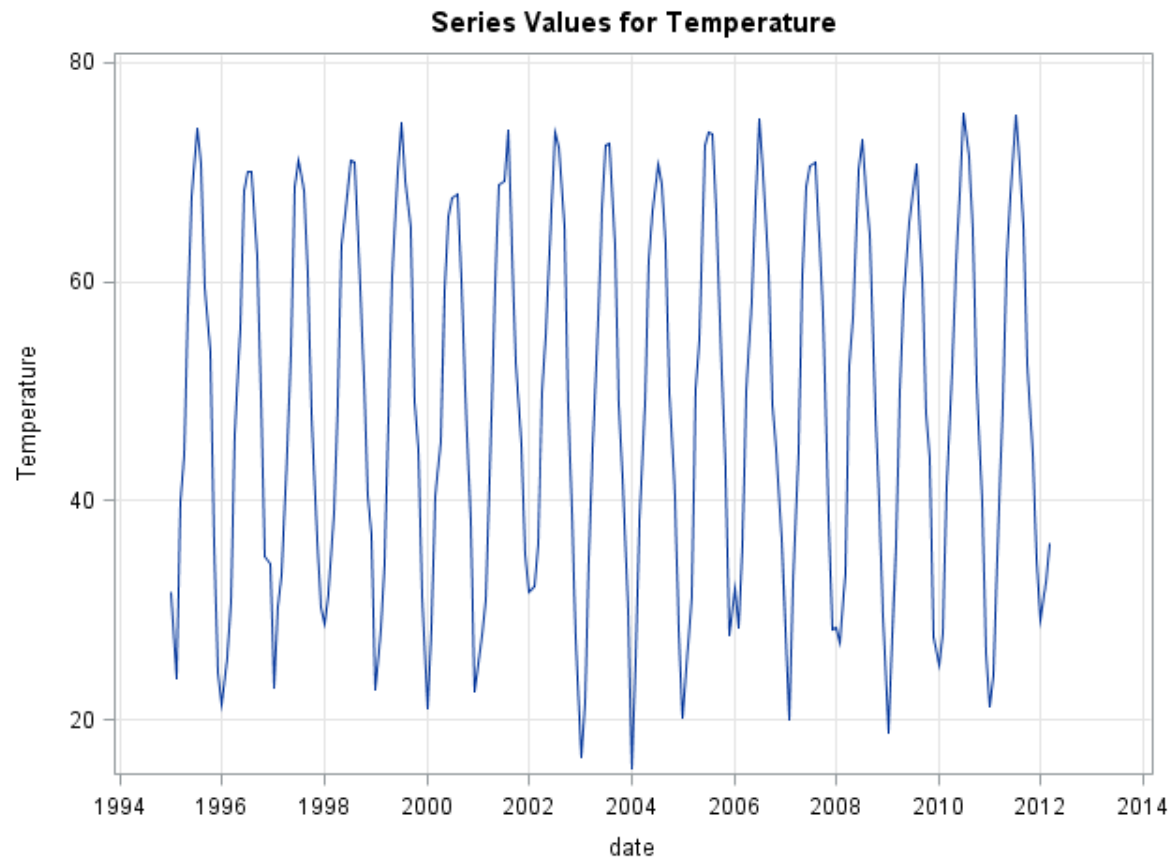
GGTTCCGCCTTCAGCCCCGCGCC  
CGCAGGGCCCGCCCCGCGCCGTC  
GAGAAGGGCCCGCCTGGCGGGCG  
GGGGGAGGCGGGGCCGCCCGAGC  
CCAACCGAGTCCGACCAGGTGCC  
CCCTCTGCTCGGCCTAGACCTGA  
GCTCATTAGGCGGCAGCGGACAG  
GCCAAGTAGAACACGCGAAGCGC  
TGGGCTGCCTGCTGCGACCAGGG



# Ordered Data: Time Series Data



- Special type of sequential data
- Temporal autocorrelation

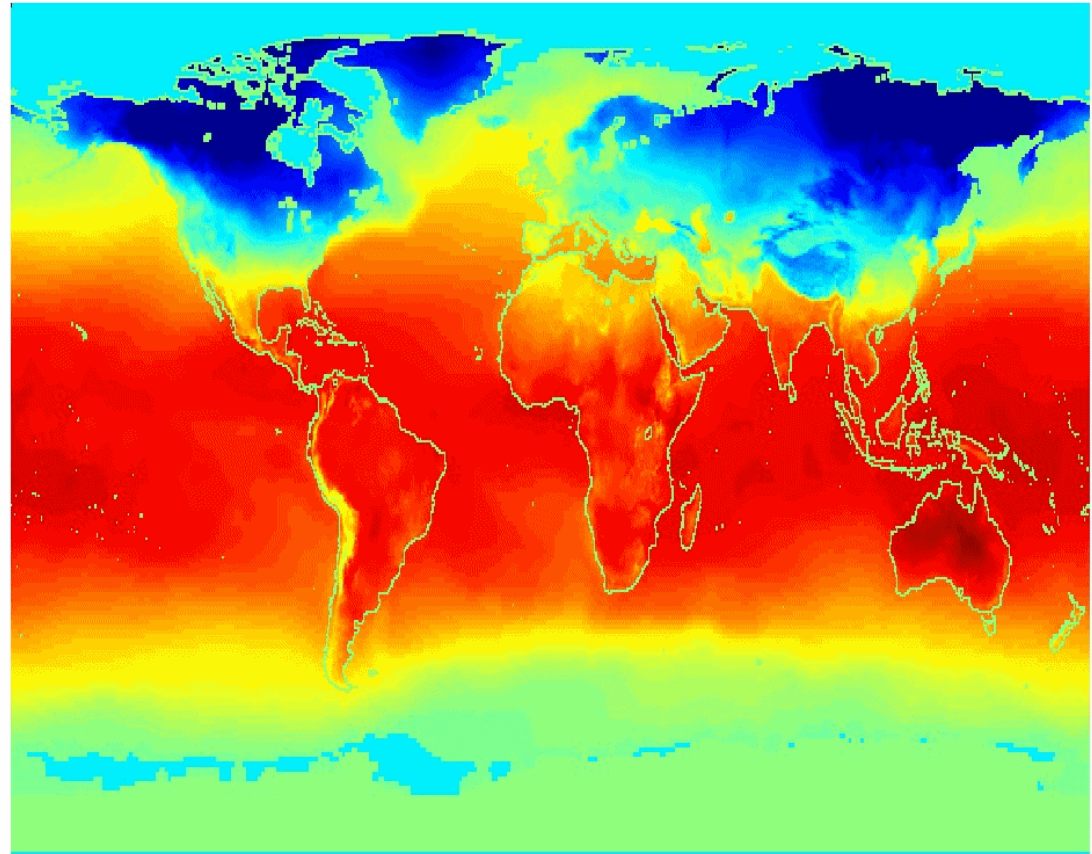


# Ordered Data: Spatio-Temporal Data



Jan

**Average Monthly  
Temperature of  
land and ocean**

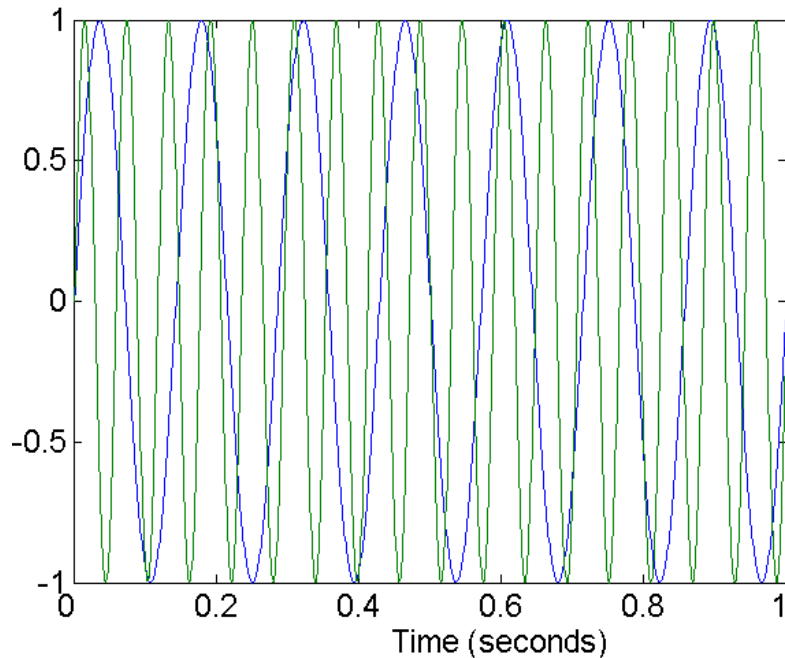


- What kinds of data quality problems?
- How can we detect problems with the data?
- What can we do about these problems?
  
- Examples of data quality problems:
  - Noise and outliers
  - missing values
  - duplicate data

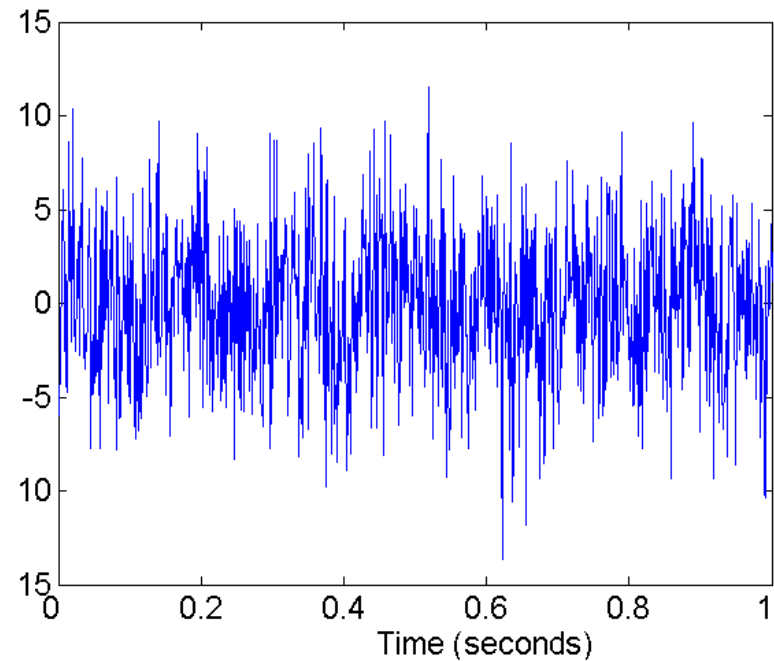
# Noise



- 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**

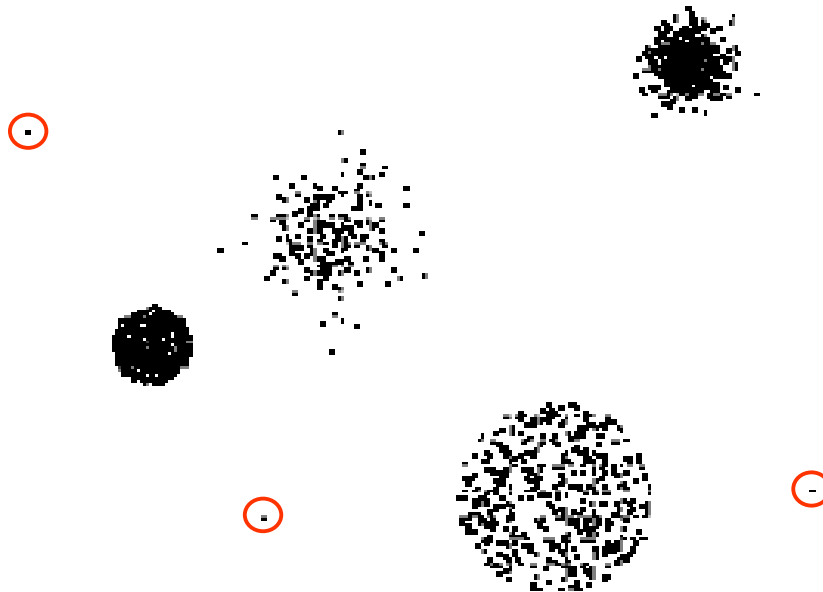


**Two Sine Waves + Noise**

# 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)
- 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

# Data Preprocessing



- Aggregation
- Sampling
- Dimensionality Reduction
- Feature subset selection
- Feature creation
- Discretization and Binarization
- Attribute Transformation



# Aggregation

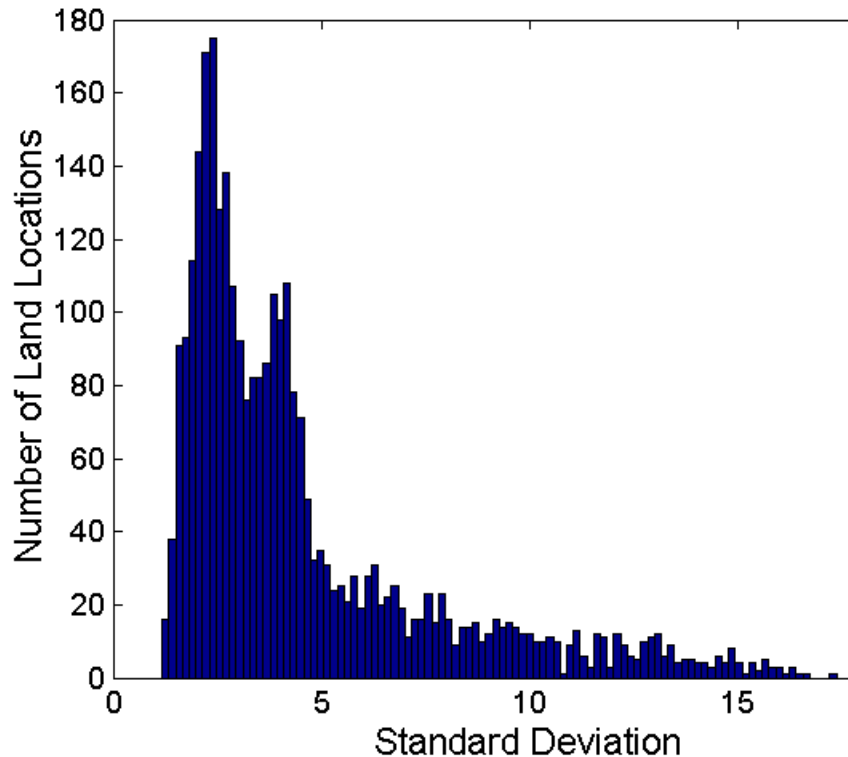


- Combining two or more attributes (or objects) into a single attribute (or object)
- Purpose
  - Data reduction
    - ◆ Reduce the number of attributes or objects
  - Change of scale
    - ◆ Cities aggregated into regions, states, countries, etc
  - More “stable” data
    - ◆ Aggregated data tends to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

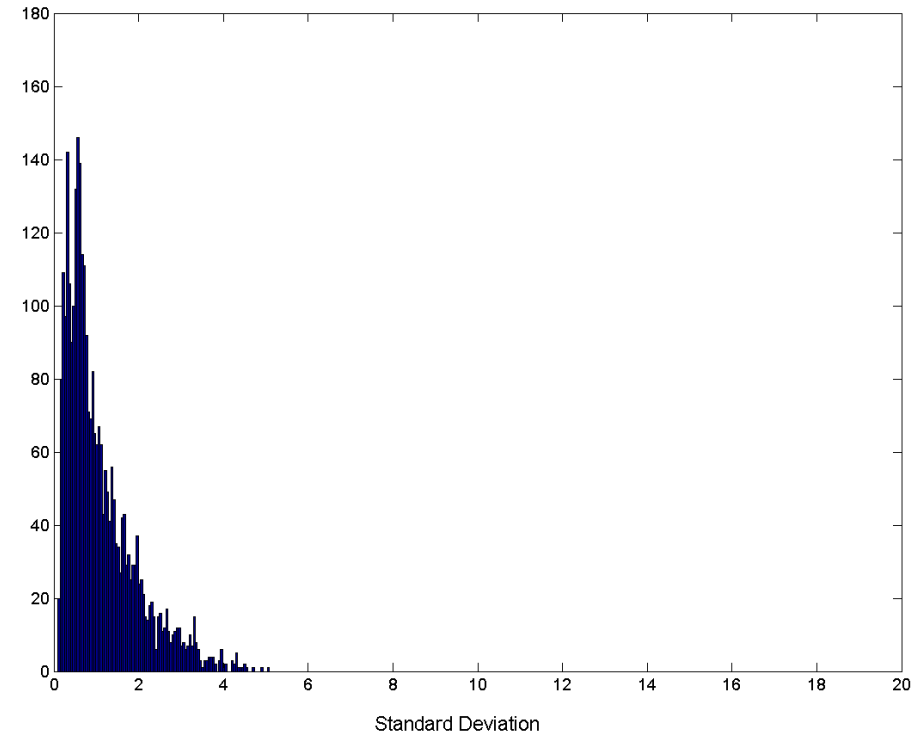
# Aggregation



## Variation of Precipitation in Australia



**Standard Deviation of Average  
Monthly Precipitation**



**Standard Deviation of Average  
Yearly Precipitation**

# 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 ...



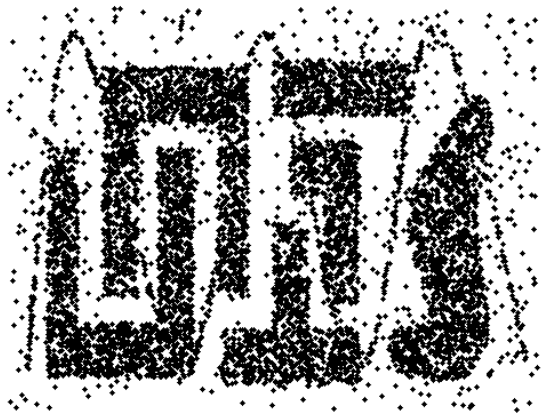
- The key principle for effective sampling is the following:
  - using a sample will work almost as well as using the entire data sets, if the sample is representative
  - A sample is representative if it has approximately the same property (of interest) as the original set of data

# Types of Sampling

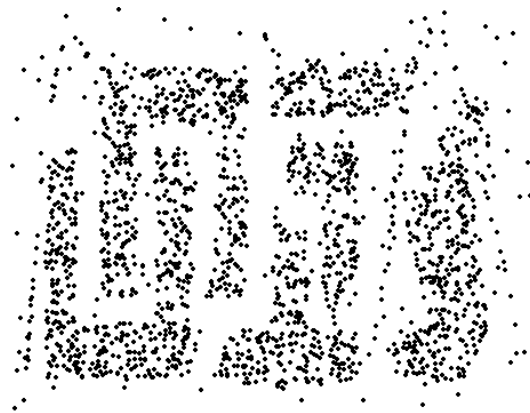


- Simple Random Sampling
  - There is an equal probability of selecting any particular item
- Sampling without replacement
  - As each item is selected, it is removed from the population
- Sampling with replacement
  - Objects are not removed from the population as they are selected for the sample.
    - ◆ In sampling with replacement, the same object can be picked up more than once
- Stratified sampling
  - Split the data into several partitions; then draw random samples from each partition

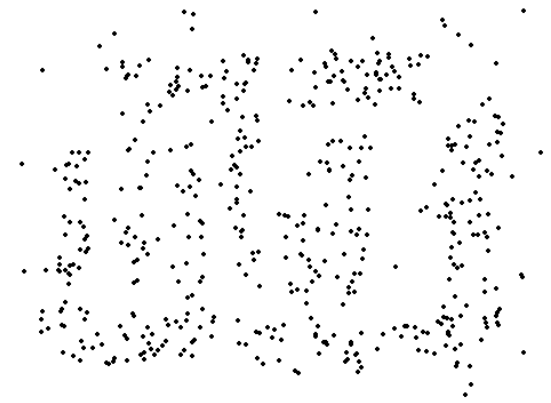
# Sample Size



8000 points



2000 Points

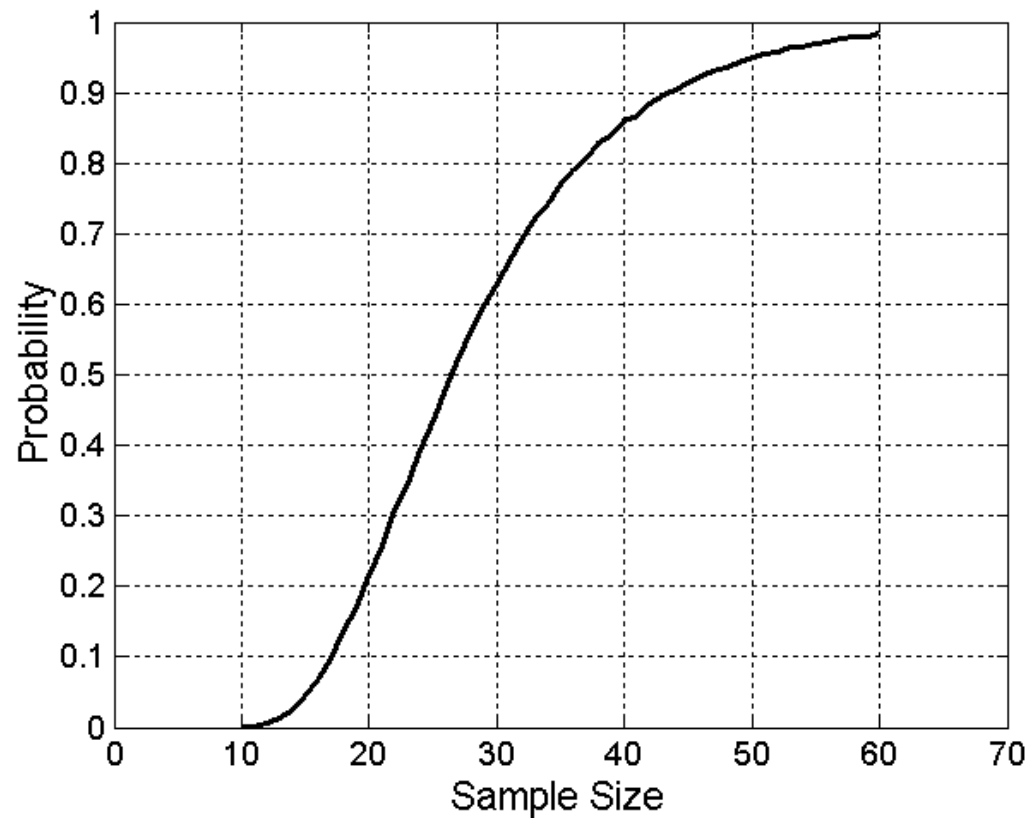


500 Points

# Sample Size



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



# Progressive Sampling



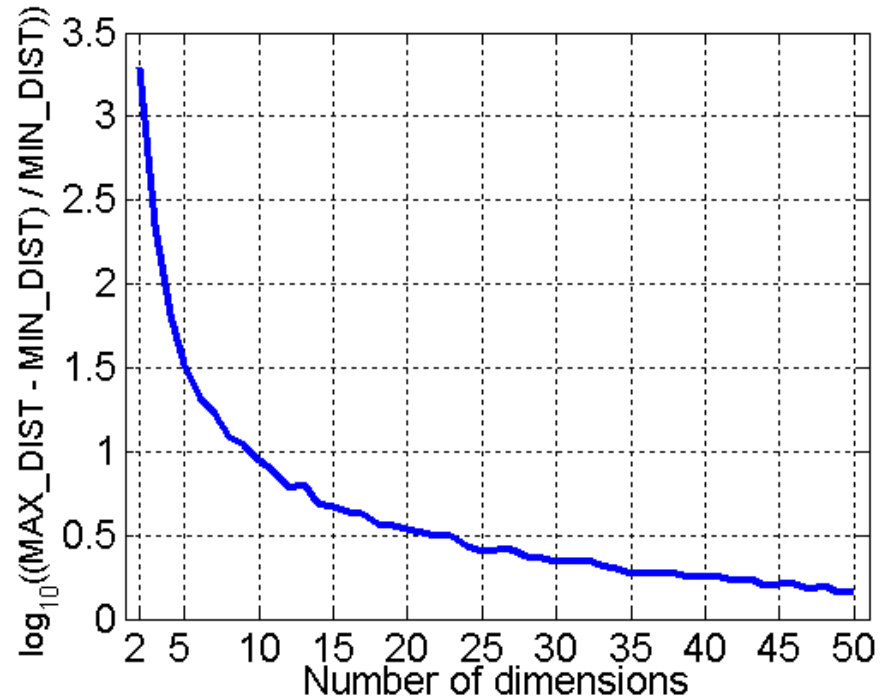
- Start with a small sample
- Increase the sample size
- Need to evaluate the sample to judge if it is large enough
- Marginal effect (边际效应)



# Curse of Dimensionality



- When dimensionality increases, data becomes increasingly sparse in the space that it occupies
- Definitions of density and distance between points, which is critical for clustering and outlier detection, become less meaningful



- Randomly generate 500 points
- Compute difference between max and min distance between any pair of points

# 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
  
- Techniques
  - Principle Component Analysis
  - Singular Value Decomposition
  - Others: supervised and non-linear techniques

# Dimensionality Reduction: PCA



- 变量之间是有一定的相关关系的
- 当两个变量之间有一定相关关系时，可以解释为这两个变量的信息有一定的重叠
- 主成分分析是对于原先提出的所有变量，将重复的变量（关系紧密的变量）删去多余，建立尽可能少的新变量，使得这些新变量是两两不相关的
- 这些新变量在反映信息方面尽可能保持原有的信息

# 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



- Techniques:
  - 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:
    - ◆ Features are selected before data mining algorithm is run
  - Wrapper approaches:
    - ◆ Use the data mining algorithm as a black box to find best subset of attributes

# Feature Creation

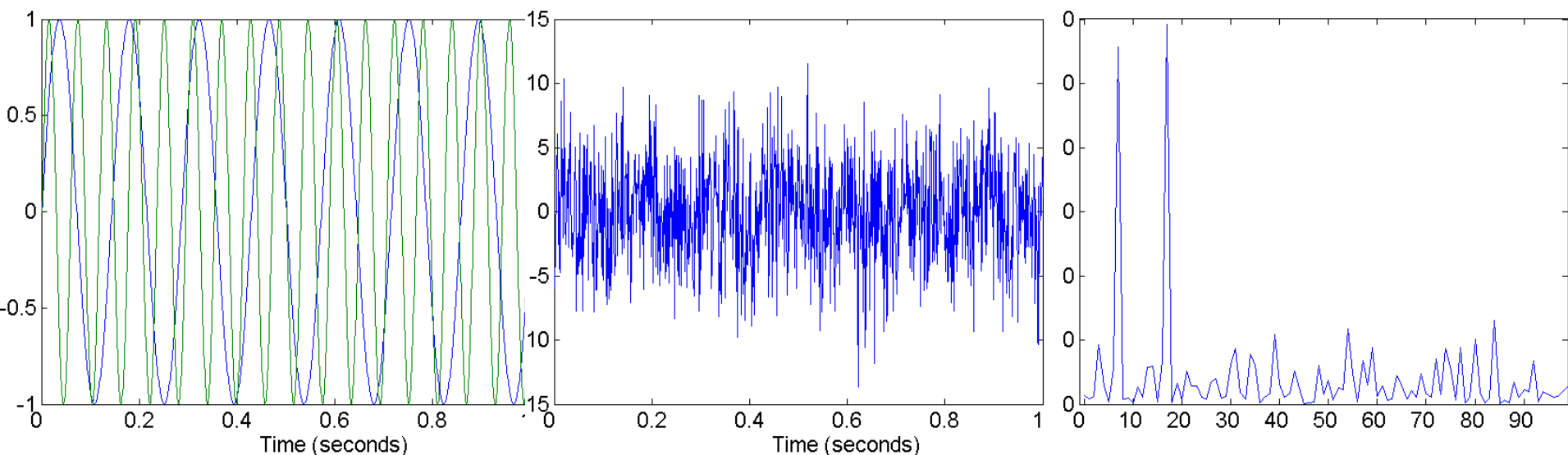


- Create new attributes that can capture the important information in a data set much more efficiently than the original attributes
- Three general methodologies:
  - Feature Extraction
    - ◆ domain-specific
  - Mapping Data to New Space
  - Feature Construction
    - ◆ combining features

# Mapping Data to a New Space



- Fourier transform
- Wavelet transform



**Two Sine Waves**

**Two Sine Waves + Noise**

**Frequency**

# Discretization and Binarization



- Discretization: Transform a continuous attribute to categorical attribute
- Binarization: Transform continuous (or discrete) attributes into one or more binary attributes

表 2-5 一个分类属性到三个二元属性的变换

分类值	整数值	$x_1$	$x_2$	$x_3$
<i>awful</i>	0	0	0	0
<i>poor</i>	1	0	0	1
<i>OK</i>	2	0	1	0
<i>good</i>	3	0	1	1
<i>great</i>	4	1	0	0



# Discretization and Binarization



表 2-6 一个分类属性到五个非对称二元属性的转换

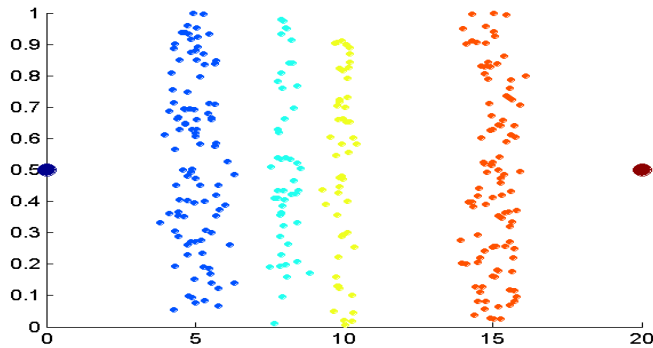
分类值	整数值	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$
<i>awful</i>	0	1	0	0	0	0
<i>poor</i>	1	0	1	0	0	0
<i>OK</i>	2	0	0	1	0	0
<i>good</i>	3	0	0	0	1	0
<i>great</i>	4	0	0	0	0	1

# Discretization

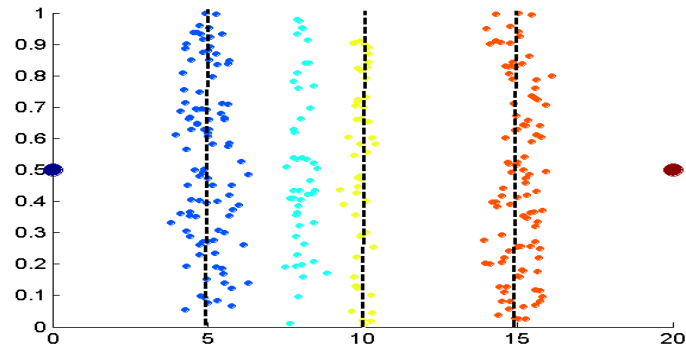


- The best discretization depends on the algorithm being used
- How many categories?
- How to map the values of continuous attributes to these categories?
- How many split points to choose and where to place them?
- Solutions
  - Unsupervised discretization
  - Supervised discretization

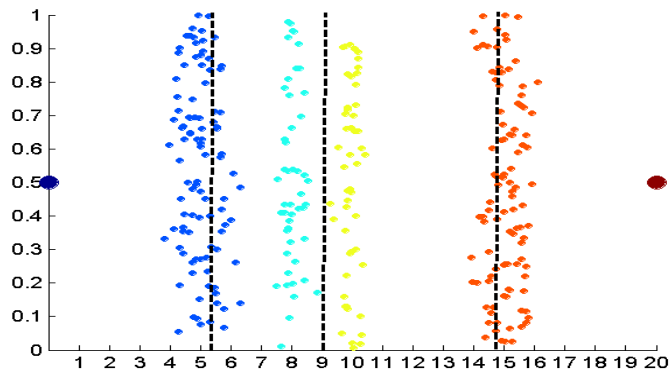
# Discretization Without Using Class Labels



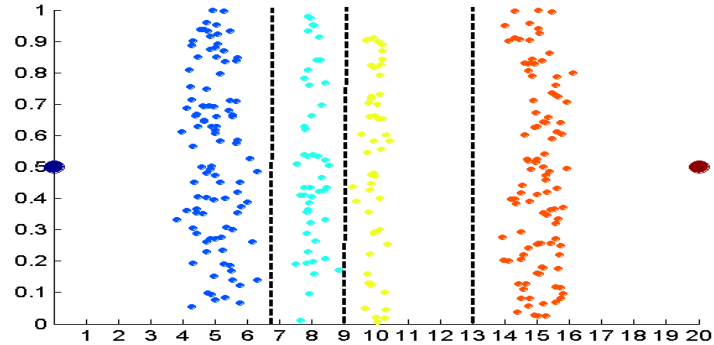
Data



Equal interval width



Equal frequency



K-means

- 基于熵的离散化方法
  - 最大化区间的纯度

$$e_i = -\sum_{j=1}^k p_{ij} \log_2 p_{ij}$$

首先，需要定义熵（entropy）。设  $k$  是不同的类标号数， $m_i$  是某划分的第  $i$  个区间中值的个数，而  $m_{ij}$  是区间  $i$  中类  $j$  的值的个数。第  $i$  个区间的熵  $e_i$  由如下等式给出

$p_{ij} = m_{ij}/m_i$  是第  $i$  个区间中类  $j$  的概率（值的比例）。

# Supervised Discretization: Entropy (熵)



## ● Entropy (熵)

- 熵的概念是由德国物理学家克劳修斯于1865年所提出。熵最初是被用在热力学方面的
- 香农1948年的一篇论文 [《A Mathematical Theory of Communication》](#) 提出了**信息熵**的概念，解决了对信息的量化度量问题，并且以后信息论也被作为一门单独的学科
- 要搞清楚一件非常非常不确定的事，就需要了解大量的信息。相反，如果我们对某件事已经有了较多的了解，我们不需要太多的信息就能把它搞清楚。
- 对于任意一个随机变量  $X$ ，熵定义如下：“变量的不确定性越大，熵也就越大，把它搞清楚所需要的信息量也就越大。”

- 世界杯谁是冠军？
- 世界杯赛后问一个知道结果的观众“哪支球队是冠军”？他不愿意直接告诉我，而要让我猜，并且我每猜一次，他要收一元钱才肯告诉我是否猜对了，那么我需要付给他多少钱才能知道谁是冠军呢？
- 我可以把球队编上号，从 1 到 32，然后提问：“冠军的球队在 1-16 号中吗？”假如他告诉我猜对了，我会接着问：“冠军在 1-8 号中吗？”假如他告诉我猜错了，我自然知道冠军队在 9-16 中。这样最多只需要五次，我就能知道哪支球队是冠军
- 谁是世界杯冠军这条消息的信息量值五块钱

# Entropy (熵)



- 不需要猜五次就能猜出谁是冠军，巴西、德国、意大利这样的球队得冠军的可能性比美国、韩国等队大的多。
- 第一次猜测时不需要把 32 个球队等分成两个组，而可以把少数几个最可能的球队分成一组，把其它队分成另一组。然后我们猜冠军球队是否在那几只热门队中。
- 重复这样的过程，根据夺冠概率对剩下的候选球队分组，直到找到冠军队。也许三次或四次就猜出结果。
- 当每个球队夺冠的可能性（概率）不等时，“谁世界杯冠军”的信息量的信息量比五比特少。

— “谁是世界杯冠军”的信息量：

$$= - (p_1 \log p_1 + p_2 \log p_2 + \dots + p_{32} \log p_{32}),$$

—  $p_1, \dots, p_{32}$  是 32 个球队各自夺冠的概率

- 课外阅读：《数学之美》第六章“信息的度量与作用”

- 基于熵的离散化方法
  - 最大化区间的纯度

$$e_i = -\sum_{j=1}^k p_{ij} \log_2 p_{ij}$$

首先，需要定义熵（entropy）。设  $k$  是不同的类标号数， $m_i$  是某划分的第  $i$  个区间中值的个数，而  $m_{ij}$  是区间  $i$  中类  $j$  的值的个数。第  $i$  个区间的熵  $e_i$  由如下等式给出

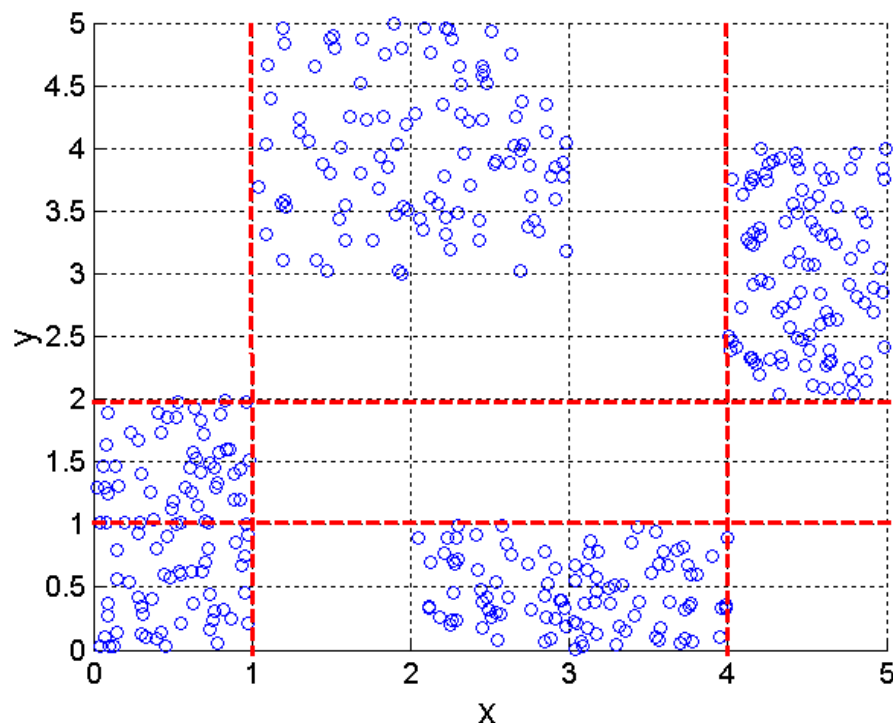
$p_{ij} = m_{ij}/m_i$  是第  $i$  个区间中类  $j$  的概率（值的比例）。



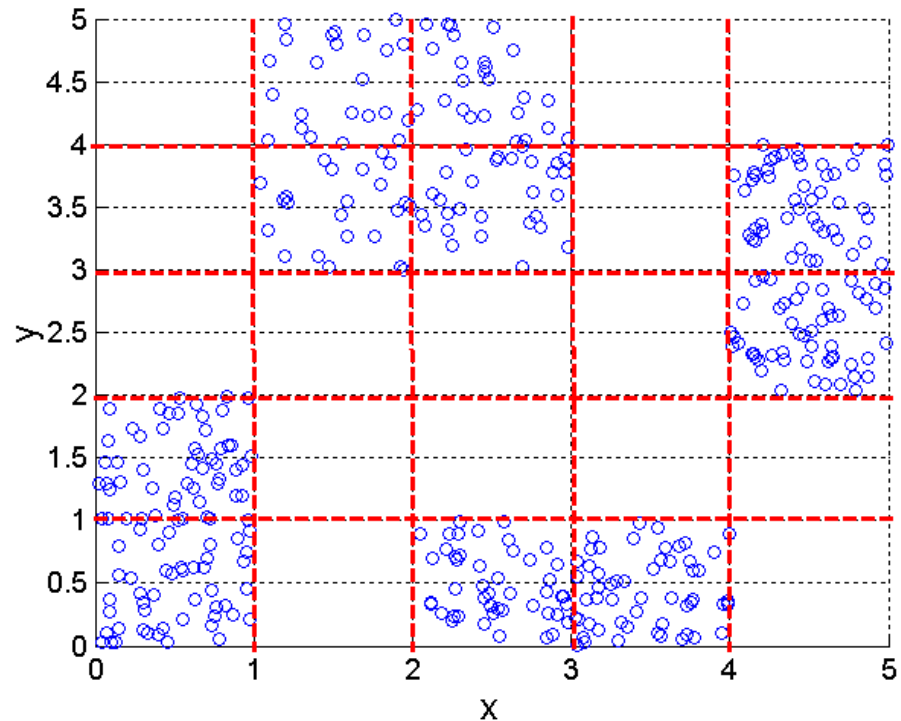
# Supervised Discretization



- 熵：区间纯度的度量
  - 只包含一个类：熵为0
  - 包含所有类，并且每类出现的概率相等：熵最大



3 categories for both x and y

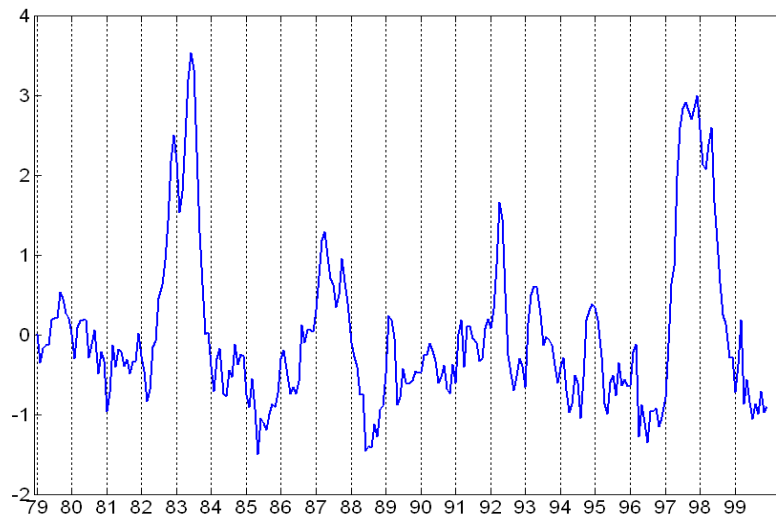


5 categories for both x and y

# 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



# 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

# 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 $	$s = -d, s = \frac{1}{1+d} \text{ or } s = 1 - \frac{d - \min\_d}{\max\_d - \min\_d}$

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

# Euclidean Distance



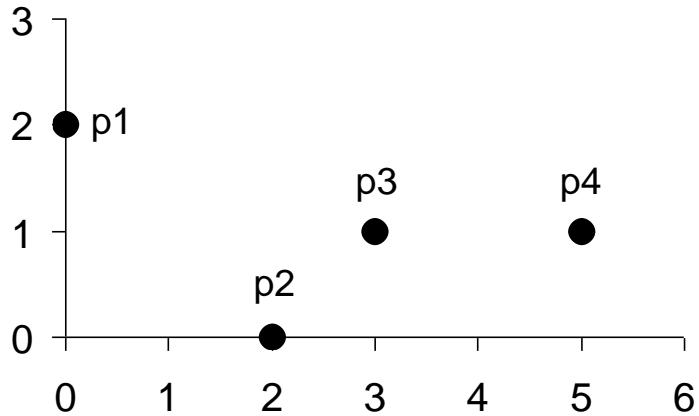
- Euclidean Distance

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

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

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

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  $q$ .

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

# 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**

# 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, 非负性)
2.  $d(p, q) = d(q, p)$  for all  $p$  and  $q$ . (Symmetry, 对称性)
3.  $d(p, r) \leq d(p, q) + d(q, r)$  for all points  $p, q$ , and  $r$ . (Triangle Inequality, 三角不等式)

where  $d(p, q)$  is the distance (dissimilarity) between points (data objects),  $p$  and  $q$ .

- A distance that satisfies these properties is a **metric** (度量)

# Example: Non-metric dissimilarities



- $A=\{1,2,3,4\}$     $B=\{2,3,4\}$
- $A-B = \{1\}$     $B-A=\emptyset$
- $\text{dis}(A,B) = \text{size}(A - B) = 1$
- $\text{dis}(B,A) = \text{size}(B - A) = 0$
- **$\text{dis}(A,B) = \text{size}(A-B) + \text{size}(B-A)$**

# Example: Non-metric dissimilarities



- Distance between time of the day:
- $d(t_1, t_2) = t_2 - t_1$  if  $t_1 \leq t_2$
- $d(t_1, t_2) = 24 + (t_2 - t_1)$  if  $t_1 > t_2$
- $d(1\text{PM}, 2\text{PM}) = 1$        $d(2\text{PM}, 1\text{PM}) = 23$

# 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



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

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

$M_{01} = 2$  (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)

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

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

# Extended Jaccard Coefficient (Tanimoto)



- Variation of Jaccard for continuous or count attributes
  - Reduces to Jaccard for binary attributes

$$T(p, q) = \frac{p \bullet q}{\|p\|^2 + \|q\|^2 - p \bullet q}$$

两个向量的交集

两个向量的并集

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



# Cosine Similarity



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

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

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

- Example:

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

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

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

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

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

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

$$A \cdot B = \|A\| * \|B\| * \cos(A, B); \quad A = d_1 / \|d_1\|, \quad B = d_2 / \|d_2\|$$

# Correlation (PCC皮尔森相关性)



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

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

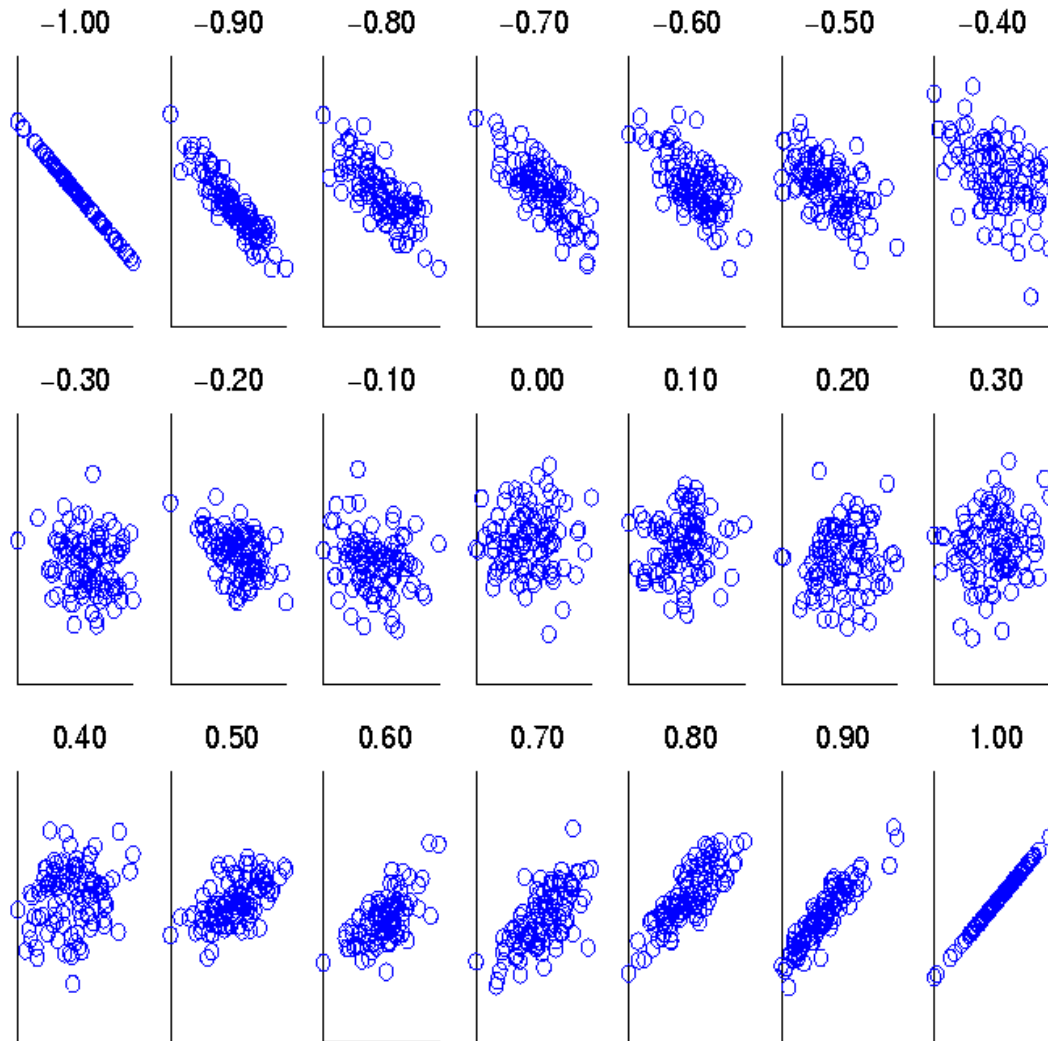
$$q'_k = (q_k - \text{mean}(q)) / \text{std}(q)$$

$$\text{correlation}(p, q) = p' \bullet q'$$

先标准化，再内积；内积代表相似性

**Cosine vs PCC:** 标准化的过程不同

# 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.

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

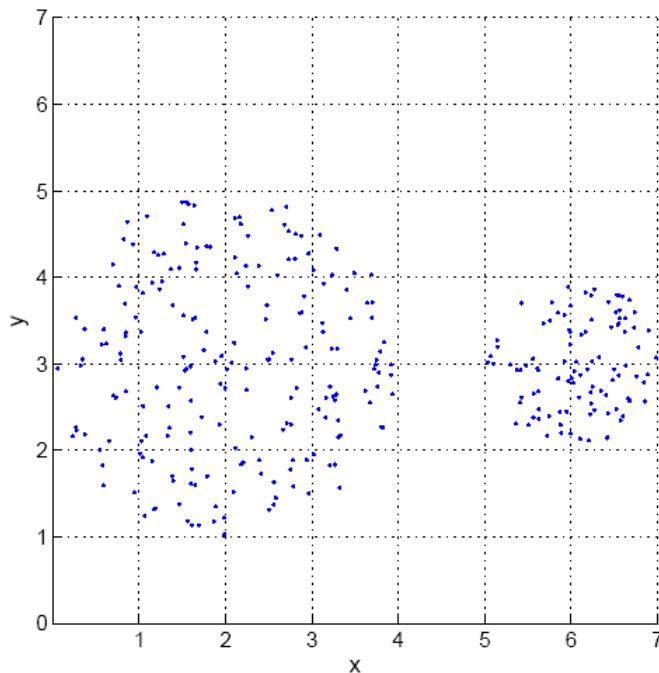
$$\text{distance}(p, q) = \left( \sum_{k=1}^n w_k |p_k - q_k|^r \right)^{1/r}$$

- Density-based clustering require a notion of density
- Examples:
  - Euclidean density
    - ◆ Euclidean density = number of points per unit volume
  - Probability density
  - Graph-based density

# Euclidean Density – Cell-based



- Simplest approach is to divide region into a number of rectangular cells of equal volume and define density as # of points the cell contains



**Figure 7.13.** Cell-based density.

0	0	0	0	0	0	0
0	0	0	0	0	0	0
4	17	18	6	0	0	0
14	14	13	13	0	18	27
11	18	10	21	0	24	31
3	20	14	4	0	0	0
0	0	0	0	0	0	0

**Table 7.6.** Point counts for each grid cell.

# Euclidean Density – Center-based



- Euclidean density is the number of points within a specified radius of the point

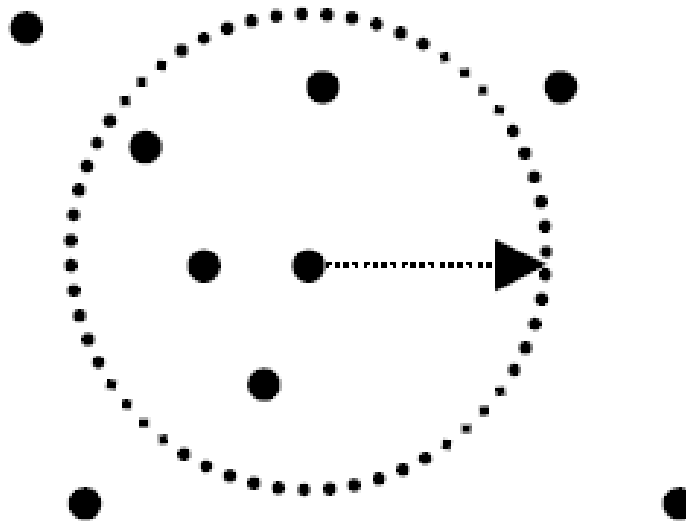


Figure 7.14. Illustration of center-based density.



# Data Preprocessing



- Aggregation
- Sampling
- Dimensionality Reduction
- Feature subset selection
- Feature creation
- Discretization and Binarization
- Attribute Transformation