



# CS 412 Intro. to Data Mining

## Chapter 3. Data Preprocessing

معالجة البيانات

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



# Chapter 3: Data Preprocessing

## □ Data Preprocessing: An Overview

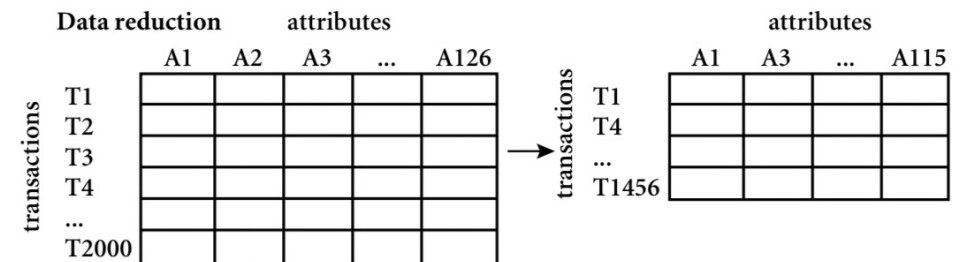
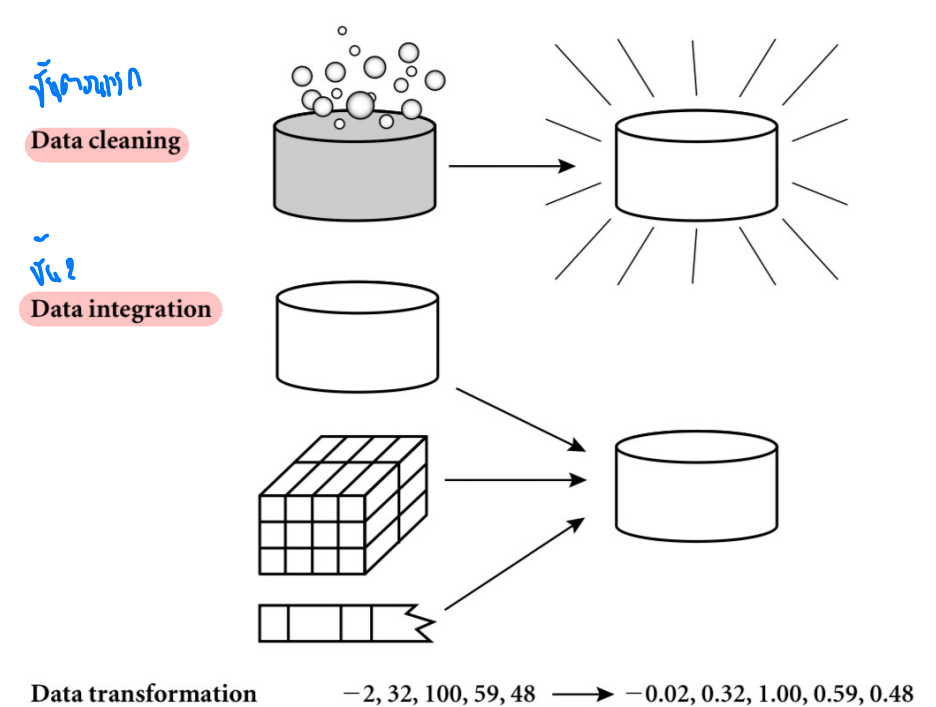
### 1 □ Data Cleaning *กำจัดข้อมูล Clean ที่มี missing*

### 2 □ Data Integration *นำ Data จากส่วนอื่น เพื่อใช้ในส่วนที่เราต้องการ*

### 3 □ Data Reduction and Transformation *ลด Data ให้เล็กลง* *นำ Data ไปปรับแก้*

## □ Dimensionality Reduction

## □ Summary



# What is Data Preprocessing? — Major Tasks

---

- ❑ **Data cleaning** *จัดการ missing มีทั้ง noisy และ outliers*
  - ❑ Handle missing data, smooth noisy data, identify or remove outliers, and resolve inconsistencies *จัดการ inconsistencies*
- ❑ **Data integration** *รวม Data ที่มาจากหลายๆ แหล่ง*
  - ❑ Integration of multiple databases, data cubes, or files
- ❑ **Data reduction** *การลดขนาดข้อมูล*
  - ❑ Dimensionality reduction
  - ❑ Numerosity reduction
  - ❑ Data compression
- ❑ **Data transformation and data discretization** *การแปลงข้อมูลให้อยู่ในลักษณะที่เหมาะสม*
  - ❑ 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
  - ❑ Completeness: not recorded, unavailable, ...
  - ❑ Consistency: some modified but some not, dangling, ...
  - ❑ Timeliness: timely update? *সময়মত ডাটা মনিটরিং*
  - ❑ 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

ดาตาสกปรกไม่สะอาด

- ❑ Data in the Real World Is Dirty: Lots of potentially incorrect data, e.g., instrument faulty, human or computer error, and transmission error เกิดจากที่มนุษย์ใส่ค่า ผิดพลาด หรือ หาไป
- ❑ Incomplete: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
  - ❑ e.g., *Occupation* = " " (missing data) ไม่กรอกข้อมูล
- ❑ Noisy: containing noise, errors, or outliers
  - ❑ e.g., *Salary* = "-10" (an error) ความผิด
- ❑ Inconsistent: containing discrepancies in codes or names, e.g.,
  - ❑ *Age* = "42", *Birthday* = "03/07/2010" ขัดแย้งกัน
  - ❑ 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

- Data is not always available *Data ไม่สมบูรณ์*
  - E.g., many tuples have no recorded value for several attributes, such as customer income in sales data
- Missing data may be due to *(เกิดจาก) 1. ข้อบกพร่อง*
  - Equipment malfunction *เครื่องเสีย*
  - Inconsistent with other recorded data and thus deleted *ไม่เข้ากันกับข้อมูลอื่น*
  - Data were not entered due to misunderstanding *ความเข้าใจผิดในการบันทึก*
  - 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?

- ❑ 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
  - ❑ a global constant : e.g., “unknown”, a new class?!
  - ❑ the attribute mean
  - ❑ 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**