



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 ที่เก็บมามีนักยังเหลือ** **เก็บมา** **sensor - เก็บต่อไปนี้ดี** noise, missing, ไม่ต้อง
- Data Integration **เอา Data จากหลายแหล่งมารวมกัน**
- Data Reduction and Transformation **ลดตัว Data** **ลดตัว Dimension (แนวตั้ง)**
- Dimensionality Reduction
- Summary

noise ไฟฟ้าขุ่นๆ ไม่ใช่ภาพ Ex. ชุดเดียวแต่ไฟฟ้าตัวเดียว
missing Data ร่องรอยที่ไม่ได้เก็บไว้ ไม่ต้อง

noise,
missing, ไม่ต้อง

What is Data Preprocessing? — Major Tasks

ដែលត្រូវបានធ្វើឡើង

- Data cleaning** រាយការណ៍ទិន្នន័យ កែវាទិន្នន័យ
ស្រួលការណ៍ ឬចុចុចទិន្នន័យដោយការណ៍ សំងសាលាប្រឈម smooth
- Handle missing data, smooth noisy data, identify or remove outliers, and resolve inconsistencies
កែវាទិន្នន័យ ស្រួលការណ៍ ឬចុចុចទិន្នន័យ ស្រួលការណ៍ សំងសាលាប្រឈម
outliers → ចំណេះចុចុចទិន្នន័យ Binning Methods
Regression
Clustering
- Data integration** រាយការណ៍ទិន្នន័យ
កែបានការណ៍ទិន្នន័យ ឬកែបានការណ៍ទិន្នន័យ ដោយការបញ្ចប់ទិន្នន័យ
- Integration of multiple databases, data cubes, or files
- Data reduction** រាយការណ៍លក់ទិន្នន័យ កែវាទិន្នន័យ ឬកែវាទិន្នន័យ ដោយការបញ្ចប់ទិន្នន័យ
- Dimensionality reduction
- Numerosity reduction
- Data compression
- Data transformation and data discretization**
min-max
- Normalization កែវាទិន្នន័យ ឬកែវាទិន្នន័យ ដោយការបញ្ចប់ទិន្នន័យ
- Concept hierarchy generation

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?

Chapter 3: Data Preprocessing

□ Data Preprocessing: An Overview

□ Data Cleaning *ព័ត៌មានការស្រួល*

□ Data Integration *គម្រោងប្រើប្រាស់*

□ Data Reduction and Transformation *លក់ប្រើប្រាស់ ផ្តល់ប្រើប្រាស់*

□ Dimensionality Reduction

□ Summary

Data Cleaning

ពេលវេលាដឹងចាំរែ

- ❑ 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?

Incomplete (Missing) Data

ទម្រង់ ឬសម្របក្នុង

- 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 កែវិក ឬកែវិក
 - Record data ឬតម្លៃការ
 - 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 ឬការពន្លាឯក្រាម
- 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 *Data record មានតើ missing កំណត់ទេ*
- Fill in the missing value manually: tedious + infeasible?
នូវការ missing ត្រូវបាន
- Fill in it automatically with
 - a global constant : e.g., “unknown”, a new class?!
 - the attribute mean *នឹង mean, ដែលត្រូវបានពីរបៀប*
 - 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**