



CS 412 Intro. to Data Mining

Chapter 3. Data Preprocessing

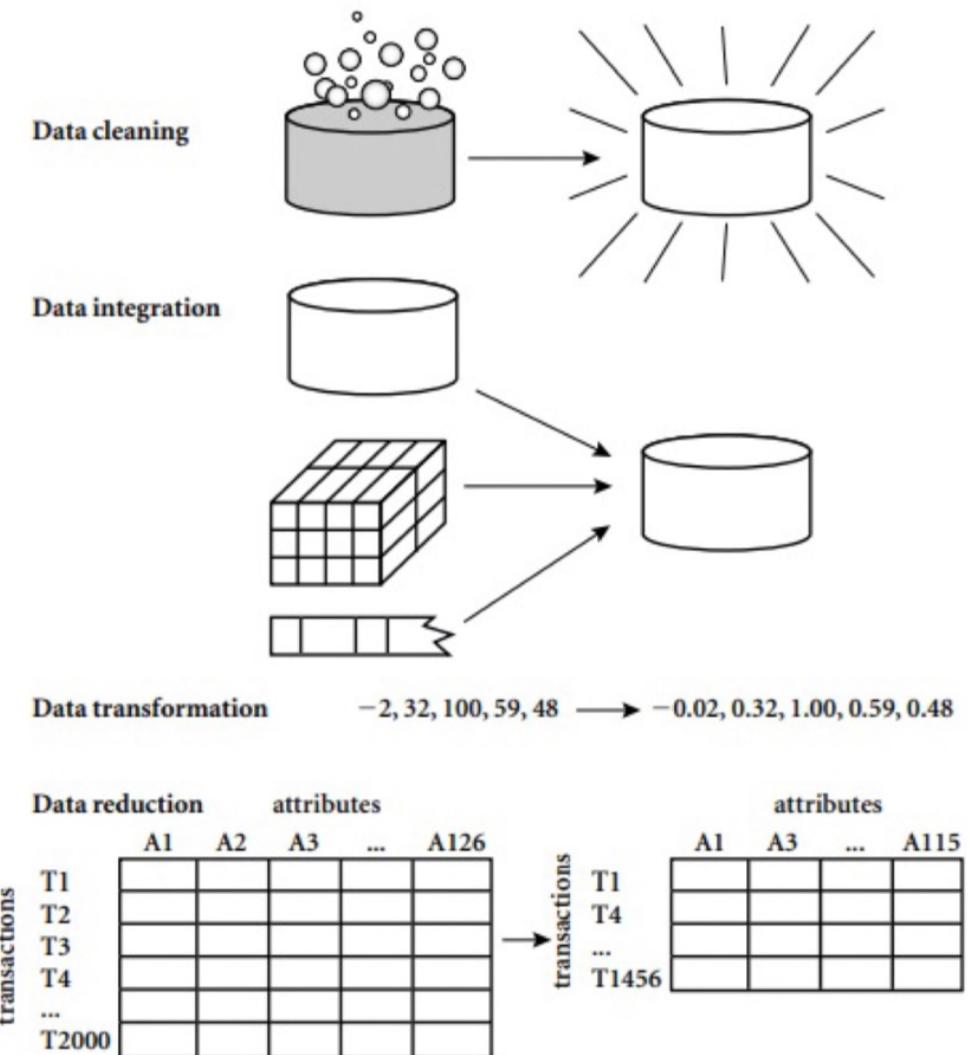
Jiawei Han, Computer Science, Univ. Illinois at Urbana-Champaign, 2017



Chapter 3: Data Preprocessing

ການຕັ້ງທຳມູນຄ່ອງມືການປະເພດຂໍ້ຕົວ

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



What is Data Preprocessing? — Major Tasks

ជំនួយនៃការសម្រេចដោយបង្កើតរាយការណ៍

□ Data cleaning

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

□ Data integration

គ្រប់គ្រាន់ទៅការបញ្ជូន ឧបាទ់

- Integration of multiple databases, data cubes, or files

□ Data reduction

- Dimensionality reduction

តាមលក្ខណៈការពិនិត្យ

- Numerosity reduction

- Data compression

ប្រើប្រាស់បច្ចុប្បន្ន ដើម្បីបង្កើតការងារបច្ចុប្បន្ន

□ Data transformation and data discretization

ព័ត៌មានត្រូវបានក្រោមការសម្រេច ឬពិនិត្យដើម្បីបង្កើតការងារបច្ចុប្បន្ន

- Normalization

- Concept hierarchy generation

របស់ជីវិត

Why Preprocess the Data? — Data Quality Issues

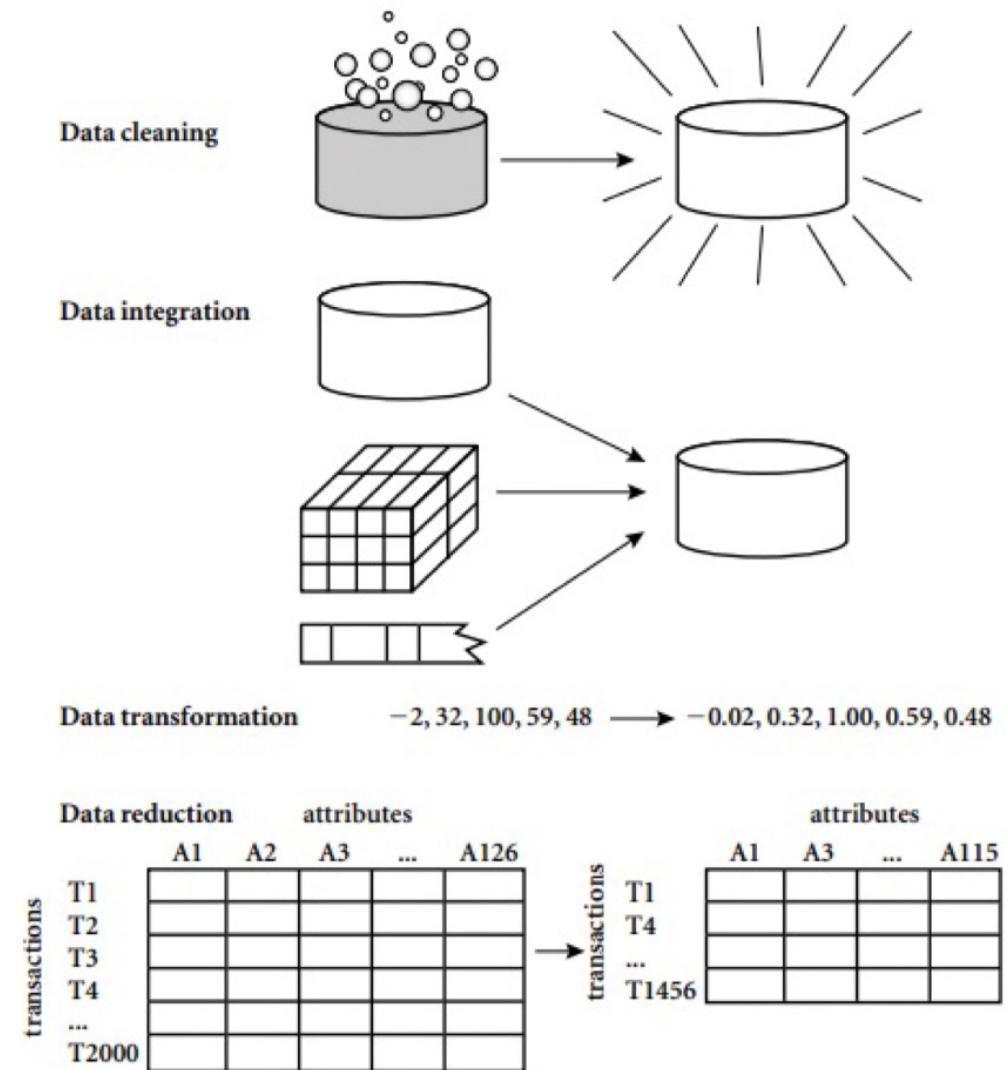
ພິບສອນໃຈ Data Preprocess

❑ Measures for data quality: A multidimensional view

- ❑ Accuracy: correct or wrong, accurate or not ອານຸຍານຕາມວັນ, ໂກງວ່າໄດ້
- ❑ Completeness: not recorded, unavailable, ... ລົມທີ່ບໍ່ໄດ້ຮັບຮັບ
- ❑ Consistency: some modified but some not, dangling, ... ຖໍ່ນormalization ມີຄືກົດຊົງທີ່ບໍ່ໄດ້ຮັບຮັບ
- ❑ Timeliness: timely update? Data ດີວຽວເປັນໄວ້ ດີວຽວເປັນໄວ້
- ❑ 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

ក្រោចការណ៍ពាណិជ្ជកម្ម

- ❑ Data in the Real World Is Dirty: Lots of potentially incorrect data, e.g., instrument faulty, human or computer error, and transmission error *ស្ថិតិភាពអាចមិនត្រឹមត្រូវបានផ្តល់ទៅបានក្នុងការបញ្ចូនការ*
- ❑ 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) *ខ្លួន ត្រូវតើម្គាល់ឡើង*
- ❑ Inconsistent: containing discrepancies in codes or names, e.g.,
 - ❑ *Age* = “42”, *Birthday* = “03/07/2010” *ចាប់ពីរយៈពេលការបង្កើត ដោយគ្រប់គ្រង*
 - ❑ Was rating “1, 2, 3”, now rating “A, B, C” *តាមរយៈពេលការបង្កើត ដោយគ្រប់គ្រង*
 - ❑ discrepancy between duplicate records
- ❑ Intentional (e.g., *disguised missing data*)
 - ❑ Jan. 1 as everyone’s birthday?

Incomplete (Missing) Data

ମୁଖ୍ୟ ହେତୁ କୁଣ୍ଡଳା

- 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

How to Handle Missing Data?

សំណង់ទិន្នន័យទូទៅ missing value

របាយការ: ពន្លាត់សម្រាប់អាជីវកម្ម

- 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 *នរណា* *សម្រាប់* *ទិន្នន័យ*
 - a global constant : e.g., “unknown”, a new class?!
 - the attribute mean *បន្ទាន់សម្រាប់* *ទិន្នន័យ*
 - the attribute mean for all samples belonging to the same class: smarter *បន្ទាន់សម្រាប់* *mean* *នៃ* *ទិន្នន័យ*
 - the most probable value: inference-based such as Bayesian formula or decision tree**

Noisy Data

សំណុំតម្រូវការ

- **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

Noisy Data
សំណុំតម្រូវការ និងសំណុំតម្រូវការដែលមានការរៀបចំឡើង និង ការចាប់ផ្តើមទាមពេលវេលា
តាមរយៈការ
- រូបភាពកំប្រឈប់បានបិទិន្ទុលការអនុវត្តន៍ឡាលូ
- ក្នុងការប្រើប្រាស់ការបង្ហើតនៃ គោលការណ៍
- ក្នុងការបង្ហើតនៃការបង្ហើតអាណាពាល
- សំណុំតម្រូវការបានបិទិន្ទុលការ

How to Handle Noisy Data?

- Binning *ការបិនតម្លៃទៅសំណើនូវការរួមចំណាំការពារ*

 - 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.

- Regression *គុណភាពការកែវិនិយោគ - ធនធានរបស់វិប័យ*

 - Smooth by fitting the data into regression functions

- Clustering *ការបិនតាមពេលវេលាប្រភេទការងារ*

 - Detect and remove outliers

- Semi-supervised: Combined computer and human inspection

 - Detect suspicious values and check by human (e.g., deal with possible outliers)