



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
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

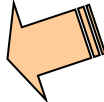
	1	12	2	5
1	2	11	7	2
2	1	15	9	3
1	0	10	1	-3
0	-1	20	12	-2
-1	1	19	6	-5
1	19	0	-3	
1	19	0	-3	
1	19	0	-3	

3D

เมทริกซ์ 1 มิติ อาจมีค่าที่ว่างหรือขาดหายไป
ส่วนเมทริกซ์ 2 มิติ matrix

	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

- Relational records
 - 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

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

no relation

- 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

	team	coach	y	pla	ball	score	game	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	3	0	0	0
Document 3	0	1	0	0	1	2	2	0	3	0	0

มีหลาย

Data structure

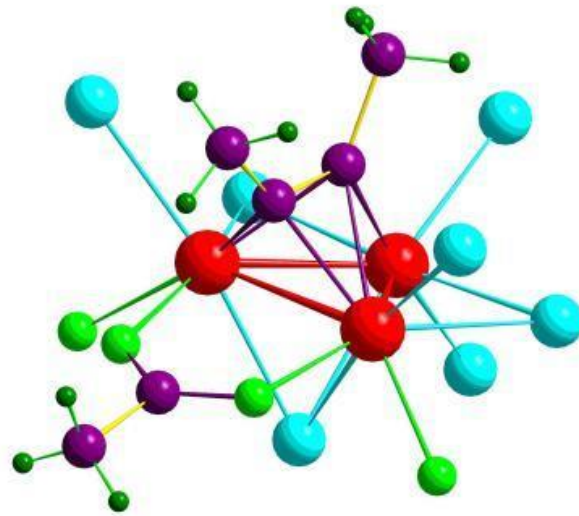
- Document data: Term-frequency vector (matrix) of text documents

normalization ทำให้ง่ายต่อการเปรียบเทียบ

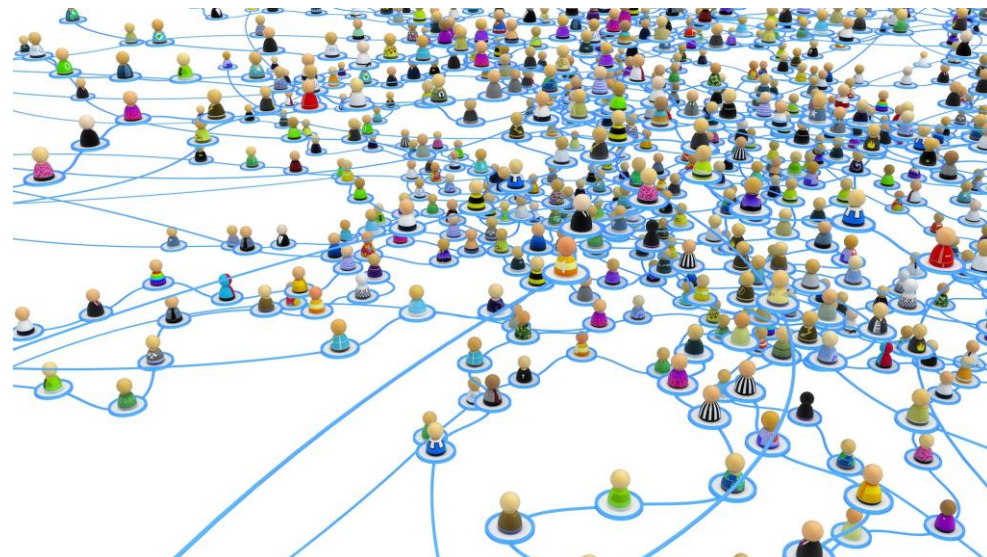
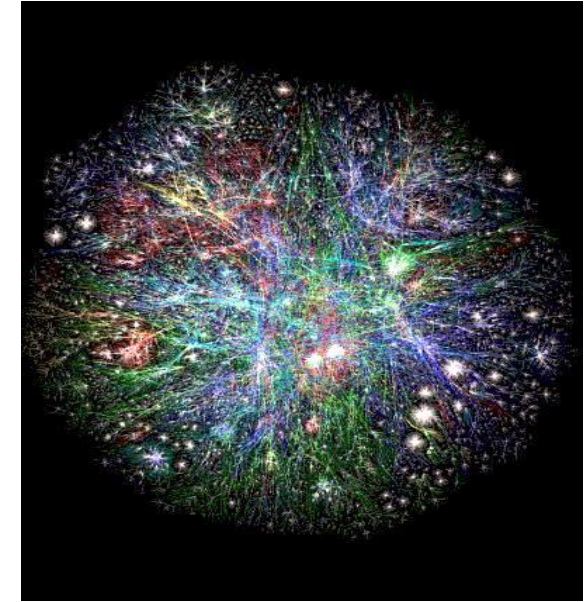
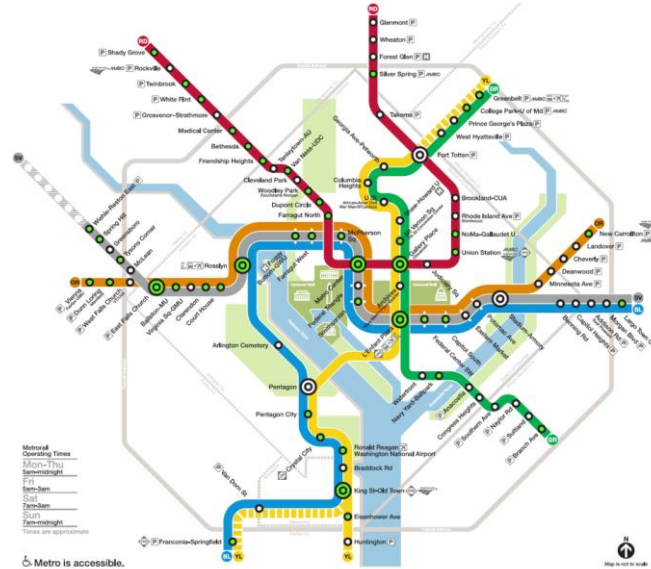
หลายรายการที่มีค่าเหมือนกัน

Types of Data Sets: (2) Graphs and Networks

- ❑ Transportation network
- ❑ World Wide Web



- ❑ Molecular Structures
- ❑ Social or information networks



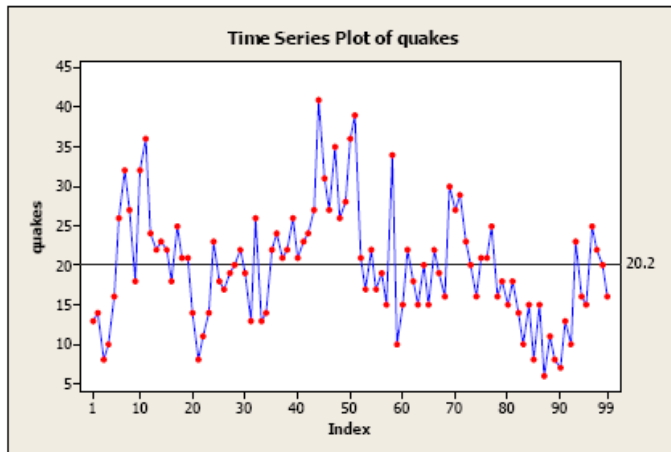
ทำไมต้องทำแบบนี้ ?!

Types of Data Sets: (3) Ordered Data

ข้อมูลมาทีละที, time series, ขึ้น, มีลำดับเวลาชัดเจน
ข้อมูลจัดเรียงตามเวลาที่ขึ้น
หน้าไปมาขึ้นต้น

❑ Video data: sequence of images

❑ Temporal data: time-series



❑ Sequential Data: transaction sequences

❑ Genetic sequence data

	Start
Human	GTTTGGAGG --- ATGTTCAACAAATGCTCCTTTTCATTCTCTATTTACAGACCTGCCGCA
Chimpanzee	GTTTGGAGG --- ATGTTCAATAAATGCTGCTTTCACTCCTCTATTTACAGACCTGCCGCA
Macaque	GTTTGGAGG --- ATGTTCAATAAATGCTCCTTTTCATTCTCTATTTACAAACTGCCGCA
Human	GACAATTCTGCTAGCAGCCTTTGTGCTATTATCTGTTTTCTAAACTTAGTAATTGAGTGT
Chimpanzee	GACAATTCTGCTAGCAGCCTTTGTGCTATTATCTGTTTTCTAAACTTAGTAATTGAGTGT
Macaque	GACAATTCTGCTAGCAGCCTTTGTGCTATTATCTGTTTTCTAAACTTAGTAATTGAGTGT
Human	GATCTGGAGACTAA - CTCTGAAATAAATAAGCTGATTATTTATTTATTTTCTCAAAACAA
Chimpanzee	GATCTGGAGACTAAACTCTGAAATAAATAAGCTGATTATTTATTTATTTTCTCAAAACAA
Macaque	TATCTGGAGACTAAACTCTGAAATAAATAAGCTGATTATTTATTTATTTTCTCAAAACAA
Human	CAGAATACGATTTAGCAAAATTACTTCTTAAGATATTATTTTACATTTCTATATTCTCCTA
Chimpanzee	CAGAATACGATTTAGCAAAATTACTTCTTAAGATATTATTTTACATTTCTATATTCTCCTA
Macaque	CAGAATATGATTTAGCAAAATTACTTCTTAAGATATTATTTTGCATTTCTATATTCTCCTA
Human	CCCTGAGTTGATGTGTGAGCAATATGTCACCTTTTCATAAAGCCAGGTATACA --- TTATG
Chimpanzee	CCCTGAGTTGATGTGTGAGCCGATATGTCACCTTTTCATAAAGCCAGGTATACA --- TTATG
Macaque	CCCTGAGTTGATGTGTGAGCAATATGTCACCTTCACAAAGCCAGGTATATATACATTACG
Human	GACAGGTAAGTAAAAACATATTATTATTCTACGTTTTGTCCAAAAATTTTAAATTTTC
Chimpanzee	GACAGGTAAGTAAAAACATATTATTATTCTACGTTTTGTCCAAAAATTTTAAATTTTC
Macaque	GACAGGTAAGTAAAAACATATTATTATTCTACGTTTTGTCCAAAAATTTTAAATTTTC
Human	AACGTGTGCGCGTGTGTTGGTAA --- TGTAAAACAAAC TCAGTACA
Chimpanzee	AACGTGTGCGCGTGTGTTGGTAA --- TGTAAAACAAAC TCAGTACA
Macaque	AACGTGTGTGCAATGTGTTGGTAA --- CBTAAAACAAATTCAGTACG

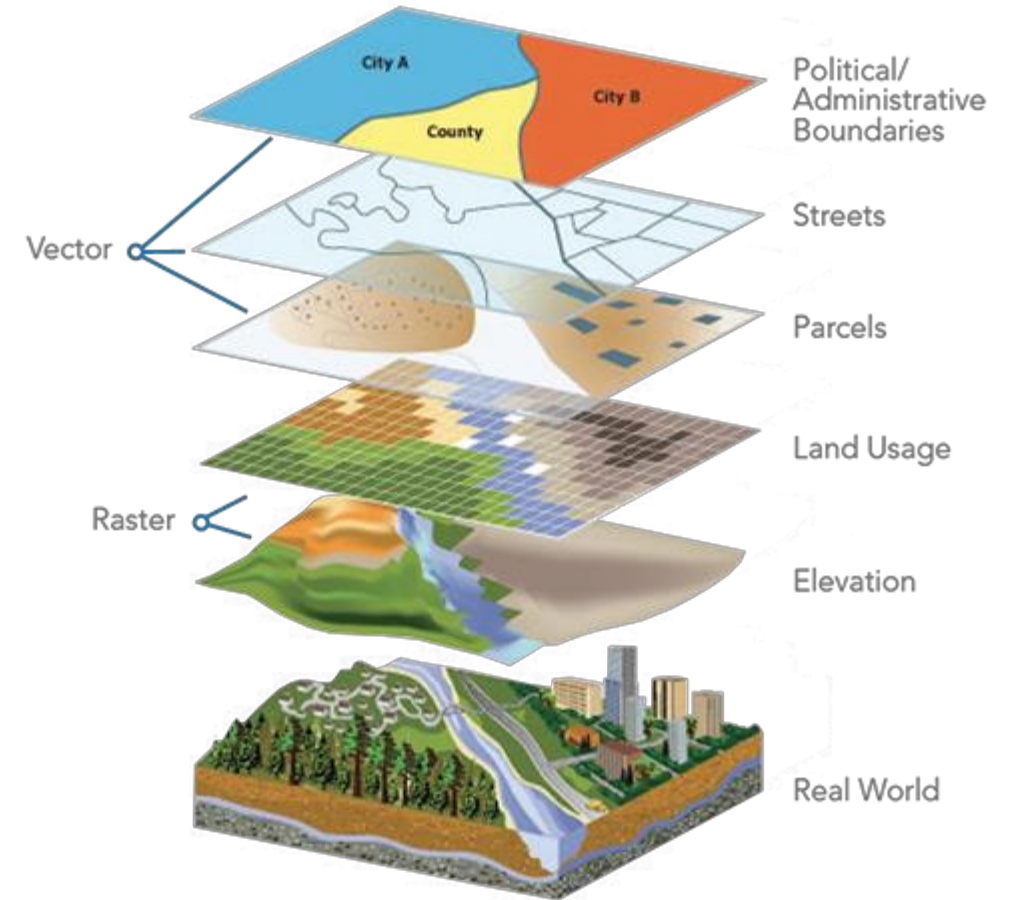
Types of Data Sets: (4) Spatial, image and multimedia Data

❑ Spatial data: maps



❑ Image data:

❑ Video data:



Important Characteristics of Structured Data

- Dimensionality มีมิติมากขึ้น ไม่ได้หมายความว่า 2, 3, 4, 5
 - Curse of dimensionality
- Sparsity สหภาพที่มีข้อมูล
 - Only presence counts
- Resolution เก็บข้อมูลได้ละเอียด
 - Patterns depend on the scale
- Distribution จัดค่าเฉลี่ยค่า ส่วนที่ 0 50
 - Centrality and dispersion

Data Objects

- ❑ Data sets are made up of data objects *ក្រុមសំណុំទិន្នន័យ: គ្របដណ្តប់ដោយ data*
- ❑ A *វត្ថុទិន្នន័យ* **data object** represents an entity
- ❑ Examples:
 - ❑ sales database: customers, store items, sales *តារាងលក់*
 - ❑ medical database: patients, treatments
 - ❑ university database: students, professors, courses
- ❑ Also called *samples, examples, instances, data points, objects, tuples*
- ❑ Data objects are described by **attributes** *កំរិតលក្ខណៈឬលក្ខណៈទិន្នន័យ attributes*
- ❑ Database rows → data objects; columns → attributes

Attributes คุณสมบัตินี้ใช้แยกข้อมูลแต่ละตัว

□ 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) ชื่อของกลุ่ม บางทีก็ไม่มีตัวเลข

□ Binary (e.g., {true, false}) ข้อมูลที่มีแค่สองค่า

□ Ordinal (e.g., {freshman, sophomore, junior, senior}) ข้อมูลเรียงลำดับ

□ Numeric: quantitative +, -, ×, ÷ ได้ผลคูณ: บวก ลบ คูณ หาร

□ 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?

Attribute Types

↓ ชนิดของ attribute

□ **Nominal:** categories, states, or “names of things”

□ *Hair_color* = {auburn, black, blond, brown, grey, red, white}

□ marital status, occupation, ID numbers, zip codes

□ **Binary** *เหมือนกับ Nominal แต่มีแค่ 2 สถานะ 0 กับ 1, 9.6 กันบ้าง*

□ Nominal attribute with only 2 states (0 and 1)

□ Symmetric binary: both outcomes equally important

□ e.g., ~~gender~~, left/right-handed, coke/pepsi, hot/cold

□ Asymmetric binary: outcomes not equally important. *ที่ 2 สำคัญกว่าอีกอัน*

□ e.g., medical test (positive vs. negative)

□ Convention: assign 1 to most important outcome (e.g., HIV positive)

□ **Ordinal** *รู้ลำดับก่อนหลังมาตามลำดับ, ก็อาจวัดค่าไม่ได้*

□ Values have a meaningful order (ranking) but magnitude between successive values is not known

□ *Size* = {small, medium, large}, grades, army rankings

Numeric Attribute Types

→ ចំណាត់ថ្នាក់លេខ
→ វិធីសាស្ត្រវាស់វែង

- 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


❑ Discrete Attribute *มีค่าจำกัด 2 ค่า ไม่ต่อเนื่องกัน อยู่ตรงกลาง เช่น เพศ อายุ 20 ไม่ได้อยู่ตรงกลาง*

- ❑ 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 *มีค่าต่อเนื่องกัน 2 ค่า สามารถเป็น 160 หรือ 161 หรือ 160.5*

- ❑ 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

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:

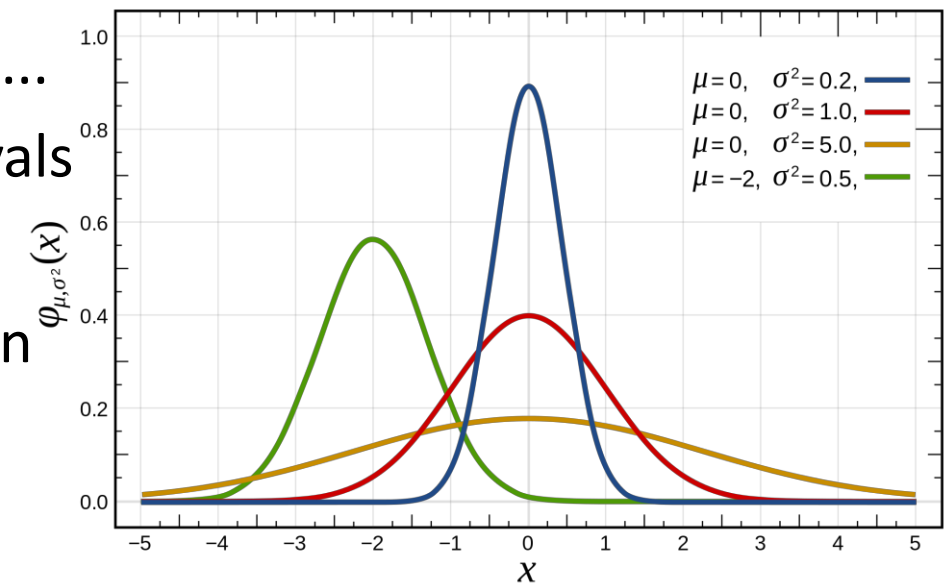
- Analyzed with multiple granularities of precision

- 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



อนาลันในข้อมูลที่มีอยู่ ๒๐ ปี
ฐานนิยมค่าเฉลี่ย