

## Chapter 3 Data Preprocessing (40% คือบทนี้ ตั้งใจเรียนด้วยนะจ๊ะ)

### Chapter 3: Data Preprocessing

- ☐ Data Preprocessing: An Overview
- ☐ Data Cleaning
- ☐ Data Integration
- ☐ Data Reduction and Transformation
- ☐ Dimensionality Reduction
- ☐ Summary

The diagram shows four main steps of data preprocessing:

- Data cleaning:** Represented by a cylinder with bubbles being removed, indicating the removal of noise or outliers.
- Data integration:** Represented by two cylinders merging into one, indicating the combination of data from different sources.
- Data reduction:** Represented by a large grid of attributes (A1 to A126) being reduced to a smaller grid (A1 to A115), indicating the selection of relevant attributes.
- Data transformation:** Represented by a row of values [-2, 32, 100, 59, 48] being transformed into a normalized row [-0.02, 0.32, 1.00, 0.59, 0.48], indicating the scaling of data.

- Data Cleaning คือ ทำการ cleaning Data เนื่องจากเก็บข้อมูลมาหลายแหล่ง เช่น แบบฟอร์มให้คนอื่นกรอก แล้ว เกิดเป็น noise คือ กรอกข้อมูลผิด เป็นต้น
- Data Integration คือ การนำ Data จากหลายแหล่งมารวมกัน ก็ขบวนการรวม อาจจะรวม เป็นตารางเพื่อนำไปทำ Data Mining ต่อ หรือรวม เพื่อเป็น Data warehouse เพื่อเรียกดูข้อมูลแนวต่างๆได้
- Data Reduction and Transformation คือ การลดจำนวนข้อมูล / สอนว่าจะแปลงข้อมูลอะไรให้เหมาะสมลดได้
- Dimensionality Reduction คือ การ ลดจำนวนข้อมูลแนวตั้ง

## 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

Data cleaning จัดการ missing, inconsistencies  
/ กำจัด noisy, outliers

Data integration รวม Data จากหลายแหล่ง ไม่จำเป็น  
ว่ามาจาก database

Data reduction ลดจำนวนข้อมูล

Data transformation เปลี่ยนข้อมูลเพื่อให้เข้ากับเพื่อน ๆ

## 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?

ทำไมต้องทำ preprocess ?

เพราะ Data ของเราจากแหล่งต่างๆ

นั้นอาจมีลักษณะการทำให้ preprocessing

Data Cleaning : Incomplete (ข้อมูลหาย),

Noisy, Inconsistent (ข้อมูลขัดแย้งกัน),

Intentional

## 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
- ❑ Missing data may need to be inferred

เช่น เราทำหอยปี 1 กรองข้อมูล ที่อยู่ในแบบฟอร์ม ของว่า ชื่อนี้ เราทำปี 1 เรากรอกว่า  
ได้รูดเงินหรือยัง แบบนี้เรารู้ว่า Missing Data เพราะไม่ได้เขากรอกในช่องนี้ เช่นต้น

## 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

ถ้า Data ใดมี Missing ก็ให้ลบมันออกไป แต่ถ้าข้อมูลใดมีเป็นหลักร้อย ๆ อาจ  
ไม่ดีกว่าไหม ถ้าใช้วิธีนี้ แต่เราก็สามารถทำได้ ถ้าเราอยากทำ