

## Getting to Know Your Data

Data คืออะไร ที่น่าสนใจ จัดการอย่างไรและประโยชน์อะไรกับข้อมูลได้บ้าง !!

nu. Data

1
2
1
0
-1
1

10

1	12	2	5
2	11	7	2
1	15	9	3
0	10	1	-3
-1	20	12	-2
1	19	6	-5

20

3D

1	12	2	5
2	11	7	2
1	15	9	3
0	10	1	-3
-1	20	12	-2
1	19	6	-5

40

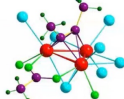
	Attribute 1	Attribute 2	Attribute 3	Attribute 4
Record 1				
⋮				
Record 6				

## 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

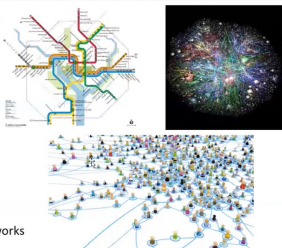
## Types of Data Sets: (2) Graphs and Networks

- ❑ Transportation network
- ❑ World Wide Web



- Molecular Structures
- Social or information networks

กราฟช่องมือ



## 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 Outdoors (cricket) China	12.00	6.00	1.00	240.00	297.00	
Active Outdoors (golf) China	12.00	6.00		333.00	329.00	
Active (cricket) China	3.00	6.00	0.00	132.00	149.00	
Active (golf) China		2.00		143.00	145.00	
Enough Pvc (indoor)	3.00	1.00	7.00	333.00	244.00	
Enough Sports (indoor)		3.00	31.00	674.00	939.00	
Enough Adult (indoor)	8.00	8.00	7.00	2.00	254.00	
Enough Youth (indoor)			1.00		74.00	77.00
Total	54.00	43.00	54.00	3.00	1,977.00	2,064.00

- Transaction data

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

	Team	Coach	Player	Ball	Score	Game	Win	Goal	Time	Season
Document 1	3	0	5	6	2	6	0	2	9	2
Document 2	9	7	6	2	1	0	0	3	9	0
Document 3	0	1	0	6	1	2	2	8	3	

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

**Person**

	First Name	City
0	Miller	London
1	Orange	Valencia
2	Blanc	Zurich
3	Blanc	Paris
4	Bertolini	Rome

no relation

**Car:**

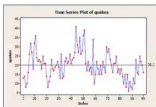
Car_ID	Model	Year	Value	Pers_ID
101	Bentley	1973	100000	0
102	Rolls Royce	1965	230000	0
103	Pagewood	1993	500	3
104	Ferrari	2005	150000	4
105	Renault	1998	2000	3
106	Renault	2001	7000	3
107	Smart	1999	2000	2

Normalization → เพื่อไม่ให้ตารางซ้ำกัน

### Types of Data Sets: (3) Ordered Data

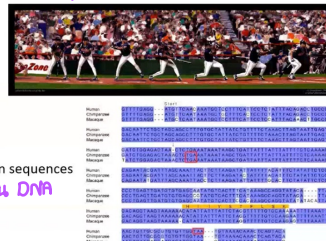
- Video data: sequence of images

- Temporal data: time-series



- ❑ Sequential Data: transaction sequences  
สลับลำดับไม่ได้ เช่น DNA

- Genetic sequence data



## Types of Data Sets: (4) Spatial, image and multimedia Data

ข้อมูลเชิงพื้นที่

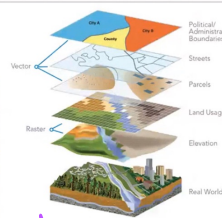
- Spatial data: maps

- Image data:



- Video data: Spatio-temporal

ข้อมูลเชิงวิดีโอ



ข้อมูลที่มีเวลาเข้ามาเกี่ยวข้อง

## Important Characteristics of Structured Data

คุณสมบัติที่สำคัญ

- Dimensionality → มิติ
  - Curse of dimensionality
- Sparsity → สันใจตารางที่มีข้อมูลเท่านั้น!! เพราะบางคอลัมน์มี 0 เยอะ/ช่องว่าง
  - Only presence counts
- Resolution → ความละเอียดในตาราง เช่น ภาพ กว้าง ขาว ก็ไม่ชัด
  - Patterns depend on the scale
- Distribution → การกระจายของข้อมูล
  - Centrality and dispersion

## Data Objects

ข้อมูลที่ได้อะไร

ประกอบด้วย Data Objects หลายๆอันมารวมกัน

- Data sets are made up of data objects
- A **data object** represents an entity ข้อมูลแต่ละตัว
- Examples:
  - sales database: customers, store items, sales
  - medical database: patients, treatments
  - university database: students, professors, courses

\*

- Also called **samples, examples, instances, data points, objects, tuples** ข้อมูล 1 ชุด
- Data objects are described by **attributes** ข้อมูลถูกอธิบายด้วย attributes
- Database rows → data objects; columns → attributes

หมายเหตุ

## Attributes

คุณสมบัติที่อธิบายข้อมูลแต่ละตัว

ชื่อเรียก

- Attribute** (or dimensions, features, variables)
  - A data field, representing a characteristic or feature of a data object.
  - E.g., *customer\_ID, name, address*
- Types: ชนิด
  - Nominal (e.g., red, blue)
  - Binary (e.g., {true, false}) มีแค่ 2 ค่า
  - Ordinal (e.g., {freshman, sophomore, junior, senior}) ข้อมูลที่เรียงลำดับ
  - Numeric: quantitative ข้อมูลตัวเลข บางทีก็ถูกแปลงแล้วใช้คำนวณ
    - Interval-scaled: 100°C is interval scales +, -, x, ÷ ได้
    - 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?

## ชนิดของ Attribute Types

- **Nominal:** categories, states, or "names of things" *หมวดกลุ่ม, สถานะ, ชื่อ*
  - *Hair\_color* = {auburn, black, blond, brown, grey, red, white}
  - marital status, occupation, ID numbers, zip codes
- **Binary** *สถานะ*
  - Nominal attribute with only 2 states (0 and 1)
  - **Symmetric binary:** both outcomes equally important → *สมมาตรกัน (คส. เท่ากัน)*
    - e.g., gender *(ชาย/หญิง / โคก) , (ร้อน/เย็น)*
  - **Asymmetric binary:** outcomes not equally important. → *ไม่สมมาตรกัน (คส. ไม่เท่ากัน)*
    - e.g., medical test (positive vs. negative)
    - Convention: assign 1 to most important outcome (e.g., HIV positive)
- **Ordinal** *เรียงลำดับได้* *ไม่ลำดับการบ่งชี้ทุกหน่วยนอกกรอบที่ใด*
  - Values have a meaningful order (ranking) but magnitude between successive values is not known
  - Size = {small, medium, large}, grades, army rankings *ตำแหน่ง/ยศ*


## Discrete vs. Continuous Attributes

- **Discrete Attribute** *ไม่ต่อเนื่อง*
  - 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** *มีค่าพรอบ*
  - Has real numbers as attribute values
    - E.g., temperature, height, or weight *(อุณหภูมิ หรือ 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

## Numeric Attribute Types

- Quantity (integer or real-valued)
- **Interval** *0 ไม่แท้ คือ มีค่ามณฑล*
  - 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** *0 แท้ คือ ไม่มีความหมาย*
  - Inherent **zero-point** *มีต้นสอ 0 แท้ = ไม่ได้นิสอ*
  - 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°). *0 K - ไม่มีความหมาย*
  - e.g., temperature in Kelvin, length, counts, monetary quantities *0 \$ = ไม่มีความหมาย*

## 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

## 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

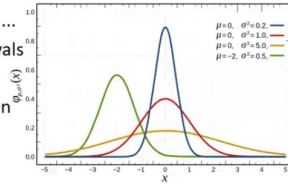
- Data dispersion:

- 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



## 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)