



Data Preprocessing

—About data—

徐华

清华大学 计算机系 智能技术与系统国家重点实验室

xuhua@tsinghua.edu.cn

1

Data Preprocessing



- ◉ **About data**
- ◉ **Why preprocess the data?**
- ◉ **Descriptive data summarization**
- ◉ **Data cleaning**
- ◉ **Data integration and transformation**
- ◉ **Data reduction**
- ◉ **Discretization and concept hierarchy generation**
- ◉ **Summary**

2



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

| 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 |

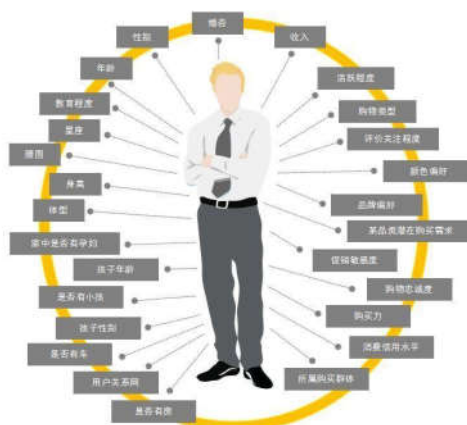
Objects

3



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



4



Types of Attributes



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

5



Properties of Attribute Values



- ◉ The type of an attribute depends on which of the following properties it possesses:
 - ◆ **Distinctness:** = ≠
 - ◆ **Order:** < >
 - ◆ **Addition:** + -
 - ◆ **Multiplication:** * /

 - ◆ **Nominal attribute:** distinctness
 - ◆ **Ordinal attribute:** distinctness & order
 - ◆ **Interval attribute:** distinctness, order & addition
 - ◆ **Ratio attribute:** all 4 properties

6



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.



7

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



8

Important Characteristics of Structured Data



Dimensionality (维度)

Curse of Dimensionality (维数灾难)

Sparsity (稀疏)

Only presence counts

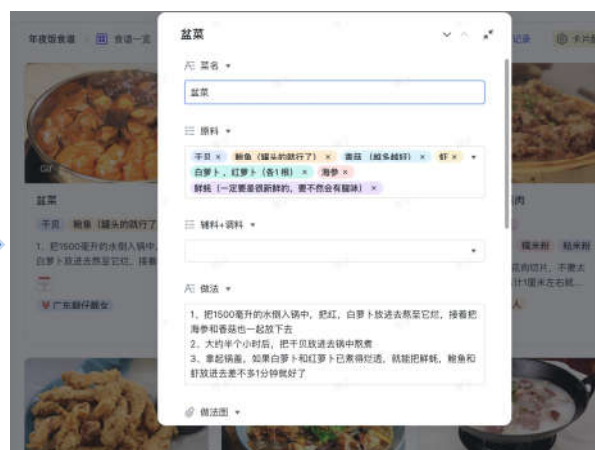
Resolution (分辨率, 解析度)

Patterns depend on the scale

| 频道 | 节目名称 | 播出时间 | 节目类型 | 备注 |
|----|------|----------------|------|-----------|
| 1 | 杜尚涛 | 20230204 19:30 | 综艺 | 01-《综艺盛典》 |
| 2 | 沈涛 | 20230204 19:30 | 综艺 | 01-《综艺盛典》 |
| 3 | 杨云萍 | 20230204 19:30 | 综艺 | 02-《综艺盛典》 |
| 4 | 高敏涛 | 20230204 19:30 | 综艺 | 02-《综艺盛典》 |
| 5 | 黄雅 | 20230204 19:30 | 综艺 | 03-《综艺盛典》 |
| 6 | 高海平 | 20230204 19:30 | 综艺 | 04-《综艺盛典》 |
| 7 | 郭涛 | 20230204 19:30 | 综艺 | 05-《综艺盛典》 |
| 8 | 魏涛 | 20230204 19:30 | 综艺 | 06-《综艺盛典》 |
| 9 | 陈涛 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |
| 10 | 魏江江 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |
| 11 | 郭涛 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |
| 12 | 郭涛 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |
| 13 | 郭涛 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |
| 14 | 郭涛 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |
| 15 | 郭涛 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |
| 16 | 郭涛 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |
| 17 | 郭涛 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |
| 18 | 郭涛 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |
| 19 | 郭涛 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |
| 20 | 郭涛 | 20230204 19:30 | 综艺 | 07-《综艺盛典》 |

9

Important Characteristics of Structured Data



10



Record Data



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

| 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 |

11



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 |

12



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 | 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 |

13



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.

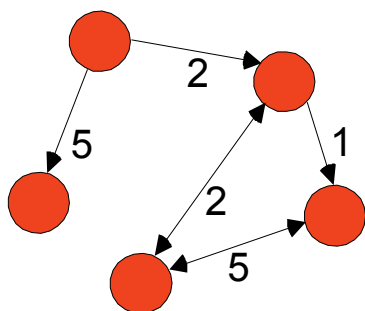
| <i>TID</i> | <i>Items</i> |
|------------|---------------------------|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |

14



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
  
```

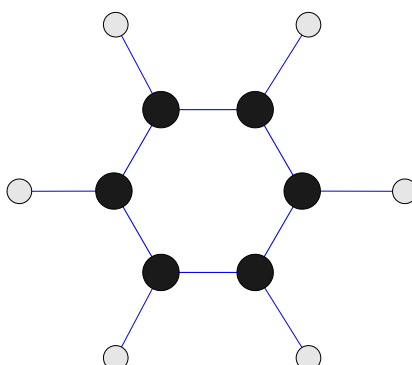


15



Chemical Data

Benzene Molecule (苯分子): C_6H_6



16

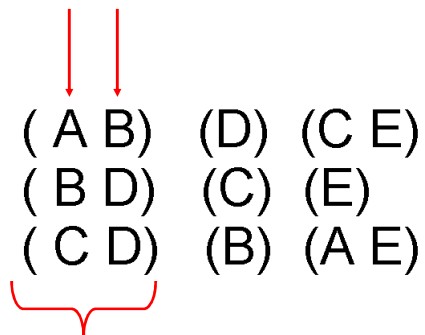


Ordered Data



Sequences of transactions

Items/Events



An element of
the sequence

17



Ordered Data



Genomic (染色体) sequence data

```
GGTTCCGCCTTCAGCCCCGCGCC
CGCAGGGCCCGCCCCGCGCCGTC
GAGAAGGGCCCGCCTGGCGGGCG
GGGGGAGGCGGGGCCGCCGAGC
CCAACCGAGTCCGACCAGGTGCC
CCCTCTGCTCGGCCTAGACCTGA
GCTCATTAGGCGGCAGCGGACAG
GCCAAGTAGAACACGCGAAGCGC
TGGGCTGCCTGCTGCGACCAGGG
```

18

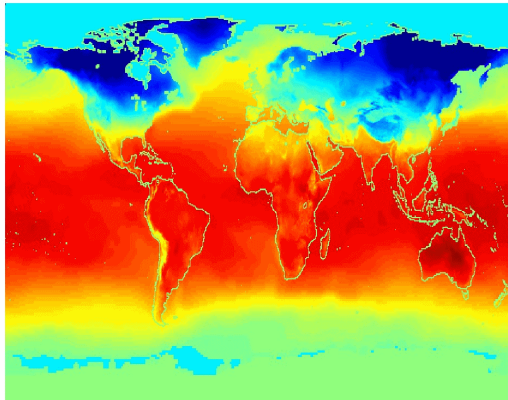


Ordered Data

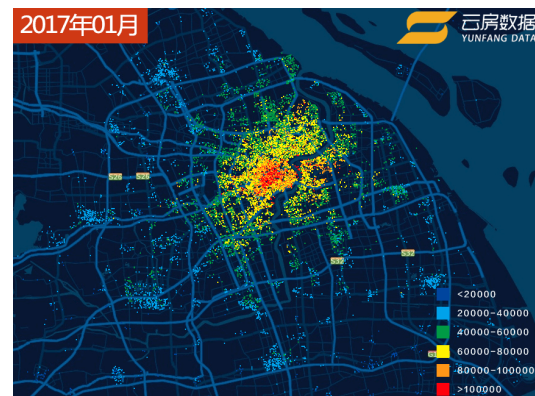


○ Spatio-Temporal Data

Jan



Average Monthly Temperature of Land and Ocean



Heat map of human flow



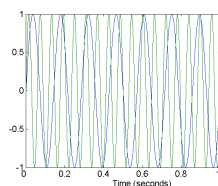
Data Quality



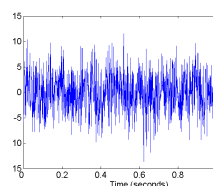
- What kinds of data quality problems?
- How can we detect problems with the data?
- What can we do about these problems?

○ Examples of data quality problems:

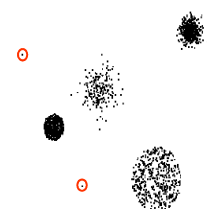
- ◆ Noise and outliers
- ◆ Missing values
- ◆ Duplicate data



Two Sine Waves



Two Sine Waves + Noise



Outliers





Thanks !

