



CS 412 Intro. to Data Mining

Chapter 3. Data Preprocessing

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



Chapter 3: Data Preprocessing

ການເຕັ້ມມັງງູບ, ດາວໂຫຼວງຂອງພະຍານືອງ

- Data Preprocessing: An Overview
- Data Cleaning
- Data Integration *(Data ອານວາຍແລ້ວການກົດ)*
 - ດັວວິນານ Data ໂດຍ Dimension*
- Data Reduction and Transformation
- Dimensionality Reduction
- Summary

What is Data Preprocessing? — Major Tasks

เป็นขั้นตอนสำหรับการตัดสินใจว่า ไม่ต้องสนใจอะไร

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

เป็นขั้นตอนการตัดสินใจว่า ไม่ต้องสนใจอะไร

Why Preprocess the Data? — Data Quality Issues

- Measures for data quality: A multidimensional view

- Accuracy: correct or wrong, accurate or not *ត្រូវតាមរយៈការសម្រេច*
- Completeness: not recorded, unavailable, ... *ក្នុងសម្រាប់ប្រើប្រាស់ទៅទំនាក់ទំនង*
- Consistency: some modified but some not, dangling, ... *នៅលីមិនមានការសម្រេចនៅក្នុងគ្រប់គ្រាន់*
- Timeliness: timely update? *តារាងនៃពេលវេលាដែលត្រូវបានដោះស្រាយឡើង*
ត្រូវបានដោះស្រាយឡើង?
- Believability: how trustable the data are correct?
- Interpretability: how easily the data can be understood?

Chapter 3: Data Preprocessing

❑ Data Preprocessing: An Overview

ମୋଟାର୍କୁଲେସନ୍

❑ Data Cleaning



ମୋଟାର୍କୁଲେସନ୍

❑ Data Integration

ମୋଟାର୍କୁଲେସନ୍

ମୋଟାର୍କୁଲେସନ୍

❑ Data Reduction and Transformation

❑ Dimensionality Reduction

❑ Summary

Data Cleaning

ការបោតការងារពីពាណិជ្ជកម្ម

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

Data Twikyru

Incomplete (Missing) Data

Դաշտականացություն

- 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 *արագակ*
 - Inconsistent with other recorded data and thus deleted *լրացրեծ տվյալների հետ մասնաւոր չէ*
 - 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 *պահպան չեղանակ*
- Missing data may need to be inferred

How to Handle Missing Data?

ស្មារម៌ Missing

- Ignore the tuple: usually done when class label is missing (when doing classification)—not effective when the % of missing values per attribute varies considerably
Data record ទិន្នន័យ missing ស្ថាគភាព
- 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 អំពីរបាយនូវរឹង
 - 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**