ชื่อ นางสาวชนัญชิดา เมธีกุลมานิต รหัสนักศึกษา 623020516-1

สรุปบทที่2.1

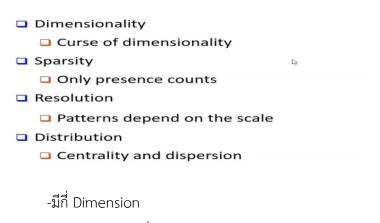
Chapter2 Getting to know Your Data

Data

										Pe	erson:									
 Relational records Relational tables, highly structured 										ers_ID	Surname Miller		Fire	First_Name		City				
										0			Paul		_		8			
								H	1	_	Ortega Alvan			Valencia		- no relation				
								\vdash	2	Huber Blanc		+	Urs Gaston		Zurich	-				
u Data	Data matrix, e.g., numerio				cal r	al matrix, crosstabs				\vdash	4	Bertolini			Saston Fabrizio		Paris Rom	-		
		Oina	England	France	Japan	USA	Total			_	-	Oction	Onnie.	1	-	_	Hotel	_		
I.	Active Quidoors Crochet Glove	CHELL	12.00	4.00	-		257.00			_	ar:	_		_		_		_	1	
17	Active Outdoors Lycra Glove		10.00	6.00		223,00	324.00				ar_(D	_	odel	+	Year	-	Value	Pers_ID 0	4	
Inflox Crecket Glove		3.00	_	8.00		132.00	149.00			-	102	Ben	tley Royce	+	1973	+	100000	0	-	
1 1	Influx (yora Glove		2.00			143,00	145.00			-	103	_	geot	+	1993	+	500	3	-	
1 1	Triumph Pro Helmet Triumph Vertige Helmet		3.00	7.00	_	333.00 474.00	344.00			-	104	Fer		+	2005	-	150000	4		
1 1	Xtrume Adult Helmet Xtrume Voulb Helmet Total		6.00	7.00			276.00				105	Ren	ault		1998		2000	3	1	
×			1.00			76.00	77.00				106	Ren	ault		2001	\neg	7000	3	1 -	
1			#3.00 34.00 3.00		1,972.00	2.084.00				107	Smart			1999		2000	2			
☐ Trans	saction da	ata						team	coach	y pia	ball	some	game	3 E	lost	tim cout	season			
1	Bread, Coke,	Milk					78													
2	2 Beer, Bread						Document 1	3	0	5	0	2	6	0	2	0	2			
3	3 Beer, Coke, Diaper, Milk 4 Beer, Bread, Diaper, Milk						Discotore 2	2 0	7	0	2	1	0	0	3	0	0			
4							Docutent 2		1	0	2	,	0	0	3	0	0			
	5 Coke, Diaper, Milk			_			Document 3	0	1	0	0	1	2	2	0	3	0			

Data Relational คือดาต้าที่มีหลายตารางและมีความสัมพันธ์กัน

Important Characteristics of Structured Data



- -สนใจเฉพาะช่องที่มีข้อมูล
- -ความแอดอัดในการเก็บข้อมูล เช่นรูปภาพว่ามีขนาดภาพกี่พิกเซล ถ้าเกิดเรสรูต่ำเราจะแทนค่าสีของแต่ ละหลายๆพิกเซลด้วย1สี จะได้เปลืองพื้นที่เก็บน้อย

-วัดค่ากลางของข้อมูล (สว่าง/มืด)

Data Objects

- Data sets are made up of data objects
- A data object represents an 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
- □ Database rows → data objects; columns → attributes

แนวนอนเรียกว่า data objects

แนวตั้ง เรียกว่า attributes

Data ที่เราได้ ว่าจะเรียกว่า Data sets ก็คือประกอบไปด้วยข้อมูลหลายๆข้อมูล

ดาต้า ออฟ คือข้อมูลแต่ละตัว เช่น ตารางเซลล์ก็จะมีการเก็บข้อมูลของคอสตอเมอร์ เซลล์ เป็นต้น

การเรียกดาต้า

Also called samples, examples, instances, data points, objects, tuples

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?

คือ คุณสมบัติที่ใช้อธิบายข้อมูลแต่ละตัว

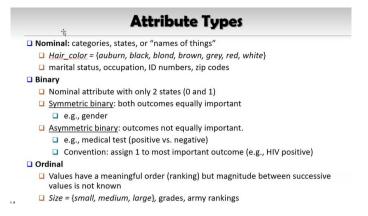
ชนิดของข้อมูล

- -Nominal ข้อมูลที่ไม่มีตัวเลข เช่น สีแดง สีน้ำเงิน
- -Binary เช่น มีข้อมูล 2 ค่า เช่น ถูกกับผิด
- -Ordinal เช่น ข้อมูลที่มีการเรียงลำดับ เช่น เฟรชชี่ หม่อ จูเนียร์ ซีเนียร์ เป็นต้น

Numeric

Interval บวกลบคูณหารแล้วมีความหมาย

Ratio ความสัมพันธ์ของจำนวนทั้งสองจำนวน



		Numeric Attribute Types	
0	Quanti	ty (integer or real-valued)	0
	Interva	ıl	
		Measured on a scale of equal-sized units	
		Values have order	
		■ E.g., temperature in <u>Gor</u> F°, calendar dates	
0		No true zero-point	0
	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	35

ไม่มีศูนย์แท้ เช่น อุณหภูมิ เกรด ความยาว เงิน เป็นต้น

มีศูนย์แท้ เช่น น้ำหนัก ระยะทาง ความสูง ยอดขาย เป็นต้น

Discrete Attribute 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

Discrete vs. Continuous 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

Discrete ความไม่ต่อเนื่อง เช่นจำนวนของสิ่งของ รหัสไปรษณีย์

Continuous มีค่าได้ทุกค่าในช่วงที่กำหนด

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

