



# **CS 412 Intro. to Data Mining**

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

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









2 rows Data

# Data

1
2
1
0
-1
1

1D

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

2D

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

3D

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

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

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

4D

# Data

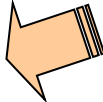
	Attribute 1	Attribute 2	Attribute 3	Attribute 4
Record 1	1	12	2	5
Record 2	2	11	7	2 <sup>I</sup>
Record 3	1	15	9	3
Record 4	0	10	1	-3
Record 5	-1	20	12	-2
Record 6	1	19	6	-5

[[Record]] คือ Data แต่ละจุด

[[Attribute]] คือ คุณสมบัติที่ใช้อธิบาย Data แต่ละจุด

# 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

אנשים ורכבים

- Relational records רשומות רצופות
- Relational tables, highly structured טבלאות רצופות, גבוהות מאוד
- Data matrix, e.g., numerical matrix, crosstabs מטריצה של נתונים, למשל

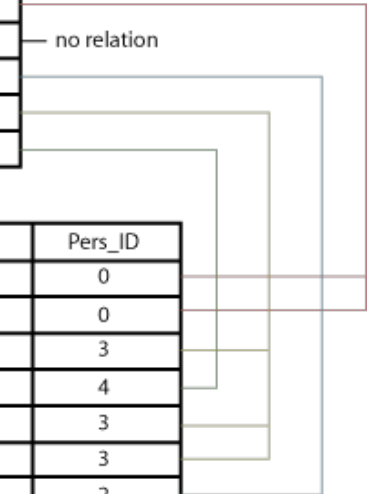
	China	England	France	Japan	USA	Total
Active Outdoors Crochet Glove		12.00	4.00	1.00	240.00	257.00
Active Outdoors Lycra Glove		10.00	6.00		323.00	339.00
InFlux Crochet Glove	3.00	6.00	8.00		132.00	149.00
InFlux Lycra Glove		2.00			143.00	145.00
Triumph Pro Helmet	3.00	1.00	7.00		333.00	344.00
Triumph 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

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 רשומות עסקאות

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

	team	coach	pla y	ball	score	game	wi n	lost	timeout	season
Document 1	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

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

מטריצה של תדירות מילים

למסמכים

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

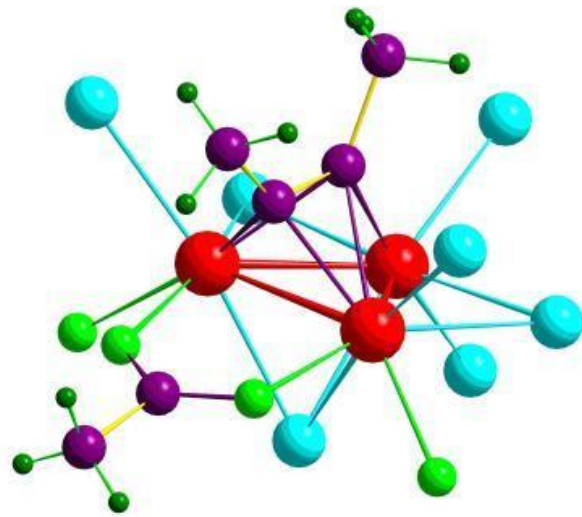
กราฟ

เครือข่าย

☐ Transportation network

เครือข่ายขนส่ง

☐ World Wide Web

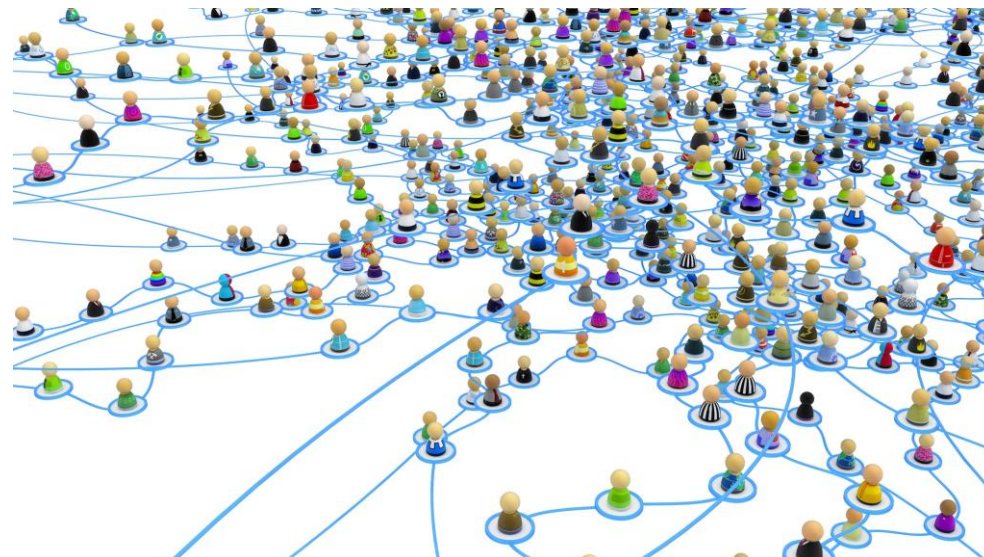
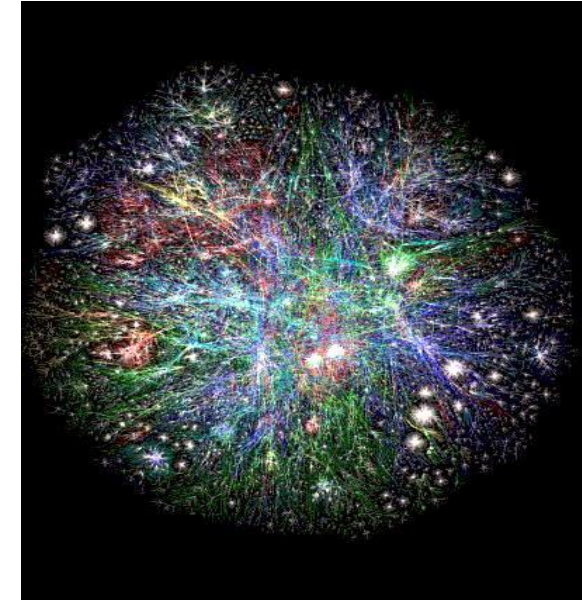
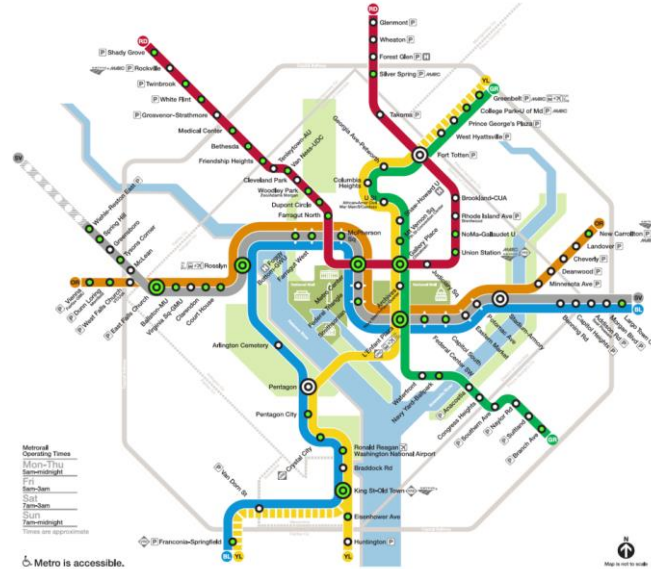


☐ Molecular Structures

โครงสร้างโมเลกุล

☐ Social or information networks

เครือข่ายสังคม

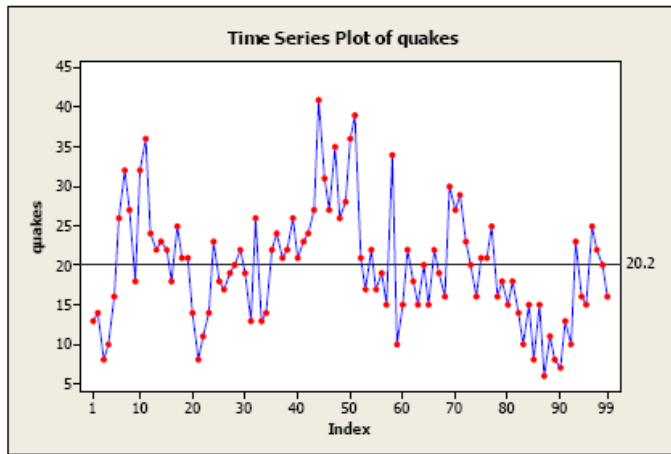




# Types of Data Sets: (3) Ordered Data

□ Video data: sequence of images (video data: sequence of images)

□ Temporal data: time-series (time-series)



□ Sequential Data: transaction sequences (transaction sequences)

□ Genetic sequence data (genetic sequence data)

	Start
Human	GTTTGGAGG --- ATGTTCAACAAATGCTCCTTTTCATTCTCTATTTACAGACCTGCCGCA
Chimpanzee	GTTTGGAGG --- ATGTTCAATAAATGCTGCTTTCACTCCTCTATTTACAGACCTGCCGCA
Macaque	GTTTGGAGG --- ATGCTCAATAAATGCTCCTTTTCATTCTCTATTTACAAACTTGCCGCA
Human	GACAATTCTGCTAGCAGCCTTTGTGCTATTATCTGTTTTCTAAACTTAGTAATTGAGTGT
Chimpanzee	GACAATTCTGCTAGCAGCCTTTGTGCTATTATCTGTTTTCTAAACTTAGTAATTGAGTGT
Macaque	GACAATTCTGCTAGCAGCCTTTGTGCTATTATCTGTTTTCTAAACTTAGTAATTGAGTGT
Human	GATCTGGAGACTAA - CTCTGAAATAAATAAGCTGATTATTTATTTATTTTCTCAAAACAA
Chimpanzee	GATCTGGAGACTAAACTCTGAAATAAATAAGCTGATTATTTATTTATTTTCTCAAAACAA
Macaque	TATCTGGAGACTAAACTCTGAAATAAATAAGCTGATTATTTATTTATTTTCTCAAAACAA
Human	CAGAATACGATTTAGCAAAATTACTTCTTAAGATATTATTTTACATTTCTATATTCTCCTA
Chimpanzee	CAGAATACGATTTAGCAAAATTACTTCTTAAGATATTATTTTACATTTCTATATTCTCCTA
Macaque	CAGAATATGATTTAGCAAAATTACTTCTTAAGATATTATTTTGCATTTCTATATTCTCCTA
Human	CCCTGAGTTGATGTGTGAGCAATATGTCACCTTTTCATAAAGCCAGGTATACA --- TTATG
Chimpanzee	CCCTGAGTTGATGTGTGAGCCGATGTCACCTTTTCATAAAGCCAGGTATACA --- TTATG
Macaque	CCCTGAGTTGATGTGTGAGCAATATGTCACCTTCACAAAGCCAGGTATATATACATTACG
Human	GACAGGTAAGTAAAAACATATTATTTATTCTACGTTTTTGCCAAAAATTTTAAATTTTC
Chimpanzee	GACAGGTAAGTAAAAACATATTATTTATTCTACGTTTTTGCCAAAGATTTTAAATTTTC
Macaque	GACAGGTAAGTAAAAA - CATATTATTTATTCTAGGTTTTTGCCAAAGATTTTAAATTTTC
Human	AACGTGTGCGCGTGTGTTGGTAA --- TGTAAAACAAAC TCAGTACA
Chimpanzee	AACGTGTGCGCGTGTGTTGGTAA --- TGTAAAACAAAC TCAGTACA
Macaque	AACGTGTGTGCATGTGTTGGTAA --- CBTAAAACAAATTCAGTACG



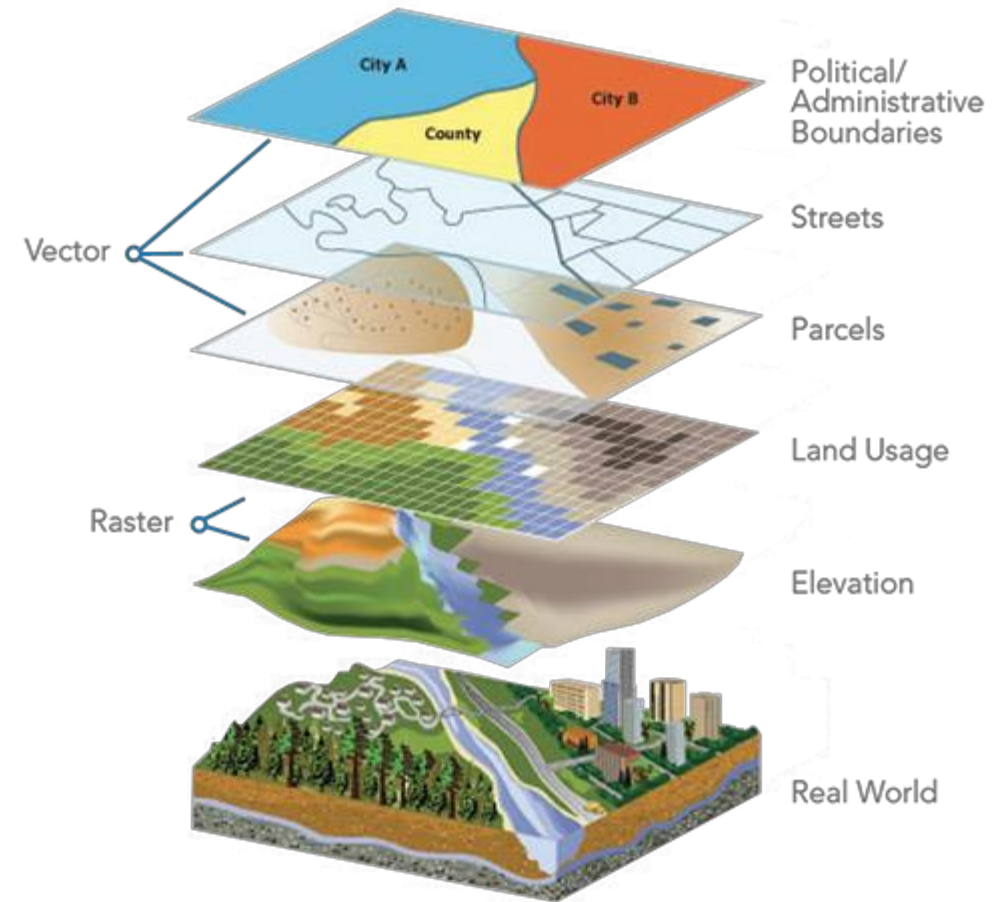
# Types of Data Sets: (4) Spatial, image and multimedia Data

□ Spatial data: maps



□ Image data:

□ Video data: spatio-temporal  
เชิงพื้นที่ เชิงเวลา



# Important Characteristics of Structured Data

- Dimensionality

Data with Dimensionality is high

- Curse of dimensionality

- Sparsity

Sparsity is high

- Only presence counts

- Resolution

- Patterns depend on the scale

- Distribution

- Centrality and dispersion

	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



# Data Objects

- <sup>အချက်အလက်အစု</sup> Data sets are made up of data objects
- A **data object** represents an entity
- Examples:
  - <sup>ရောင်းချမှု</sup> sales database: customers, store items, sales <sup>အချက်အလက်များ စာရင်းများ</sup>
  - medical database: patients, treatments
  - university database: students, professors, courses



Also called *samples, examples, instances, data points, objects, tuples* <sup>ဒီလို Data ၁ခုစီ (အချက်အလက်)</sup>

- Data objects are described by **attributes** <sup>အချက်အလက်ကိုဖော်ပြသော attributes</sup>
- Database rows → data objects; columns → attributes  
<sup>row → object                      column → attributes</sup>

# Attributes

- 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)
  - Binary (e.g., {true, false})
  - **Ordinal** (e.g., {freshman, sophomore, junior, senior})
  - Numeric: quantitative
    - **Interval-scaled**: 100°C is interval scales
    - **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?

หรือ

↓ คือคุณลักษณะที่ใช้อธิบายข้อมูล (ตัว)

หรือชื่อที่เฉพาะเจาะจง

หรือลำดับขั้น

ไม่ 0 หรือ

หรือ 0 หรือ

Nominal



# Attribute Types

- **Nominal:** categories, states, or “names of things”
  - *Hair\_color* = {auburn, black, blond, brown, grey, red, white}
  - marital status, occupation, ID numbers, zip codes
- **Binary**
  - Nominal attribute with only 2 states (0 and 1)
  - Symmetric binary: both outcomes equally important
    - e.g., gender
  - Asymmetric binary: outcomes not equally important.
    - e.g., medical test (positive vs. negative)
    - Convention: assign 1 to most important outcome (e.g., HIV positive)
- **Ordinal**
  - 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

*0 to 1m*

- 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 to 1m*

- 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** *2 סוגי תכונות: 1. סופית, 2. אינסופית*
  - ❑ 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** *2 סוגי תכונות: 1. ממשי, 2. מרוכב*
  - ❑ 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