



# CS 412 Intro. to Data Mining

## Chapter 3. Data Preprocessing

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





ការពារិភាពទូទៅ របៀបនាំរាយទូទៅ

## Chapter 3: Data Preprocessing

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

# What is Data Preprocessing? — Major Tasks

---

- **Data cleaning** *ប៉ុន្តែងតាមរយៈមធារដែលអ្នកត្រូវដោយសារពីការបញ្ចូលការ*
  - Handle missing data, smooth noisy data, identify or remove outliers, and resolve inconsistencies
- **Data integration** *រួម Data ដើម្បីទីផ្សារ / ទូរចាប់ផ្តើមទិន្នន័យដូចជាអេក្រង់*
  - Integration of multiple databases, data cubes, or files
- **Data reduction** *ការសកម្មភាពរួមមុន*
  - Dimensionality reduction
  - Numerosity reduction
  - Data compression
- **Data transformation and data discretization** *រំភេទអ្នកត្រូវដោយសារពីការបញ្ចូលការ*
  - 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, ... กານກົດຂູ້ລືໄສເກມຕໍ່ານ
  - Consistency: some modified but some not, dangling, ... ກຳ Normalization ໃນເຄື່ອງຈຸດທີ່ໃນປະເທດ
  - Timeliness: timely update? Data ສິ້ນການການໃຈຕໍ່ານ ໂດຍໄດ້ໂຈງຕັ້ງມາດທີ່ໄດ້ສຳເນົາຮູ້ໃຫຍ່ນອຸນຫວຼາງ
  - 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       *ការកំណត់ទិន្នន័យពីរបាយ*
- ❑ Data Integration
- ❑ Data Reduction and Transformation
- ❑ Dimensionality Reduction
- ❑ Summary

# Data Cleaning

ການເຄືອຂາຍສະຫອງທົບນູ້

- ❑ Data in the Real World Is Dirty: Lots of potentially incorrect data, e.g., instrument faulty, human or computer error, and transmission error *ເວັບໄດ້ຜິດ*
  - ❑ 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) *ຕ່າງໆຈາງເກີນ ເຮັນເລືອດກີ່ຕົງຄົນ*
      - ❑ Inconsistent: containing discrepancies in codes or names, e.g.,
        - ❑ *Age* = “42”, *Birthday* = “03/07/2010” *ທຸກໆກຳເຄົາໄຟເຕວກັນ*
        - ❑ 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?

ទុកដាក់នៃតម្លៃទូទៅ

តម្លៃទូទៅ

## Incomplete (Missing) Data

- Data is not always available ទុកដាក់នៃតម្លៃទូទៅ
  - 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 ឧបករណ៍ ការងារអិលម្មតី
  - 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

វិធានសំគាល់ការបង្កើតផ្តល់នូវទម្រង់

## 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      ព័ត៌មានមួយចំណាំ      សរុប
  - the attribute mean for all samples belonging to the same class: smarter ព័ត៌មានមួយចំណាំ សរុប ស្ថិតិយោគ
  - **the most probable value: inference-based such as Bayesian formula or decision tree**

# Noisy Data

ចំណាំទឹកសង្គម / ទីតាំងរាយការណ៍

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

# How to Handle Noisy Data?

---

- Binning *តាមលក្ខណៈនីមួយៗរបស់វា ចិត្តភាព និងចំណែកជាបន្ទូរ*
  - First sort data and partition into (equal-frequency) bins
  - Then one can **smooth by bin means, smooth by bin median, smooth by bin boundaries, etc.**
- Regression *សម្រាប់នូវការរាយការណ៍ទិន្នន័យ និងការពិនិត្យការងារ*
  - Smooth by fitting the data into regression functions
- Clustering
  - Detect and remove outliers
- Semi-supervised: Combined computer and human inspection
  - Detect suspicious values and check by human (e.g., deal with possible outliers)