# CS 412 Intro. to Data Mining

## Chapter 3. Data Preprocessing

การจัดการ Data ก่อนที่จะไปประมวลผล

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

# Chapter 3: Data Preprocessing

- Data Preprocessing: An Overview

- Data Cleaning   Data ที่เก็บมามันหลายแหล่ง

  เก็บเอง

  sensor - เก็บจากโนสถ์ ม noise, missing เหมือนกัน

  ข้อมูลที่ไม่ได้กรอก

  เป็น ข้อมูลที่ไม่ เข้ากับพวก

- Data Integration   เอา Data จากหลายแหล่งมารวมกัน

  ลดจน. Data     ลดด้าน Dimension ลดด้าน ตัวเลข

- Data Reduction and Transformation

- Dimensionality Reduction

- Summary

# What is Data Preprocessing? — Major Tasks

- **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

# Why Preprocess the Data? — Data Quality Issues

- Measures for data quality: A multidimensional view

  - Accuracy: correct or wrong, accurate or not
    เพราะข้อมูลที่ใส่ผ่านมาเป็นพันๆ ข้อมูลที่ถูก/ผิด

  - Completeness: not recorded, unavailable, …

  - Consistency: some modified but some not, dangling, …

  - Timeliness: timely update? → อัพเดทตามกาลเวลา
    กรอกวันเวลานั้น

  - Believability: how trustable the data are correct?
    น่าเชื่อถือจนจริงไหม?

  - Interpretability: how easily the data can be understood?

# 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
  - <u>Incomplete</u>: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
    - e.g., *Occupation* = " " (missing data)
  - <u>Noisy</u>: containing noise, errors, or outliers
    - e.g., *Salary* = "–10" (an error)
  - <u>Inconsistent</u>: 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
  - <u>Intentional</u> (e.g., *disguised missing* data)
    - Jan. 1 as everyone's birthday?

# Incomplete (Missing) Data

*Data ไม่สมบูรณ์*

- 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    *เกิดจากการไม่กรอกข้อมูล จึงเกิด Missing*
  - Equipment malfunction
  - Inconsistent with other recorded data and thus deleted    *ไม่รอดคล้อง*
  - Data were not entered due to misunderstanding    *เข้าใจผิด*
  - Certain data may not be considered important at the time of entry
  - Did not register history or changes of the data    *ข้อมูลเปลี่ยนแปลง*
- Missing data may need to be inferred
  *บางข้อมูลอาจจะประกบค่ากันได้*

# 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 ไหนที่ Missing เอาก้อนนั้นออก*

- Fill in the missing value manually: tedious + infeasible?

- Fill in it automatically with  *ไม่รู้*

  - a global constant : e.g., "unknown", a new class?!

  - the attribute mean  *เอาค่า Mean มาแทน ค่า missing*

  - the attribute mean for all samples belonging to the same class: smarter  *เหมือน mean ที่ อยู่ class เดียวกัน*

  - **the most probable value: inference-based such as Bayesian formula or decision tree**