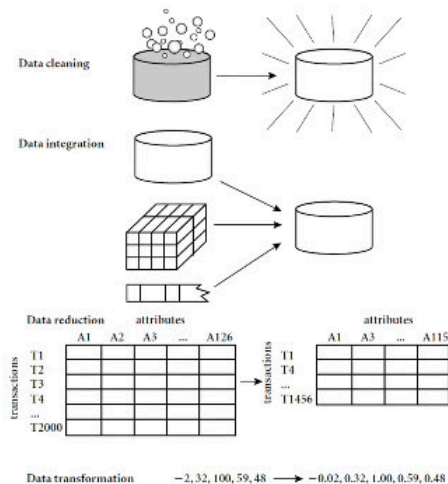


Chapter 3: Data Preprocessing

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



- Data Cleaning គឺជាការ cleaning data ដើម្បីកាត់បន្ថយការខុសឆ្គង
ដោយ ប្រើប្រាស់វិធីសាស្ត្រ កាត់ ដើម្បីកាត់បន្ថយការខុសឆ្គង ក្នុងទិន្នន័យ

- Data Integration គឺជា Data តាមលក្ខណៈប្រភេទផ្សេងៗគ្នា ត្រូវបានប្រមូល ទៅក្នុង
ប្រព័ន្ធនៃទិន្នន័យ Data Mining ដើម្បីធ្វើការស្វែងរក ក្នុង Data Warehouse ដើម្បីធ្វើការស្វែងរក
ទិន្នន័យ

- Data Retraction and transformation គឺជា ការកាត់បន្ថយ / ការកែប្រែទិន្នន័យ ដើម្បី
ធ្វើការស្វែងរក

- Dimensionality Reduction គឺជា ការកាត់បន្ថយទិន្នន័យ

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

- Data cleaning រំលង missing, inconsistencies / កំហុស noisy, outlier
- Data integration ឧទាហរណ៍ Data ចម្បងៗ រួមបញ្ចូលទៅក្នុង database
- Data reduction កាត់បន្ថយទិន្នន័យ
- Data transformation ប្តូរទិន្នន័យ ឲ្យស្របគ្នា

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?

- หน้าที่ของ Preprocess คือ: Data anomaly วนซ้ำ

ขั้นตอนที่สำคัญในการ Preprocessing

Data Cleaning: Incomplete (ไม่สมบูรณ์), Noisy, Inconsistent (ขัดแย้งกัน), Intentional

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

เช่น เราได้ข้อมูล 1
จากข้อมูลทั่วไปในแบบฟอร์ม พอที่มันเราพอได้พอออกได้ 100% 100% 100% แบบนี้เรียกว่า Missing data
เพราะมันมีค่าที่หายไปเป็น 0 เป็น 1

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
- ❑ Fill in the missing value manually: tedious + 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**

អ្នក Data ដែល Missing ក៏មិនសំខាន់ទៅ ទៅកាន់ក្នុងប្រព័ន្ធដែលយើង ខ្លាចដែលក្នុងការកែច្នៃ ទៅហៅសម្រាប់
ក្នុងការស្រាវជ្រាវ