



Data Preprocessing

——Why pre-process 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



Why Data Preprocessing(1)



◎ Data in the real world is dirty.

- ◆ **Incomplete(不完整)**: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
 - e.g., occupation= ""
- ◆ **Noisy (有噪音)**: containing errors or outliers
 - e.g., Salary= "-10"
- ◆ **Inconsistent (不一致)**: containing discrepancies (冲突) in codes or names:
 - External discrepancies
 - e.g., Age= "42" Birthday= "03/07/1997"
 - e.g., Was rating "1,2,3" , now rating "A, B, C"
 - e.g., discrepancy between duplicate records
 - Internal discrepancies
 - e.g., IngrA(10)+IngrB(3)+IngrC(4) -> Germ(70%)
 - IngrA(13)+IngrB(2)+IngrC(4) -> Germ(65%)

3



Why Data Preprocessing(2)



◎ Incomplete data comes from

- ◆ Different consideration between the time when the data was collected and when it is analyzed.
- ◆ Human/hardware/software problems

◎ Noisy data comes from the process of data

- ◆ Collection
- ◆ Entry
- ◆ Transmission
- ◆ Conflict with common sense

◎ Inconsistent data comes from

- ◆ Different data sources (Web , Manual Collections, Special Equipments, Database)
- ◆ Actual Experiment Equipments (Sensors)
- ◆ Different Environment Conditions (IOT Equipments)

4



Why Data Preprocessing(3)



◎ The form of data/information needs to be transformed.

◆ Different data need to be integrated.

- e.g. In Table A: Age = "" ;
- e.g. In Table B: Weight= ""

◆ Different data need to be transformed.

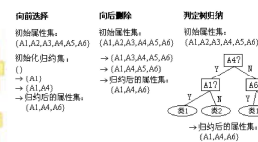
- e.g. questionnaire data

◆ Different data need to be discretized.

◆ Different data need to be reduced.



姓名	性别	年龄	身高	体重	血型	职业	学历	收入	消费	资产	负债	保险	投资	其他
张三	男	25	1.75	70	A	程序员	本科	10000	5000	200000	100000	重疾险	股票	房产
李四	女	30	1.60	55	B	教师	硕士	8000	3000	150000	50000	医疗险	基金	无
王五	男	45	1.80	85	O	工程师	本科	12000	6000	250000	120000	重疾险	股票	房产
赵六	女	28	1.65	60	A	设计师	本科	9000	4000	180000	80000	医疗险	基金	无
孙七	男	35	1.70	65	B	销售经理	本科	11000	5500	220000	110000	重疾险	股票	房产



5

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

6



Multi-Dimensional Measure of Data Quality



- ◉ **A well-accepted multidimensional view:**
 - ◆ Accuracy (准确的)
 - ◆ Completeness (完整的)
 - ◆ Consistency (一致的)
 - ◆ Timeliness (合时的)
 - ◆ Believability (可信的)
 - ◆ Value added (有附加价值的)
 - ◆ Interpretability (可解释的)
 - ◆ Accessibility (可存取的)
- ◉ **Broad categories:**
 - ◆ Intrinsic(本质的), contextual (相关的) , representational (代表性的) , and accessibility (可存取的) .

7



Major Tasks in Data Preprocessing



- ◉ **Data cleaning**
 - ◆ Fill in missing values, smooth noisy data, identify or remove outliers, and resolve inconsistencies
- ◉ **Data integration**
 - ◆ Integration of multiple databases, data cubes, or files
- ◉ **Data transformation**
 - ◆ Normalization and aggregation
- ◉ **Data reduction**
 - ◆ Obtains reduced representation in volume but produces the same or similar analytical results
- ◉ **Data discretization**
 - ◆ Part of data reduction but with particular importance, especially for numerical data

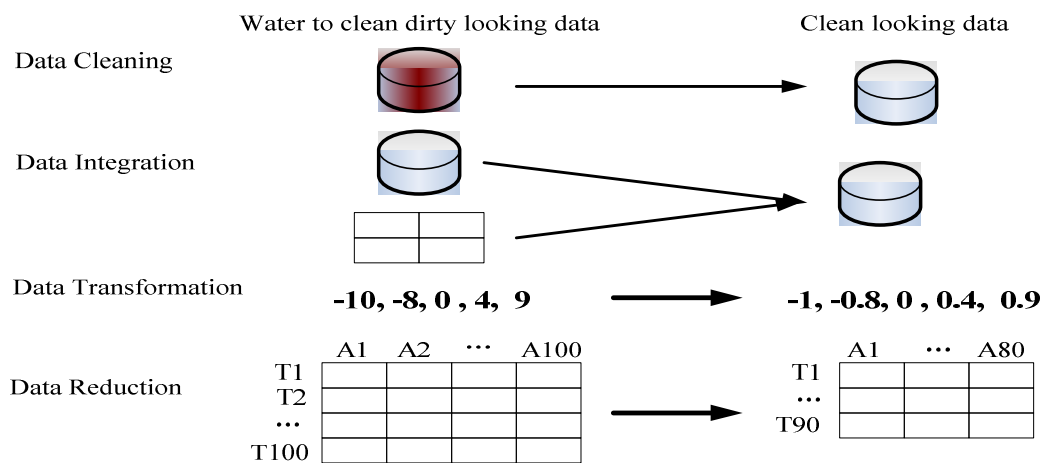
8



Form of data preprocess



Key Steps: Grasp and understand the data.



9



Thanks !

10

