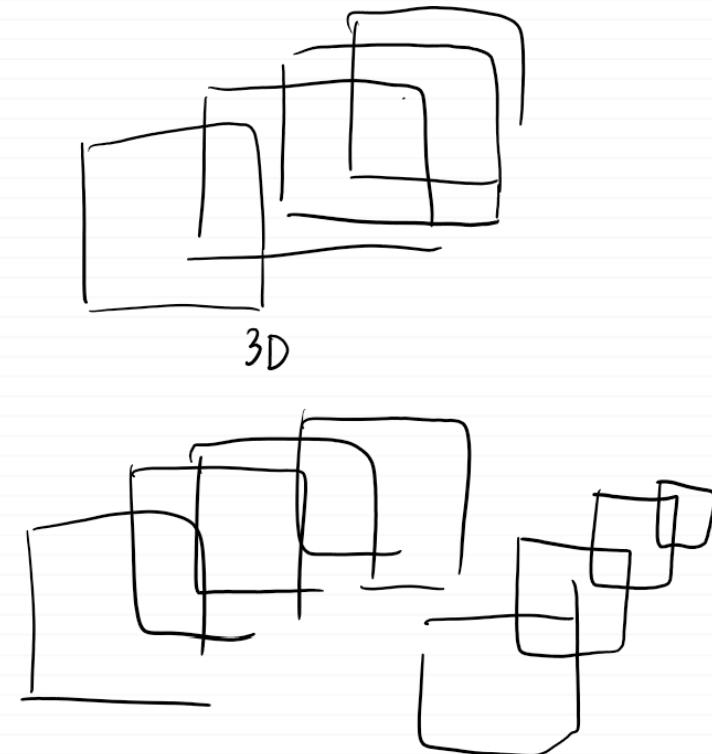
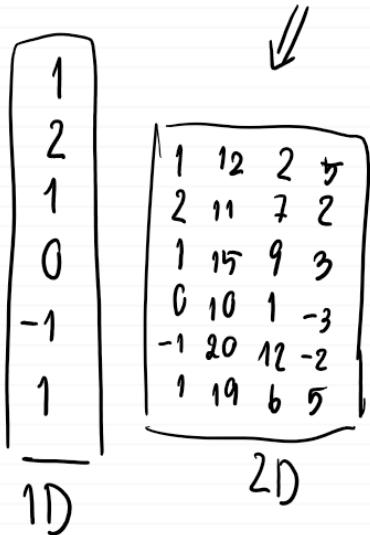


Data ជាអំពី  
រូប និង លក្ខណៈ



4D  
ផ្លូវតាម attribute នៃការ  
រួមចំណាំ

Record1

Record2

⋮

Record6

Attribute 1      Attribute 2      ...      Attribute n



**CS 412 Intro. to Data Mining**

## **Chapter 2. Getting to Know Your Data**

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



## **Chapter 2. Getting to Know Your Data**

- ❑ Data Objects and Attribute Types
  - ❑ Basic Statistical Descriptions of Data
  - ❑ Data Visualization
  - ❑ Measuring Data Similarity and Dissimilarity
  - ❑ Summary

## Cross-tab of Attribute

## **Types of Data Sets: (1) Record Data**

- ❑ Relational records
    - ❑ Relational tables, highly structured
  - ❑ Data matrix, e.g., numerical matrix, crosstabs

	China	England	France	Japan	USA	Total
Active Holders French Glass	12,000	1,000	1,000	1,000	1,000	19,000
Active Holders Japan Glass	3,000	6,000	6,000	3,000	3,000	21,000
SellBack French Glass	3,000	6,000	6,000	3,000	3,000	18,000
SellBack Japan Glass	3,000	6,000	6,000	3,000	3,000	18,000
Through Prod Direct	3,000	1,000	1,000	3,000	3,000	12,000
Through Verilog Network	3,000	22,000	22,000	47,000	47,000	120,000
Xpresso Adult Indirect	9,000	8,000	7,000	2,000	25,000	60,000
Xpresso Youth Indirect	1,000	1,000	1,000	1,000	1,000	5,000
Total	14,000	43,000	54,000	3,000	3,072,000	2,081,000

- #### ❑ Transaction data

<b>TID</b>	<b>Items</b>
<b>1</b>	Bread, Coke, Milk
<b>2</b>	Beer, Bread
<b>3</b>	Beer, Coke, Diaper, Milk
<b>4</b>	Beer, Bread, Diaper, Milk
<b>5</b>	Coke, Dianer, Milk

- ❑ Document data: Term-frequency vector (matrix) of text documents

	Name	Document 1	Document 2	Document 3	Document 4	Document 5	Document 6	Document 7	Document 8	Document 9	Document 10
1	John	3	0	5	0	2	6	0	2	0	2
2	Jane	0	7	0	2	1	0	0	3	0	0
3	Mike	0	1	0	0	1	2	2	0	3	0

Pers_ID	Surname	First_Name	City
0	Miller	Paul	London
1	Orenga	Alvaro	Valetina
2	Hutten	Urs	Zurich
3	Blanc	Gaston	Paris
4	Bonnotte	Eduardo	Bogot

— no relation

Car ID	Model	Year	Value	Pers. ID
101	Bentley	1973	100000	0
102	Rolls Royce	1965	330000	0
103	Peugeot	1993	500	3
104	Ferrari	2005	150000	4
105	Renault	1998	2000	3
106	Renault	2001	7000	3
108	Kia	2000	1000	0

} term-frequency matrix  $\Rightarrow$  Text

Normalization → ເສັ່ນໃຈກົດຕາວ່າ ສູງເປັນກົງ.



## Data Objects

- Data sets are made up of data objects → **data ที่เก็บรวบรวม = data set**
- A **data object** represents an entity **เจ้าของ**
- Examples:
  - sales database: customers, store items, sales **ขาย Sell เก็บ ขายต่อ**
  - medical database: patients, treatments
  - university database: students, professors, courses
- Also called *samples, examples, instances, data points, objects, tuples* **Data ตัวอย่าง**
- Data objects are described by **attributes** **เจ้าของ**
- Database rows → data objects; columns → attributes

- ## Attributes
- คุณลักษณะที่อยู่ในตัวอย่าง
- Attribute (or dimensions, features, variables)**
    - A data field, representing a characteristic or feature of a data object.
    - E.g., *customer\_ID, name, address*
  - Types:** **ชนิดของคุณลักษณะ**
    - Nominal (e.g., red, blue) **นามบัญญัติ / ชื่อของกลุ่ม/ชนิด** ⇒ ๑, ๐, ✕, ✗ **ไม่ต่อ**
    - Binary (e.g., {true, false}) **มีแค่ 2 ตัว**
    - Ordinal (e.g., {freshman, sophomore, junior, senior}) **ชั้นมัธยมต่อตัวอื่น**
    - Numeric: quantitative **ตัวเลขตัวตน**
      - Interval-scaled: 100°C is interval scales **บาก, กลาง, ด้านหลัง, ด้านหน้า** **ไม่มีความหมาย**
      - Ratio-scaled: 100°K is ratio scaled since it is twice as high as 50°K
  - Q1: Is student ID a nominal, ordinal, or interval-scaled data?
  - Q2: What about eye color? Or color in the color spectrum of physics?  
**สีฟ้า ฟ้าอมเขียว ฟ้าเขียว** **สีแดง - สีขาว**  
**Numeric ไม่ต่อ**

35%

## Attribute Types

- Nominal:** categories, states, or “names of things” **บุคลิก, สถานะ, ชื่อของสิ่ง**
  - Hair\_color = {auburn, black, blond, brown, grey, red, white}* **สีผม**
  - marital status, occupation, ID numbers, zip codes
- Binary** **binary**: **ตัวสอง**
  - Nominal attribute with only 2 states (0 and 1) **มี 2 ตัว**
  - Symmetric binary: both outcomes equally important ⇒ **สมมาตรการกัน (ความสำคัญเท่ากัน)**
    - e.g., **gender** (**เพศเด็ก/เด็ก, ผู้ชาย/ผู้หญิง**) / (**เด็ก/ผู้ใหญ่**) / (**เด็ก/เด็ก**)
  - Asymmetric binary: outcomes not equally important. ⇒ **ไม่สมมาตรการกัน (ความสำคัญต่างกัน)**
    - e.g., medical test (positive vs. negative)
    - Convention: assign 1 to most important outcome (e.g., HIV positive)
- Ordinal** **เรียงลำดับได้** **ไม่สามารถรู้ว่า มากกว่า น้อยกว่า**
  - Values have a meaningful order (ranking) but magnitude between successive values is not known
  - Size = {small, medium, large}, grades, army rankings*

## Numeric Attribute Types

- Quantity** (integer or real-valued)
- Interval** → **ไม่มี ๐ หรือ = ไม่มีความหมาย**
  - Measured on a scale of **equal-sized units**
  - Values have order
    - E.g., temperature in C° or F°, calendar dates
  - No true zero-point
- Ratio** → **มี ๐ หรือ = ไม่มี** **ไม่มี ๐ หรือ < ไม่มี ค่านอก**
  - Inherent **zero-point**
  - We can speak of values as being an order of magnitude larger than the unit of measurement (10 K° is twice as high as 5 K°).
    - e.g., temperature in Kelvin, length, counts, monetary quantities

## Discrete vs. Continuous Attributes

### Discrete Attribute

- Has only a finite or countably infinite set of values
- E.g., zip codes, profession, or the set of words in a collection of documents
- Sometimes, represented as integer variables
- Note: Binary attributes are a special case of discrete attributes

### Continuous Attribute

- Has real numbers as attribute values
- E.g., temperature, height, or weight Ex อุณหภูมิ, ส่วนสูง, น้ำหนัก
- Practically, real values can only be measured and represented using a finite number of digits
- Continuous attributes are typically represented as floating-point variables

13

## Chapter 2. Getting to Know Your Data

### Data Objects and Attribute Types

### Basic Statistical Descriptions of Data

### Data Visualization

### Measuring Data Similarity and Dissimilarity

### Summary



14



## Basic Statistical Descriptions of Data



### Motivation

- To better understand the data: central tendency, variation and spread

### Data dispersion characteristics

- Median, max, min, quantiles, outliers, variance, ...

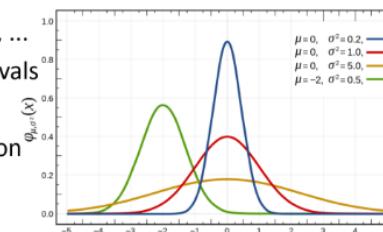
### Numerical dimensions correspond to sorted intervals

- Data dispersion:

- Analyzed with multiple granularities of precision
- Boxplot or quantile analysis on sorted intervals

### Dispersion analysis on computed measures

- Folding measures into numerical dimensions
- Boxplot or quantile analysis on the transformed cube



15

## Measuring the Central Tendency: (1) Mean

### Mean (algebraic measure) (sample vs. population):

Note:  $n$  is sample size and  $N$  is population size.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \mu = \frac{\sum x}{N}$$

### Weighted arithmetic mean:

$$\bar{x} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

### Trimmed mean:

- Chopping extreme values (e.g., Olympics gymnastics score computation)

16

## Measuring the Central Tendency: (2) Median

- Median:
  - Middle value if odd number of values, or average of the middle two values otherwise
- Estimated by interpolation (for *grouped data*):

age	frequency
1–5	200
6–15	450
16–20	300
21–50	1500
51–80	700
81–110	44

Approximate median

$$\text{median} = L_1 + \left( \frac{n/2 - (\sum \text{freq})_l}{\text{freq}_{\text{median}}} \right) \text{width}$$

Low interval limit

Sum before the median interval

Interval width ( $L_2 - L_1$ )

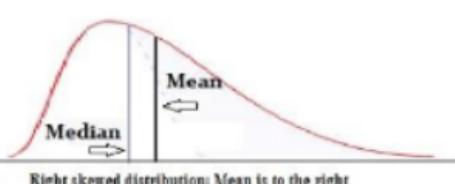
## Measuring the Central Tendency: (3) Mode

- Mode: Value that occurs most frequently in the data

- Unimodal

- Empirical formula:

$$\text{mean} - \text{mode} = 3 \times (\text{mean} - \text{median})$$



- Multi-modal

- Bimodal

- Trimodal

