

Chapter 3: Data Preprocessing

พิช 407. วุฒิวิจัย

□ Data Preprocessing: An Overview

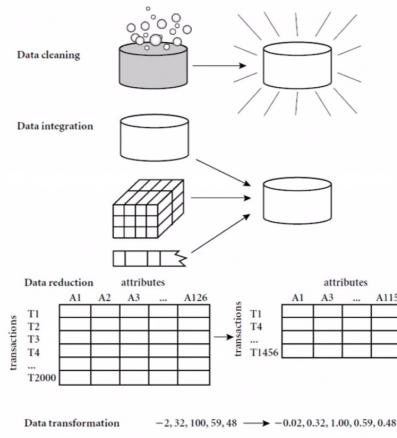
□ Data Cleaning

□ Data Integration

□ Data Reduction and Transformation

□ Dimensionality Reduction

□ Summary



ต่อไป Data ที่รับเข้ามายัง

What is Data Preprocessing? — Major Tasks

□ Data cleaning คือ

- Handle missing data, smooth noisy data, identify or remove outliers, and resolve inconsistencies ข้อมูลหายไป, ข้อมูลเสียงด้วย, ข้อความไม่ถูกต้อง

□ Data integration

- Integration of multiple databases, data cubes, or files รวม Data ที่มีอยู่ในรูปแบบต่างๆ

□ Data reduction

ลดขนาด หรือ ลดจำนวน

- Dimensionality reduction
- Numerosity reduction
- Data compression

□ Data transformation and data discretization

เปลี่ยนแปลงรูปแบบข้อมูลให้เข้ากับการดำเนินการประมวลผลได้

- 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? *— Data ที่น่าเชื่อถือหรือไม่*

- Interpretability: how easily the data can be understood?

คุณภาพของข้อมูล

เวลาอันที่ใช้ในการคำนวณ

Inconsistency ข้อความไม่ถูกต้อง

Data Cleaning

ព័ត៌មានក្នុងពេលវេលា មានស្ម័គ្រប់របស់ខ្លួន

- ❑ 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” *ចំណាំនាង់សម្រាប់ការ*
 - ❑ Was rating “1, 2, 3”, now rating “A, B, C”
 - ❑ discrepancy between duplicate records — *Data mining* នឹងស្វែងរកជាបញ្ជី នៃការសម្រាប់ការ
 - ❑ **Intentional** (e.g., *disguised missing data*)
 - ❑ Jan. 1 as everyone’s birthday? *ចំណាំនាង់សម្រាប់ការបង្កើតការបង្កើត*

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 *ការអនុវត្តន៍របស់ពេលវេលា*
 - ❑ Equipment malfunction
 - ❑ Inconsistent with other recorded data and thus deleted *ការអនុវត្តន៍របស់ពេលវេលាដែលបានលើក*
 - ❑ 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

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 *សារធានាអាស់នៅក្នុងពេលវេលា*
- ❑ Fill in the missing value manually: *tedious + infeasible?* *ស្ថានិភ័យ*
- ❑ Fill in it automatically with *ដែលវាន់ប្រាក់ដែលត្រូវបានដែឡើង*
 - ❑ a global constant : e.g., “**unknown**” a new class?!
 - ❑ the attribute mean *ម៉ឺនីតុលី*
 - ❑ the attribute mean for all samples belonging to the same class: smarter *ឥន្ទៃដែឡើងដែលស្ថានិភ័យ* *mean របស់ទីផ្សារ*
 - ❑ the most probable value: inference-based such as Bayesian formula or decision tree