



# CS 412 Intro. to Data Mining

ມະຊຸມ Data ກອນທີ່ຈະນຳໄປຫຼາຍາພວ

## 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 ที่เก็บมาไม่ถูกต้อง -
  - เก็บ wrong
  - Sensor : เก็บข้อมูลจาก sensor
  - noise : ข้อมูลที่ไม่ใช่ปกติ (เสียงรบกวน)
  - missing : ข้อมูลที่ไม่ได้กรอก
- Data Integration
- 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
  - จัดการข้อมูล missing  
→ จัดการข้อมูล noisy  
→ จัดการข้อมูล outliers  
→ จัดการข้อมูล inconsistencies
- Data integration
  - Integration of multiple databases, data cubes, or files
- Data reduction → ลดจำนวนข้อมูล
  - Dimensionality reduction
  - Numerosity reduction
  - Data compression
  - ลดขนาดของข้อมูล
- 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  
☞ เปราะ: ข้อมูลที่ถูกต้องแม่นยำหรือไม่
  - Completeness: not recorded, unavailable, ...
  - Consistency: some modified but some not, dangling, ...
  - Timeliness: timely update? → อัปเดตตามกำหนดเวลา  
↳ ก่อให้เกิดข้อผิดพลาด
  - Believability: how trustable the data are correct?  
↳ น่าเชื่อถือหรือไม่?
  - Interpretability: how easily the data can be understood?

<#>

## Chapter 3: Data Preprocessing

- Data Preprocessing: An Overview
- Data Cleaning → Data ที่เก็บมาแล้วเหลือ
  - เก็บตรง
  - Sensor: เก็บอัตโนมัติ
  - noise: ข้อมูลที่ไม่ใช่ข้อมูล
  - missing: ข้อมูลที่ไม่ได้ก่อตัว
- Data Integration  
↳ ลดจำนวน Data
- Data Reduction and Transformation
- Dimensionality Reduction
  - ลด Data จำนวนมากเหลือจำนวนที่น้อยลง
  - ลด Dimension ลดจำนวนคุณลักษณะ
- Summary

<#>

# Data Cleaning

ເກົ່າງມື້ນີ້ຢັ້ງໃຈກວ

- Data in the Real World Is Dirty: Lots of potentially incorrect data, e.g., instrument faulty, human or computer error, and transmission error  
↳ ໂຄງຫຍານ  
• Incomplete: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data  
↳ ປູ້ໄດ້ກອກຂໍ້ມູນ  
• e.g., *Occupation* = “ ” (missing data)  
• Noisy: containing noise, errors, or outliers  
• e.g., *Salary* = “-10” (an error) → ອະນຸຍິດ  
• 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?

<#>

ພາບ Data ຍັງສະບຽບ

## 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 → ເກົ່າງຈາກທີ່ໄປກອກຂໍ້ມູນ ຈົ້າເກີດ Missing
  - Equipment malfunction  
↳ ໄດ້ສອດດັບ
  - Inconsistent with other recorded data and thus deleted
  - Data were not entered due to misunderstanding → ເນັ້າໃຈຜົດ
  - Certain data may not be considered important at the time of entry
  - Did not register history or changes of the data → ຂໍ້ມູນມີການປັບປຸງຂັ້ນແລ້ວ
- Missing data may need to be inferred



ບາງອັນຫາຈາກປະການຄໍາໄລ່

<#>

# 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 → หากพบว่า Data ตัวใดมี Missing จะทำการลบออก
- 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 → ใช้ค่า Mean ของแทนที่ Missing
  - 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**

→ แทน Mean ที่อยู่ใน class เดียวกัน

<#>

## Noisy Data

- Noise:** random error or variance in a measured variable
- Incorrect attribute values** may be due to
  - Faulty data collection instruments
  - Data entry problems
  - Data transmission problems
  - Technology limitation
  - Inconsistency in naming convention
- Other data problems**
  - Duplicate records
  - Incomplete data
  - Inconsistent data

<#>