



CS 412 Intro. to Data Mining

Chapter 3. Data Preprocessing

Jiawei Han, Computer Science, Univ. Illinois at Urbana-Champaign, 2017



Chapter 3: Data Preprocessing

☐ Data Preprocessing: An Overview



☐ Data Cleaning

☐ Data Integration

☐ Data Reduction and Transformation

☐ Dimensionality Reduction

☐ Summary

What is Data Preprocessing? — Major Tasks

☐ **Data cleaning**

- ☐ Handle missing data, smooth noisy data, identify or remove outliers, and resolve inconsistencies

☐ **Data integration**

- ☐ Integration of multiple databases, data cubes, or files

☐ **Data reduction**

- ☐ Dimensionality reduction
- ☐ Numerosity reduction
- ☐ Data compression

☐ **Data transformation and data discretization**

- ☐ Normalization
- ☐ Concept hierarchy generation

Why Preprocess the Data? — Data Quality Issues

- Measures for data quality: A multidimensional view
 - Accuracy: correct or wrong, accurate or not
 - Completeness: not recorded, unavailable, ...
 - Consistency: some modified but some not, dangling, ...
 - Timeliness: timely update?
 - Believability: how trustable the data are correct?
 - Interpretability: how easily the data can be understood?

ความถูกต้อง
ความสมบูรณ์
การแก้ไขสัมพันธ์กัน
normalization
ไม่ได้!

5

Chapter 3: Data Preprocessing

- Data Preprocessing: An Overview
- Data Cleaning
- Data Integration
- Data Reduction and Transformation
- Dimensionality Reduction
- Summary

6

Data Cleaning

ข้อมูลสกปรก → ผิด / เกิน

- Data in the Real World Is Dirty: Lots of potentially incorrect data, e.g., instrument faulty, human or computer error, and transmission error
- Incomplete: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
 - e.g., *Occupation* = " " (missing data)
- Noisy: containing noise, errors, or outliers
 - e.g., *Salary* = "-10" (an error)
- Inconsistent: containing discrepancies in codes or names, e.g.,
 - Age* = "42", *Birthday* = "03/07/2010"
 - Was rating "1, 2, 3", now rating "A, B, C"
 - discrepancy between duplicate records
- Intentional (e.g., *disguised missing data*)
 - Jan. 1 as everyone's birthday?

ไม่สมบูรณ์

การผิด

ไม่สัมพันธ์กัน

→ default

7

Incomplete (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 were not entered due to misunderstanding
 - Certain data may not be considered important at the time of entry
 - Did not register history or changes of the data
- Missing data may need to be inferred

- เครื่องเสีย
- เครื่องไม่สัมพันธ์กัน/ตมกัน
- คนไม่บอก

- data เปลี่ยนแปลงไป

8

How to Handle Missing Data?

- ❑ Ignore the tuple: usually done when class label is missing (when doing classification)—not effective when the % of missing values per attribute varies considerably
- record 1 tuple missing ลบออก
- ❑ Fill in the missing value manually: tedious + infeasible? - 0101 แล้ว แก้วๆ 40 นาที
- ❑ Fill in it automatically with - ใส่ class ในนั้น
 - ❑ 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**