# CS 412 Intro. to Data Mining

## Chapter 2. Getting to Know Your Data

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

# Data

## Data

| 1D |
|---|
| 1 |
| 2 |
| 1 |
| 0 |
| -1 |
| 1 |

**2D**

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

**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 |

3

|  | Attribute 1 | Attribute 2 | Attribute 3 | Attribute 4 |
|---|---|---|---|---|
| Record 1 | 1 | 12 | 2 | 5 |
| Record 2 | 2 | 11 | 7 | 2 |
| Record 3 | 1 | 15 | 9 | 3 |
| Record 4 | 0 | 10 | 1 | -3 |
| Record 5 | -1 | 20 | 12 | -2 |
| Record 6 | 1 | 19 | 6 | -5 |

# 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: (1) Record Data

หลายตารางที่มีความสัมพันธ์กัน

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

Ex. ตารางสต็อกของ

|  | China | England | France | Japan | USA | Total |
|---|---|---|---|---|---|---|
| **Active Outdoors Crochet Glove** |  | 12.00 | 4.00 | 1.00 | 240.00 | 257.00 |
| **Active Outdoors Lycra Glove** |  | 10.00 | 6.00 |  | 323.00 | 339.00 |
| **InFlux Crochet Glove** | 3.00 | 6.00 | 8.00 |  | 132.00 | 149.00 |
| **InFlux Lycra Glove** |  | 2.00 |  |  | 143.00 | 145.00 |
| **Triumph Pro Helmet** | 3.00 | 1.00 | 7.00 |  | 333.00 | 344.00 |
| **Triumph Vertigo Helmet** |  | 3.00 | 22.00 |  | 474.00 | 499.00 |
| **Xtreme Adult Helmet** | 8.00 | 8.00 | 7.00 | 2.00 | 251.00 | 276.00 |
| **Xtreme Youth Helmet** |  | 1.00 |  |  | 76.00 | 77.00 |
| **Total** | 14.00 | 43.00 | 54.00 | 3.00 | 1,972.00 | 2,086.00 |

Person:

| Pers_ID | Surname | First_Name | City |
|---|---|---|---|
| 0 | Miller | Paul | London |
| 1 | Ortega | Alvaro | Valencia |
| 2 | Huber | Urs | Zurich |
| 3 | Blanc | Gaston | Paris |
| 4 | Bertolini | Fabrizio | Rom |

— no relation

Car:

| 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 |
| 107 | Smart | 1999 | 2000 | 2 |

- **Transaction data**

| TID | Items |
|---|---|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |

คำ →

|  | team | coach | play | ball | score | game | win | lost | timeout | season |
|---|---|---|---|---|---|---|---|---|---|---|
| **Document 1** (บทความ) | 3 | 0 | 5 | 0 | 2 | 6 | 0 | 2 | 0 | 2 |
| **Document 2** | 0 | 7 | 0 | 2 | 1 | 0 | 0 | 3 | 0 | 0 |
| **Document 3** | 0 | 1 | 0 | 0 | 1 | 2 | 2 | 0 | 3 | 0 |

→ จำนวนคำ

→ Data ใช้แทน ข้อความ

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

4

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

กราฟ

- ❑ Transportation network

- ❑ World Wide Web



- ❑ Molecular Structures

- ❑ Social or information networks

5

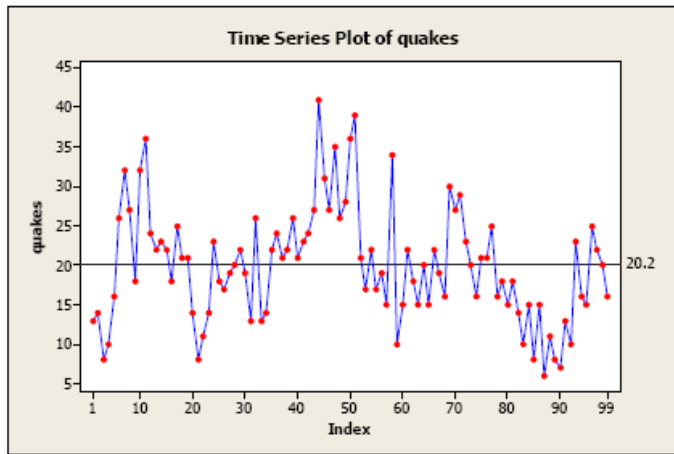# Types of Data Sets: (3) Ordered Data

มีเวลามาเกี่ยวข้อง, time series, หุ้น, มีลำดับมีตามสำคัญ

ข้อมูลวิดีโอมีเวลา เข้ามาเกี่ยวข้อง, นำรูปมาซ้อนๆกัน (การ์ตูนมิกกี้เมาส์ในอดีต)

❑ Video data: sequence of images

❑ Temporal data: time-series



**Time Series Plot of quakes**

❑ Sequential Data: transaction sequences



❑ Genetic sequence data

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

ข้อมูลเชิงพื้นที่มี 2 ด้าน กว้างกับยาว x,y

- ❑ Spatial data: maps



- ❑ Image data:

- ❑ Video data: spatio-temporal



Political/Administrative Boundaries

Streets

Parcels

Land Usage

Elevation

Real World

Vector

Raster

City A

City B

County

# Important Characteristics of Structured Data

❑ Dimensionality → มีไดเมนชันกี่ไดเมนชันก์ 2,3,4,5

  ❑ Curse of dimensionality

❑ Sparsity → สนใจตรวที่มีข้อมูล, ข้อมูลเป็น 0 เยอะไม่สนใจ

  ❑ Only presence counts

❑ Resolution → เก็บข้อมูลใส่สีเป็นจุดๆ

  ❑ Patterns depend on the scale

❑ Distribution → จิตต่ำกลางที่สีนี้สว่างหรือมืด

  ❑ Centrality and dispersion

Sparsity

Ex. ตารางสต๊อกของ

| | China | England | France | Japan | USA | Total |
|---|---|---|---|---|---|---|
| Active Outdoors Crochet Glove | | 12.00 | 4.00 | 1.00 | 240.00 | 257.00 |
| Active Outdoors Lycra Glove | | 10.00 | 6.00 | | 323.00 | 339.00 |
| InFlux Crochet Glove | 3.00 | 6.00 | 8.00 | | 132.00 | 149.00 |
| InFlux Lycra Glove | | 2.00 | | | 143.00 | 145.00 |
| Triumph Pro Helmet | 3.00 | 1.00 | 7.00 | | 333.00 | 344.00 |
| Triumph Vertigo Helmet | | 3.00 | 22.00 | | 474.00 | 499.00 |
| Xtreme Adult Helmet | 8.00 | 8.00 | 7.00 | 2.00 | 251.00 | 276.00 |
| Xtreme Youth Helmet | | 1.00 | | | 76.00 | 77.00 |
| Total | 14.00 | 43.00 | 54.00 | 3.00 | 1,972.00 | 2,086.00 |

Land Usage

Raster

Elevation

8

# Data Objects

□ Data sets are made up of data objects → กลุ่มของข้อมูลประกอบด้วยหลายๆ data

→ ข้อมูลแต่ละตัว

□ A **data object** represents an entity

□ Examples:

  □ sales database:  customers, store items, sales

  □ medical database: patients, treatments

  □ university database: students, professors, courses

คือ data

□ Also called *samples , examples, instances, data points, objects, tuples*

□ Data objects are described by **attributes** → ข้อมูลจะถูกอธิบายเรียก attributes

□ Database rows → data objects; columns → attributes

# Attributes

คุณสมบัติที่ใช้เรียกข้อมูลแต่ละตัว

- **Attribute (**or **dimensions, features, variables**)
  - A data field, representing a characteristic or feature of a data object.
  - *E.g., customer _ID, name, address*
- Types:
  - Nominal (e.g., red, blue) → ชื่อของกลุ่มชื่อของชนิด ไม่ใช่ตัวเลข
  - Binary (e.g., {true, false}) → ข้อมูลมีแค่สองค่า
  - Ordinal (e.g., {freshman, sophomore, junior, senior}) → ข้อมูลเรียงลำดับ
  - Numeric: quantitative → $+, -, \times, \div$ ได้แล้วมีความหมาย
    - Interval-scaled: $100^{\circ}$C is interval scales
    - Ratio-scaled: $100^{\circ}$K is ratio scaled since it is twice as high as $50^{\circ}$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?

Numeric

10

# Attribute Types

*→ ชนิดของ attribute*

- **Nominal:** categories, states, or "names of things"
    - *Hair_color = {auburn, black, blond, brown, grey, red, white}*  *สีผม*
    - marital status, occupation, ID numbers, zip codes  *สถานะ*
- **Binary** *→ เหมือน Nominal แต่มีแค่ 2 เช่น 0 กับ 1 , ใช่ กับ ไม่ใช่*

    - Nominal attribute with only 2 states (0 and 1)
    - Symmetric binary: both outcomes equally important  *สมมาตร*

        - e.g., ~~gender~~ *, Left / Right-handed, Coke / Pepsi, Hot/Cold*
    - Asymmetric binary: outcomes not equally important.  *→ มี 2 ค่าความสำคัญไม่เท่ากัน เช่น ตรวจ covid สัดส่วนเป็น*

        - e.g., medical test (positive vs. negative)  *2 คน ไม่เป็น 1000 คน*
        - Convention: assign 1 to most important outcome (e.g., HIV positive)  *ในมหาลัย*
- **Ordinal** *→ ไม่สามารถนำมารวมกันได้ , สามารถเรียงลำดับได้*

    - Values have a meaningful order (ranking) but magnitude between successive values is not known
    - *Size = {small, medium, large},* grades, army rankings

# Numeric Attribute Types

*ข้อมูลที่เป็นตัวเลข*

- Quantity (integer or real-valued)

- **Interval**    *0 แท้ เรามีดินสอ 0 แท่ง คือ 0 แท้*    *→ มีศูนย์แท้กับไม่แท้*

    - 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**

    - 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

_↗ ไม่ต่อเนื่อง_

- **Discrete Attribute** → _ระหว่างค่า 2 ค่า ไม่มีค่าที่อยู่ตรงกลาง เช่น แมว พยาบาล ไม่มีค่าอะไรอยู่ตรงกลาง_

    - 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** → _มีค่ากลางระหว่าง 2 ค่า เช่น สูง 180 กับ สูง 181 แต่มีคนสูง 180.5_

    - Has real numbers as attribute values

        - E.g., temperature, height, or weight

    - 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

# **Basic Statistical Descriptions of Data**

- ❑ <u>Motivation</u>

    เชี่ยวขนจากต่ำลางมากเท่าไหร่

  - ❑ To better understand the data: central tendency, <u>variation</u> and spread

- ❑ <u>Data dispersion characteristics</u>

    แตกต่าวจากต่ำลางมากน้อยแค่ไบน

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

- ❑ <u>Numerical dimensions</u> correspond to sorted intervals

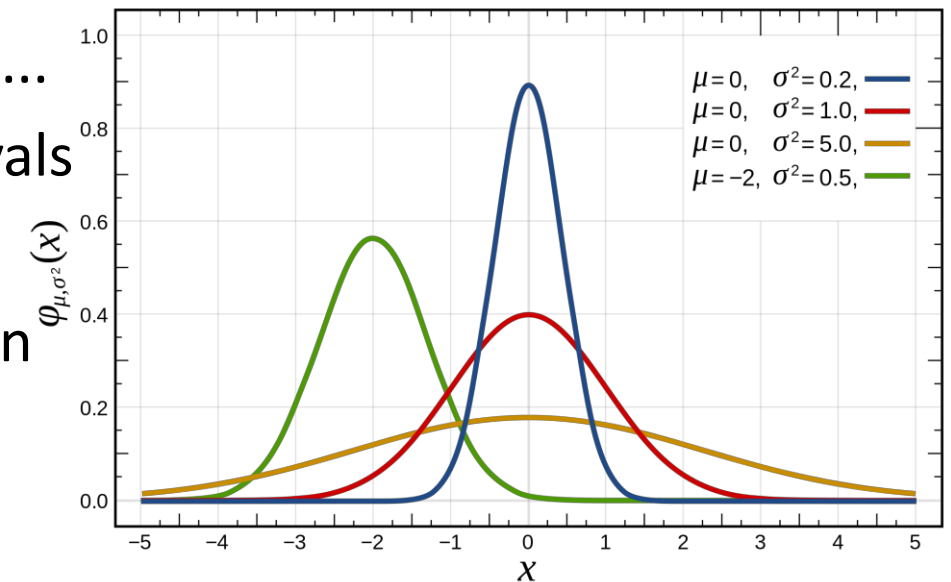  - ❑ Data dispersion:

    - ❑ Analyzed with multiple granularities of precision

  - ❑ Boxplot or quantile analysis on sorted intervals

- ❑ <u>Dispersion analysis on computed measures</u>

  - ❑ Folding measures into numerical dimensions

  - ❑ Boxplot or quantile analysis on the transformed cube



คนส่วนใหญ่ในน้ัวน้อายุ 20 ปี
– ฐานนิยม (ซ้ำเยอะ)