

Chapter 2

2.1

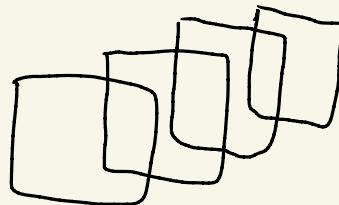
Datam
↓
18,5,4,1,3,0

1
2
1
0
-1
1

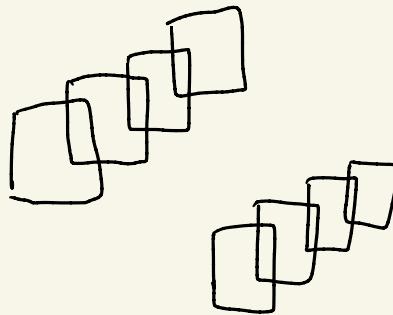
1D

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



3D



4 D

Attribute 1 Attribute 2 ... Attribute n

Records 1

Records 2

:

Records b

Attribute 1 Attribute 2 ... Attribute n

Chapter 2

၃၂၁

Types of Data Sets: (1) Record Data

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

Crossstab of Attribute 2 w/o	Other	English	French	Spanish	USA	Total
Active Students (Initial Guess)	12,491	4,091	1,076	245,013	233,730	
Active Students (Final Guess)	13,081	4,381	1,161	245,013	233,730	
Inactive Growth Guess	3,000	9,000	9,000	132,000	149,000	
Inactive Growth Actual	3,000	9,000	9,000	143,000	143,000	
Strength High School	3,000	6,000	7,000	233,000	244,000	
Strength Middle School	3,000	6,000	7,000	233,000	244,000	
Strength Adult Inactive	6,000	6,000	6,000	209,250	216,250	
Strength Youth Inactive	5,000	5,000	5,000	74,000	72,000	
Total	14,400	43,400	44,000	639,250	639,250	

Person	Per_ID	Surname	First Name	City
0	Miller	Peter	London	
1	Ortega	Alejandra	Valencia	→ no relation
2	Hansen	Karen	Bremen	
3	Blanc	Gabriel	Paris	
4	Bertoldi	Fabrizio	Rom	

Car	Car_ID	Model	Year	Value	Per_ID
0	101	Beatley	1973	100000	0
1	102	Rolls Royce	1965	330000	0
2	103	Perssonet	1993	500	3
3	104	Ferrari	2001	150000	4
4	105	Renault	1998	2000	1
5	106	Renault	2001	7000	3
6	107	Smart	1999	2000	2

- | Transaction data | | Term | | | | | | | | | | | |
|------------------|----------------------------|------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| TID | Item | WED | THU | FRI | SAT | SUN | MON | TUE | WED | THU | FRI | SAT | SUN |
| 1 | Bread, Coke, Milk | 3 | 8 | 5 | 0 | 2 | 8 | 2 | 9 | 3 | 7 | 0 | |
| 2 | Bread, Beer | 0 | 7 | 0 | 2 | 1 | 0 | 3 | 0 | 0 | 0 | 0 | |
| 3 | Beer, Coke, Popcorn, Milk | 3 | 0 | 1 | 0 | 1 | 2 | 0 | 3 | 0 | 0 | 0 | |
| 4 | Bread, Coke, Diapers, Milk | 5 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 5 | Coke, Diapers, Milk | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | |

Document 1

Document 2

Document 3

Document 4

Document 5

ມັນຕ່າງໆ ນຳດຳ
ກໍ່າຍໃຈແລ້ວໂສຫ້ເຈັບ ໃນອອນພໍາໄຕ

四

Term-frequency matrix \rightarrow Text

ມັນຈໍານວນເພື່ອມ ບອກກັບບາດກາມ

Types of Data Sets: (2) Graphs and Networks

100

- #### Transportation network



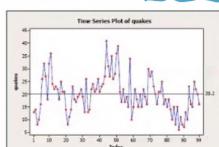
- ❑ World Wide Web
 - ❑ Molecular Structures
 - ❑ Social or information networks



Types of Data Sets: (3) Ordered Data

ຊົນລົງເຈົ້າເກມທີບາຂອງ

- ❑ Video data: sequence of images
 - ❑ Temporal data: time series



- ❑ Sequential Data: transaction sequences
↳ DNA សម្រាប់ការប្រើប្រាស់
 - ❑ Genetic sequence data

→ DNA សំគាល់ត្រូវបានដាក់

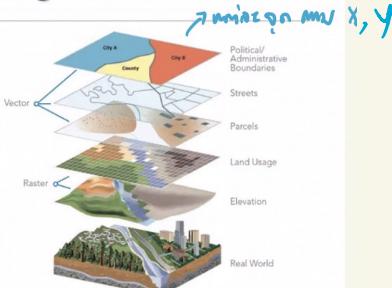
- ## Genetic sequence data

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

ពិភេស្ត

□ Spatial data: maps

□ Image data:



□ Video data: Spatio-temporal

ឯកសារ ពេលវេលា

ពិភេស្តធម៌សែនសម្រាប់វិដ៉ូ Spatio-temporal នឹងត្រួតពេលវេលា

របៀបរាយការណ៍ទីតាំង

Important Characteristics of Structured Data

- Dimensionality → ឈាន់
• Curse of dimensionality
- Sparsity → សំខាន់សំខាន់របស់ពីរ
• Only presence counts
- Resolution → រចនាផ័ត៌មាននៃរឿង
• Patterns depend on the scale ក្នុងម៉ោង pixels
- Distribution → ការចែកចាយពេលវេលា
• Centrality and dispersion

Ex: *mnist dataset*

↳ Ex: *automobile Data* បានរាយការណ៍

Data Objects

- Data sets are made up of data objects. → *ពិភេស្តនៃរឿង = data set*
- A **data object** represents an entity. *ឯកសារ*
- Examples: *ឯកសារពេន្ធអំពីរ*
- sales database: customers, store items, sales *ឯកសារ សំលើរូបរាង*
- medical database: patients, treatments
- university database: students, professors, courses *ឯកសារ សាកលវិទ្យាល័យ*
- Also called *samples*, *examples*, *instances*, *data points*, *objects*, *tuples*. *Data 1on*
- Data objects are described by **attributes**. *ឯកសារពេល*
- Database rows -> data objects; columns -> attributes.

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) សម្រាប់ប្រភេទដែលមិនមែនតម្លៃបាន ស្ថិត H.R. , នាយកដំបូង
 - Binary (e.g., {true, false}) គឺមានតម្លៃ 1 as yes no
 - Ordinal (e.g., {freshman, sophomore, junior, senior}) គឺមានរារាំងចំណាំ នាក់រវាងចំណាំរាយ
 - Numeric: quantitative គឺជាដឹងទិន្នន័យ ដែលមានលក្ខណៈ
 - Interval-scaled: 100°C is interval scales $100^\circ\text{C} + x \neq 100^\circ\text{C}$ នៅពេលគុណភាពនៃរាយបានផ្តល់មួយ
 - 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

↳ సమానం - సమీక్షలు

66 Numeri
Yāma

ເອົ້າຮຽນທີ່ໄດ້ໃຫຍ່ Data

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) \rightarrow 2 อย่าง
 - **Symmetric binary:** both outcomes equally important \rightarrow สมมาตร (ล้วนๆ)
 - **Asymmetric binary:** outcomes not equally important. \rightarrow ไม่สมมาตร (ล้วน ≠)
 - **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** → បុរាណ = តាមរយៈលក្ខណៈ
Measured on a scale of **equal-sized units**
Values have order
E.g., temperature in C° or F° , calendar dates
No true zero-point
Inherent **zero-point**
 - Ratio** → បុរាណ = បានពីរនៅទីមួយ = មិនមែនស្រី
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

• $\{T_1, T_2, \dots, T_n\}$ $T_{min} = min(T_i)$

សៀវភៅ និងការ រំលែក

សុវត្ថិ នៅ ការបែងចែក

ବ୍ୟବ ଶକ୍ତିମାତ୍ର ୦ ୧୯୦ = ୨୫୪

Discrete vs. Continuous Attributes

- Discrete Attribute** ດຳເນີນເອົ້າ 40000 20001 ຍັງແກ່ໄດ້ຮັບອະນຸຍາດ
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**
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

- ຈຳນວນໃຫຍ່ ມີລາຍລະອຽດມີຄວາມປົກຕົວຂອງສິ່ງ 20 ຈຳນວນທີ່ມີຄວາມປົກຕົວຂອງສິ່ງ 19
Max = 21
Min = 19
ກົດລົງທຶນທີ່ມີຄວາມປົກຕົວຂອງສິ່ງ
- Data Objects and Attribute Types
 - Basic Statistical Descriptions of Data 
 - Data Visualization
 - Measuring Data Similarity and Dissimilarity
 - Summary

Basic Statistical Descriptions of Data

- Motivation** ໝາຍາ ໃຊ້ສູນເພົ່າຫຼັກສົດ
- To better understand the data: central tendency, variation and spread
- Data dispersion characteristics** ມັນຕະລິມາດຕະລິ
- median, max, min, quantiles, outliers, variance, etc.
- Numerical dimensions** ພະຍານຕະຫຼາມ ຢູ່ມືນາ
- Data dispersion: analyzed with multiple granularities of precision
- Boxplot or quantile analysis on sorted intervals
- Dispersion analysis on computed measures** ພະຍານຕະຫຼາມ ຢູ່ມືນາ
- Folding measures into numerical dimensions
- Boxplot or quantile analysis on the transformed cube

ຕົວລະບົບ

Measuring the Central Tendency

- Mean (algebraic measure) (sample vs. population):** $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ $\mu = \frac{\sum x}{N}$
Note: n is sample size and N is population size.
- Weighted arithmetic mean:
- Trimmed mean: chopping extreme values
- Median:** $\bar{x} = \frac{\sum w_i x_i}{\sum w_i}$
- Middle value if odd number of values, or average of the middle two values otherwise
- Estimated by interpolation (for grouped data):
- Mode** $median = L_1 + \left(\frac{n/2 - (\sum freq)}{freq} \right) width$
- Value that occurs most frequently $freq_{max}$
- ຈິງລົງ** Unimodal, bimodal, trimodal
- Empirical formula:

$$mean - mode = 3 \times (mean - median)$$

age	frequency
1-5	200
6-15	450
16-20	300
21-50	1500
51-80	700
81-110	44