



CS 412 Intro. to Data Mining

Chapter 3. Data Preprocessing

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





Chapter 3: Data Preprocessing

□ Data Preprocessing: An Overview

□ Data Cleaning

Data ที่เก็บมาจากการซื้อขาย มีน้ำเสียง

□ Data Integration

融通 Data จากหลายแหล่งมาไว้ด้วยกัน

□ Data Reduction and Transformation

□ Dimensionality Reduction

□ Summary



เก็บมา

Sensor - เก็บอัตโนมัติ

noise, missing เหลือบกัน

ข้อมูลที่ไม่จำเป็น

ลดจำนวนข้อมูล

ลด Dimension

What is Data Preprocessing? — Major Tasks

□ Data cleaning

ຈັດຕະ

ກຳທົດຖົວພັນ Outlier

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

→ ຕາມໄຟສົດຄະກົນ

□ Data integration

- Integration of multiple databases, data cubes, or files

□ Data reduction

- Dimensionality reduction
- Numerosity reduction
- Data compression

ແປ່ນແກ້ຂັ້ນດູງດູ

□ Data transformation and data discretization

- Normalization
- Concept hierarchy generation

Why Preprocess the Data? — Data Quality Issues

- Measures for data quality: A multidimensional view

စွမ့်ဆက္စည်ရေးနှင့်သာများနှင့်အနုဂါးလုပ်လုပ်ငန်း

- Accuracy: correct or wrong, accurate or not
- Completeness: not recorded, unavailable, ...
- Consistency: some modified but some not, dangling, ...
- Timeliness: timely update? → နဲ့ကြတ်ပေါ်လောက်ခဲ့လောက်ခဲ့
- 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

កែវងសំដែងអំពីរាជរាជ

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

ទម្រងទៅលាង
អាណាពាណិជ្ជកម្ម

តុលាងនឹងបូន្ថែមព័ត៌មាន
នៅចំណែកភាពអំពីរាជរាជ

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 *វិធាននៃការផ្តល់ព័ត៌មាន*
- 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 *Pata record မျှသော Missing တိုက်ခွဲမှု*
- 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 *ပေါ်စိမ်း mean အဖြတ်များ၏ missing*
 - the attribute mean for all samples belonging to the same class: smarter *mean n' အတွက်အတွက်*
 - **the most probable value: inference-based such as Bayesian formula or decision tree**