



CS 412 Intro. to Data Mining

Chapter 3. Data Preprocessing



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



Chapter 3: Data Preprocessing

សាស្ត្រឃើមអំពីរក្សា, នរណ៍វិធានភាពអំពីរក្សា

□ Data Preprocessing: An Overview

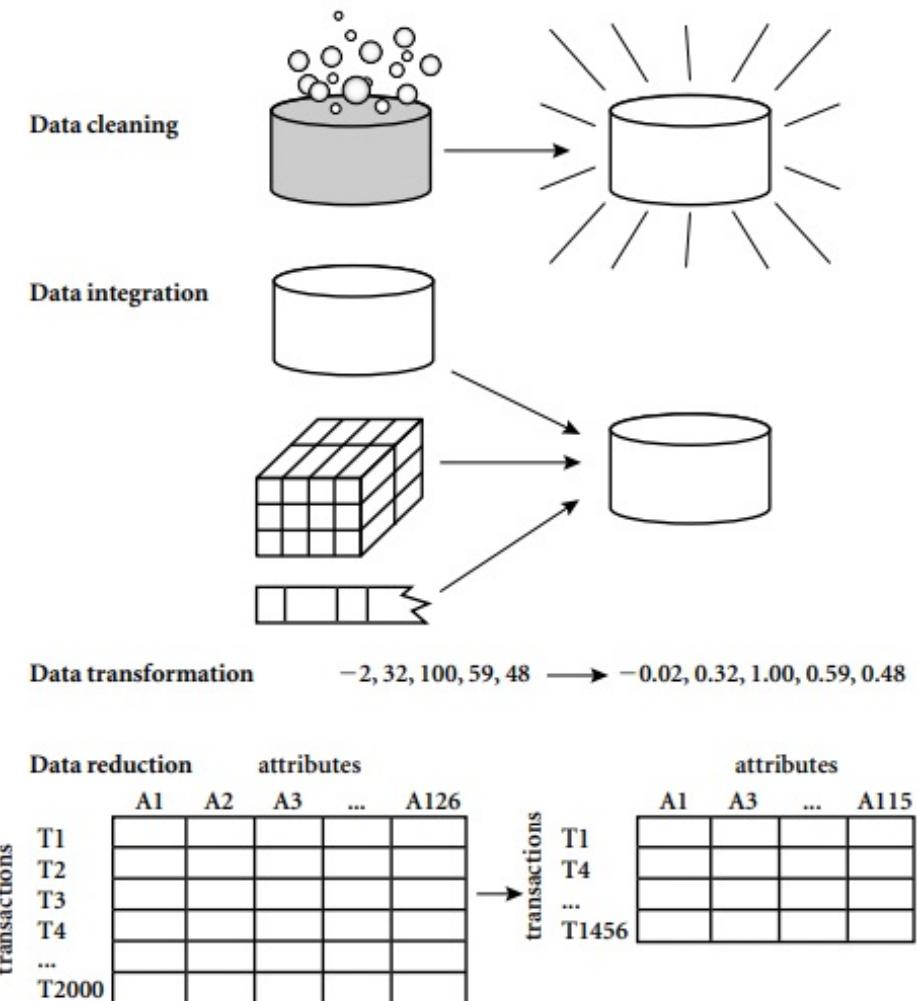
□ Data Cleaning

□ Data Integration

□ Data Reduction and Transformation

□ Dimensionality Reduction

□ Summary



What is Data Preprocessing? — Major Tasks

Why Preprocess the Data? — Data Quality Issues

ពិនិត្យការពារ
Data Preprocess

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

Chapter 3: Data Preprocessing

- Data Preprocessing: An Overview

- Data Cleaning

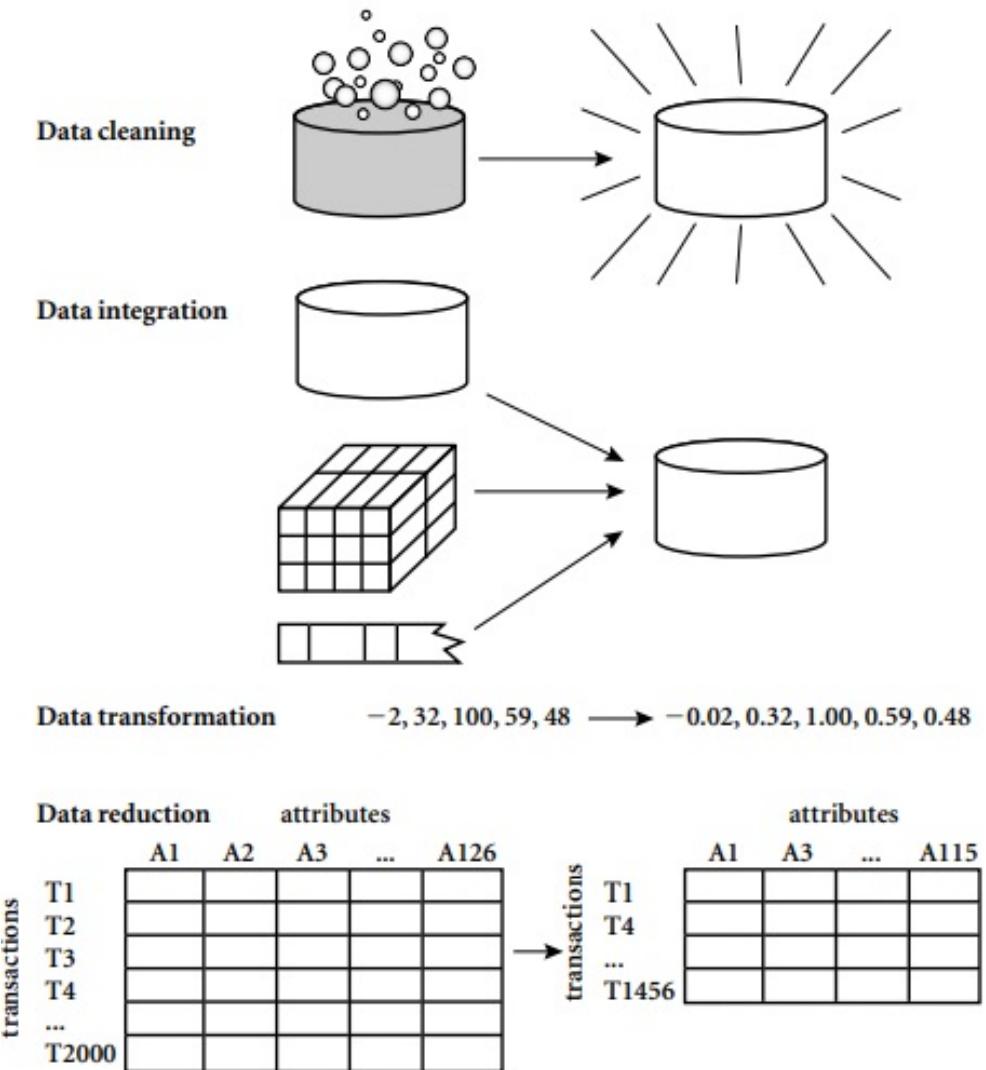


- Data Integration

- Data Reduction and Transformation

- Dimensionality Reduction

- Summary



Data Cleaning

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



Incomplete (Missing) Data

វត្ថុទិន្នន័យមាត្រា

- ❑ 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 **កំណត់តារា**
 - ❑ Equipment malfunction **កេរ៉ាសៀវភៅ**
 - ❑ Inconsistent with other recorded data **ទីតាំងកំណត់តារា** 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



How to Handle Missing Data?

ເກົ່າໃຫຍ້ກໍ່າຂອງ missing value

ລົບອອກເຈນ

ແຕ່ກໍ່າ data ພ້ອມຂັ້ນຈະຄຳບາດ

- 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 ກວດກາທິນໂຄສະນາ ດາວໂຫຼວດ
 - a global constant : e.g., “unknown”, a new class?!
 - the attribute mean ເອົາກໍ່າກາງວາງນໍາຕໍ່missing
 - the attribute mean for all samples belonging to the same class: smarter ແກ້ນດັ່ງ mean ກົດຢູ່ໃນ class ເຕັມກັນ
 - the most probable value: inference-based such as Bayesian formula or decision tree**



Noisy Data

- Noise:** random error or variance in a measured variable
- Incorrect attribute values** may be due to
 - Faulty data collection instruments
 - Data entry problems
 - Data transmission problems
 - Technology limitation
 - Inconsistency in naming convention
- Other data problems**
 - Duplicate records
 - Incomplete data
 - Inconsistent data



How to Handle Noisy Data?

