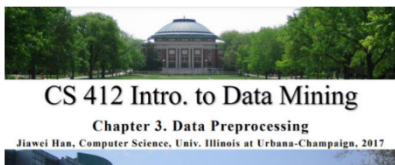


สรุป Chapter 3

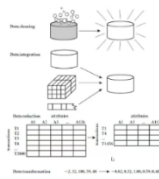


Data Preprocessing คือ การจัดการข้อมูลก่อนไปประมวลผล

Pre แปลว่า ก่อน

Chapter 3: Data Preprocessing

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



- ขั้นตอนก็จะอยู่ฝั่งซ้าย รูปร่างหน้าตามันก็อยู่ฝั่งขวา

- เริ่มมาก็จะทำการ Cleaning Data

(นอยด์ คือ ข้อมูลที่ไม่จริงข้อมูลที่กรอกผิด ไม่เข้ากับพวก) (Missing คือ ข้อมูลที่ไม่ได้กรอก ไม่ได้เก็บ หายไปพอดี)

- ขั้นที่ 2 Data Integration คือ การรวมข้อมูลมาจากหลายแหล่ง เช่น รวมเป็นตาราง เป็นต้น

- Data Reduction (การลดข้อมูลจำนวน) **เป็นเทคนิคโบราณ and Transformation

Transformation (แปลงข้อมูลยังไงให้ประมวลผลได้) -> จำนวนเป็น Data point แนวตั้งเป็น พิวเจอร์

- Dimensionality (การลดข้อมูลแนวตั้ง)

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
 - 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?

- ทำไมถึงต้องทำ Data Preprocess เพราะว่า ข้อมูลที่ใส่เข้ามา มีทั้ง ข้อมูลที่ผิดและข้อมูลที่ถูกต้องเป็นข้อมูลที่ไม่น่าเชื่อถือ มาจากหลายแหล่ง คนกรอกบ้าง เครื่องบ้าง

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
 - ❑ 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

- ❑ 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

