



Data Preprocessing

—Data Cleaning—

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



Data Cleaning



Importance

- ◆ “Data cleaning is one of the three biggest problems in data warehousing” —Ralph Kimball
- ◆ “Data cleaning is the number one (No.1) 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

3



Missing Data



Data is not always available

- ◆ E.g., many tuples have no recorded value for several attributes, such as customer income in sales data

Missing data may be due to

- ◆ equipment malfunction
- ◆ inconsistent with other recorded data and thus deleted
- ◆ data not entered due to misunderstanding
- ◆ not register certain data may not be considered important at the time of entry
- ◆ history or changes of the data

Missing data may need to be inferred.

4



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(冗余) + 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

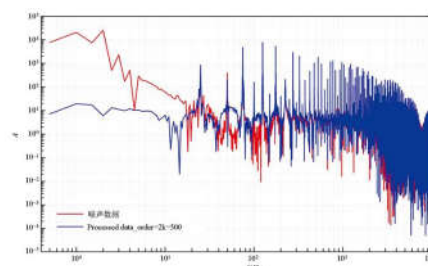
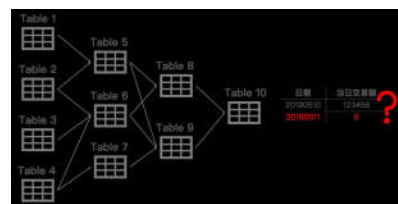
5



Noisy Data



- Noise: random error or variance in a measured variable
- Incorrect attribute values may due to
 - faulty data collection instruments
 - data entry problems
 - data transmission problems
 - technology limitation
 - inconsistency in naming convention (命名约定)
- Other data problems which requires data cleaning
 - duplicate records
 - incomplete data
 - inconsistent data



6

How to handle noisy data?



- ◉ Binning (分箱)
 - ◆ 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.
- ◉ Clustering
 - ◆ detect and remove outliers
- ◉ Combined computer and human inspection
 - ◆ detect suspicious values and check by human (e.g., deal with possible outliers)
- ◉ Regression
 - ◆ smooth by fitting the data into regression functions

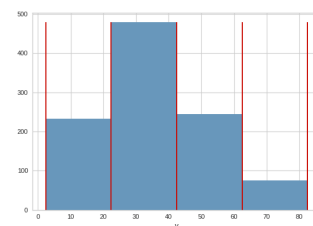
7



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.



8

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 (equal-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

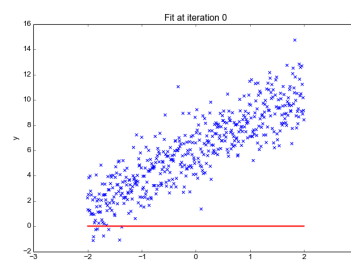
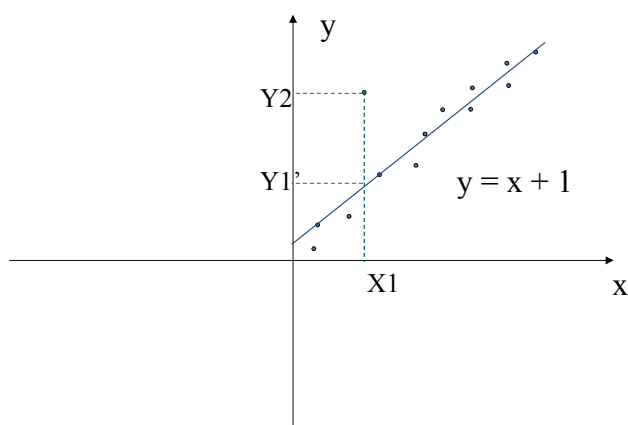
- ◆ Smoothing by bin boundaries:

- Bin 1: 4, 4, 4, 15
- Bin 2: 21, 21, 25, 25
- Bin 3: 26, 26, 26, 34

9



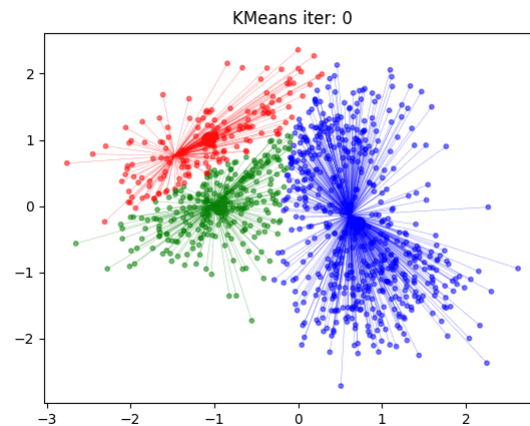
Regression



10



Cluster Analysis



11



Data Cleaning as a Process



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

12





Thanks !

