

Data

ເບີໂທລະຫວ່າງແມວ Matrix ສໍາຄັນໄດ້				
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		1	

1D

2D

3D

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

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
Inflex Crochet Glove	3.00	6.00	8.00	132.00	149.00
Inflex 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	201.00	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

Cross-tab

Person	Pen ID	Surname	First Name	City
1	Miller	Reed	London	UK
2	Hansen	Urs	Zurich	Switzerland
3	Bleier	Gertie	Paris	France
4	Breitkreis	Ulf	Rome	Italy

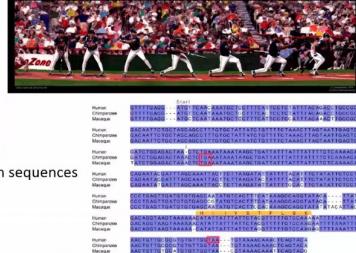
Car	Car ID	Model	Year	Value	Pen ID
101	Scallop	1971	119000	0	
102	Navi Xray	2000	330000	2	
103	Scallop	1991	20000	2	
104	Accraff	2009	150000	4	
105	Jewell	1998	2000	2	
106	Scallop	2009	20000	2	
107	Scarf	1999	2000	2	

TID	Items
1	Bread, Coke, Milk
2	Bread, Bread
3	Bread, Coke, Diaper, Milk
4	Bread, Bread, Diaper, Milk
5	Coke, Diaper, Milk
Total	14.00 43.00 94.00 3.00 3,173.00 2,086.00

↳ ແກ້ໄຂຂ່າຍສຸກທີ່ກິ່ງ text ອີ່ຕົ້ນກວາມຄ່າ ໂດຍອ່ານວ່າຈະບໍ່ມີຄວາມຄ່າ

Types of Data Sets: (3) Ordered Data

- ❑ Video data: sequence of images



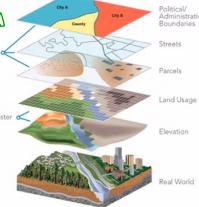
- ❑ Temporal data: time-series

- ❑ Sequential Data: transaction sequences

- ❑ Genetic sequence data

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

- ❑ Spatial data: maps



- ❑ Image data:

- ❑ Video data: Spatio-temporal

- Spatial data ດີ່ຈີ່ກຳນົດ (ກວ້າງ, ຖົກ) (X, Y)

Important Characteristics of Structured Data

- ❑ Dimensionality

- ❑ Curse of dimensionality

- ❑ Sparsity

- ❑ 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
Inflex Crochet Glove	3.00	6.00	8.00	132.00	149.00
Inflex 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	201.00	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

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?

Nominal : ចំណាំតម្លៃទីតាំង (បែងការលក់អគ្គន៍)

Binary : មាត្រា 2 ភាព ខ្លោន T, F

Ordinal : លេខរូបិយប័ណ្ណ ក្រឡើងតំបន់ ខ្លោន freshman, sophomore, junior, senior

Numeric : ចំណាំតម្លៃទីតាំងដែលមានលក្ខណៈជាអនុវត្តន៍

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
 - Measured on a scale of equal-sized units
 - Values have order *ពីរិបាល តាម ចុះឈាម*
 - 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

0 ឬក \rightarrow ខ្លោន 0 ឬក = ខ្លោនឯងកំណុច

0 ឬក \rightarrow ឯករាយ 0 $^{\circ}\text{C}$ \neq ឯករាយឯករាយ ឬក = ឯករាយ

Discrete vs. Continuous Attributes

- Discrete Attribute *ផ្ទាល់ខ្លួន ឬនៅលើពីរ*
 - Has only a finite or countably infinite set of values *ស្ថាបន្ទាប់ឯង*
 - 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 *ផ្ទាល់ខ្លួន ឬនៅលើពីរចុះឈាម*
 - Has real numbers as attribute values *រាយកើត, ចំណាំ*
 - 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 ໜັງລູກ
▢ To better understand the data: central tendency, variation and spread
- ▢ Data dispersion characteristics ສົ່ວນຕະຫຼາດຂອງຂາຍທີ່ມີມຸງ
▢ 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

