



CS 412 Intro. to Data Mining

Chapter 2. Getting to Know Your Data

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





Data

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

1	12	2	5
2	1	15	9
2	1	0	1
1	0	-1	20
0	-1	19	6
1	19	6	-5
1	19	6	-5

1	12	2	5
1	2	1	15
2	1	0	1
1	0	-1	20
0	-1	19	6
1	19	6	-5
1	19	6	-5

1	12	2	5
1	2	1	15
2	1	0	1
1	0	-1	20
0	-1	19	6
1	19	6	-5
1	19	6	-5

4D = 3D + 1 dimension

Data

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

÷

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

1D = กວດหนึ่งรายการ

2D = กວດทั้งหมด

3D = 20 สำนักงาน

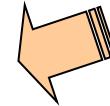
Data

	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

ตัว Data ที่ Attribute มาก็จะมีผลลัพธ์เป็นผลเดียวกัน

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

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

	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

ตารางเชื่อมต่อ กัน

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

รายการ

บันทึก

	team	coach	y	pla	ball	score	game	n	wi	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	0	2	
Document 2	0	7	0	2	1	0	0	0	3	0	0	
Document 3	0	1	0	0	1	2	2	2	0	3	0	

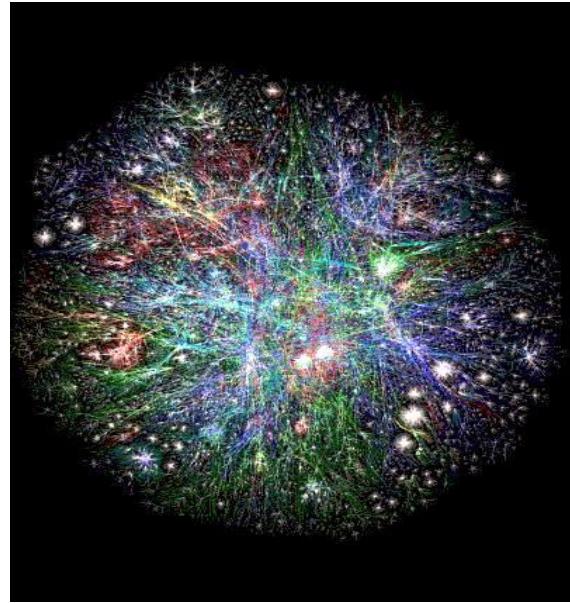
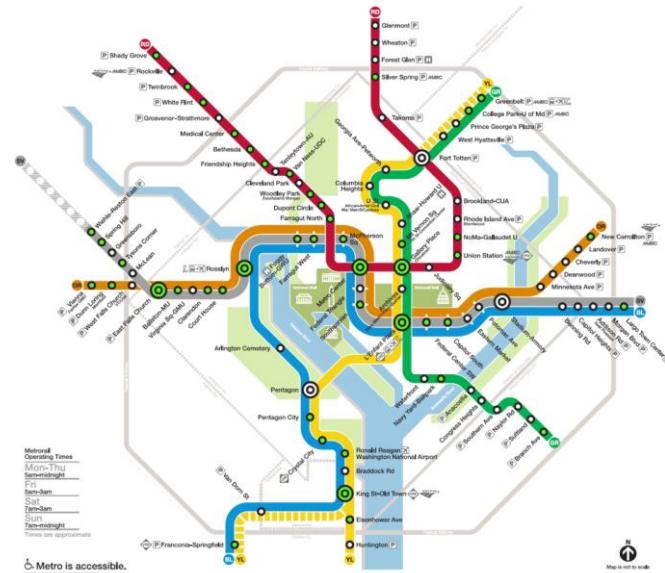
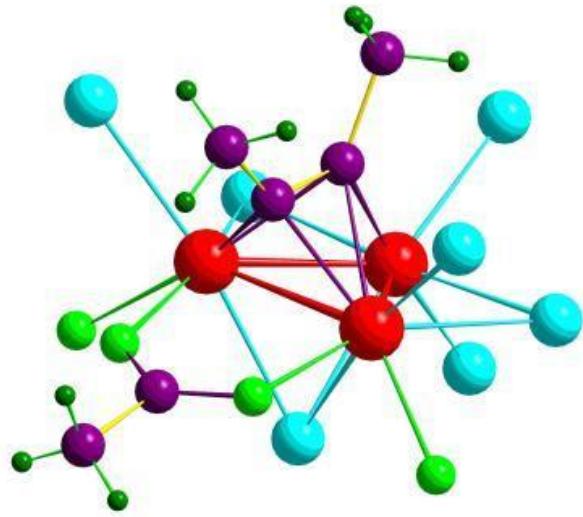
แปลงข้อมูล ให้เป็นค่าเรขา

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

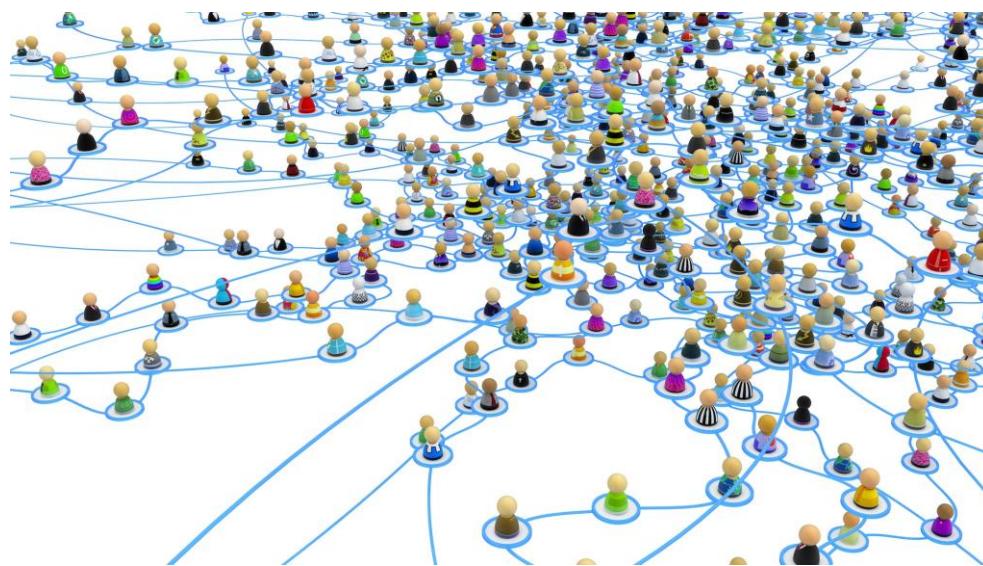
Types of Data Sets: (2) Graphs and Networks

၁၅၈

- ❑ Transportation network
 - ❑ World Wide Web



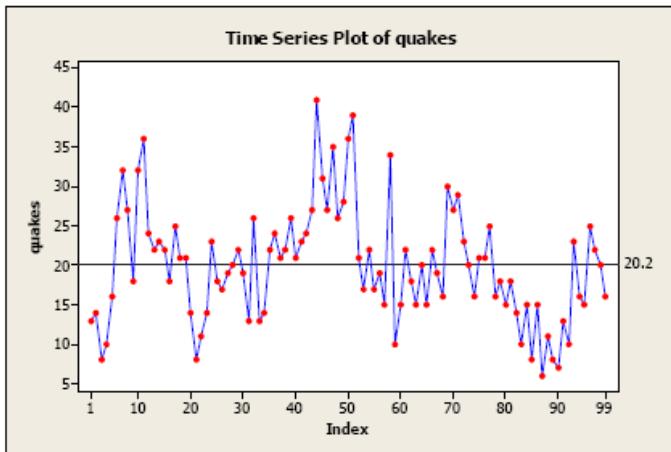
- ❑ Molecular Structures
 - ❑ Social or information networks



Types of Data Sets: (3) Ordered Data

សូមរាជនៅទីនេះបានកំណត់ចុងក្រោយ, ដើម្បីរាជការការណ៍

- ## ☐ Video data: sequence of images



- ## ❑ Sequential Data: transaction sequences

- ## Genetic sequence data

	Start	
Human	GT	TTTGAGG
Chimpanzee	GT	TTTGAGG
Macaque	GT	TTTGAGG
	- - -	ATGTTCAACAAATGCTCCTTCATTCCCTATTTACAGACCTGGCGCA
Human	GACAATTCTGCTAGCAGCCTTGTCTATTATCTGTTTCTAAACCTAGTAATTGAGTGT	
Chimpanzee	GACAATTCTGCTAGCAGCCTTGTCTATTATCTGTTTCTAAACCTAGTAATTGAGTGT	
Macaque	GACAATTCTGCTAGCAGCCTTGTCTATTATCTGTTTCTAAACCTAGTAATTGAGTGT	
	↓	
Human	GATCTGGAGACTAA-CTCTGAAAATAAAAGCTGATTATTATTATTCTCAAAACAA	
Chimpanzee	GATCTGGAGACTAAACTCTGAAAATAAAAGCTGATTATTATTATTCTCAAAACAA	
Macaque	TATCTGGAGACTAAACTCTGAAAATAAAAGCTGATTATTATTATTCTCAAAACAA	
Human	CAGAATACGATTAGCAAATTACTTCTTAAGATATTATTACATTTCATATTCTCTA	
Chimpanzee	CAGAATACGATTAGCAAATTACTTCTTAAGATACTATTACATTTCATATTCTCTA	
Macaque	CAGAATATGATTAGCAAATTACCTCTTAAGATATTATTGCACCTCTATATTCTCTA	
Human	CCCTGAGTTGATGTGTGAGCAATATGTCACTTTCATAAAGCCAGGTATACAC- - -TTATG	
Chimpanzee	CCCTGAGTTGATGTGTGAGCCGTATGTCACTTTCATAAAGCCAGGTATACAC- - -TTATG	
Macaque	CCCTGAGTTGATGTGTGAGCAATATGTCACCTCACAAGGCCAGGTATATACATTACG	
	H I Y S T F L S K	
Human	GACAGGTAAGTAAAAAACATATTATTATTCTACGTTTTGTCCAAAATTTAAATTTC	
Chimpanzee	GACAGGTAAGTAAAAAACATATTATTATTCTACGTTTTGTCCAAAATTTAAATTTC	
Macaque	GACAGGTAAGTAAAAA-CATATTATTATTCTAGGTTTTGTCCAAAAGAGTTAAATTTC	
Human	AACTGTTGCGCGTGTGTTGGTAA- - -TGTAAAAACAAACTCAAGTACA	
Chimpanzee	AACTGTTGCGCGTGTGTTGGTAA- - -TGTAAAAACAAACTCAAGTACA	
Macaque	AACTGTTGCGCGTGTGTTGGTAA- - -CGTAAAAACAAATTCAGTACG	

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

សំណុះផ្លូវការ

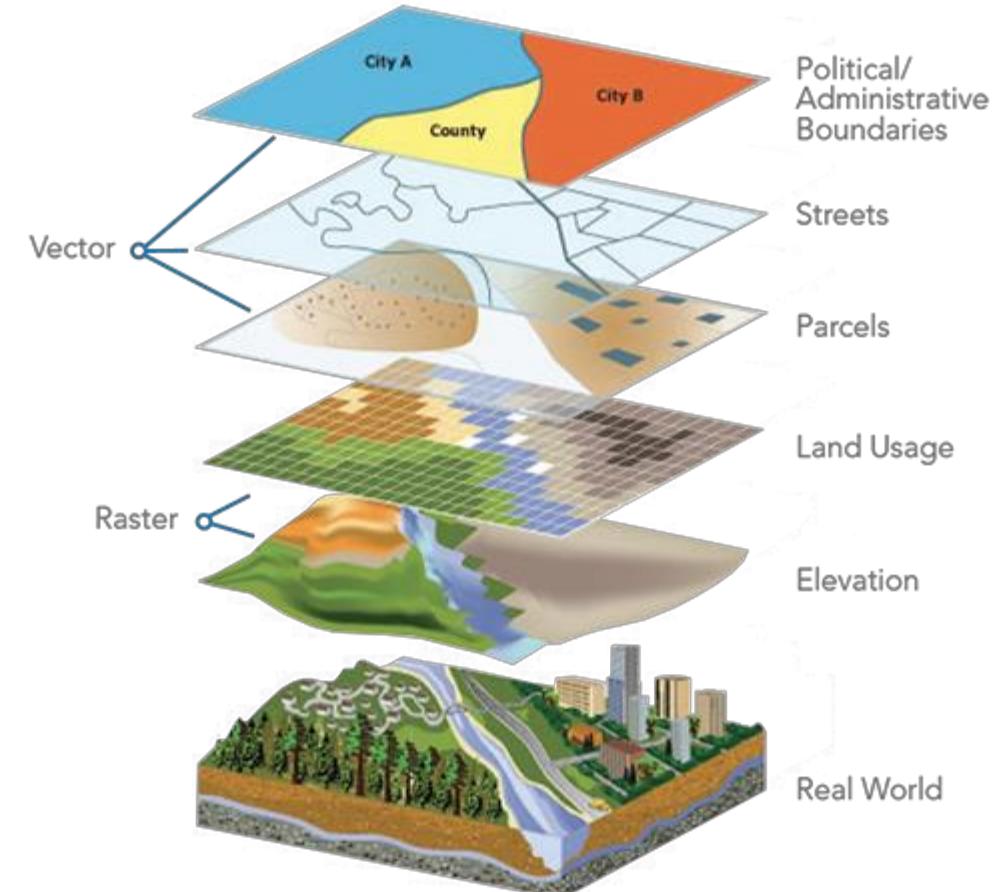
□ Spatial data: maps



សំណុះផ្លូវការ

□ Image data:

□ Video data: Spatio - temporal



Important Characteristics of Structured Data

- Dimensionality *չորս Dimension*
 - Curse of dimensionality
- Sparsity *សង្គមទូទាត់បានចរប់ដែលខ្ពស់* (ការងារអ្វីរបៀន 0 យោងគឺមានការងារ)
- Only presence counts
- Resolution *ក្រុមចំណាំតិច*
 - Patterns depend on the scale
- Distribution *តួនាទីការងារ* (*ក្នុងការងារ / ខ្លួន*)
 - Centrality and dispersion

Data Objects

- ❑ Data sets are made up of data objects ក្នុងទូរស័ព្ទមានចំណាំទូទាត់ខាងក្រោម
- ❑ A **data object** represents an entity ជំនួយណាមីតិ៍ entity
- ❑ 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*
- ❑ Data objects are described by **attributes** ជំនួយនេះគឺជាបាយការណា attributes
- ❑ Database rows → data objects; columns → attributes

Attributes

គម្រោងបច្ចុប្បន្ន

- 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) *ជំនួយក្នុងសម្រាប់ព័ត៌មានទំនាក់ទំនាក់*
 - Binary (e.g., {true, false}) *សម្រាប់ព័ត៌មាន 2 តំបន់*
 - Ordinal (e.g., {freshman, sophomore, junior, senior}) *ដំឡើងតាមលំដាប*
 - Numeric: quantitative *សម្រាប់ +, -, ×, ÷ រៀងចាយ*
 - 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? *Numeric*

Attribute Types

ផ្នែក

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

- 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

ពីតាមរយៈ

Discrete Attribute មានតាមរយៈរាយ

- 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

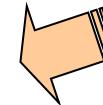
សំខាន់

Continuous Attribute មានតាមរយៈបញ្ហា

- 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

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

លោកស្រីសាស្ត្រ

□ 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

ពីរុយបេសនាការតាមរយៈការងារ

និរនោតាមការងារនៃការងារ

