

#### แบบฝึกหัด

1. กำหนดแตริบว็ดดังต่อไปนี้ จงจำแนกประเภทของแตริบว็ดว่าเป็นแตริบว็ดเชิงคุณภาพ (qualitative) ชนิดใด (nominal, ordinal) หรือเป็นแตริบว็ดเชิงปริมาณ (quantitative) ชนิดใด (interval, ratio) นอกจากนี้ให้จำแนกด้วยว่าแตริบว็ดแต่ละตัวมีชนิดเป็น Binary, Discrete, หรือ Continuous
  - (ก) เวลาในรูปแบบ AM หรือ PM
  - (ข) ความสว่างที่วัดโดยมิเตอร์วัดแสง
  - (ค) ความสว่างที่วัดจากความรู้สึกของมนุษย์

- (ง) มุมที่วัดเป็นองศาระหว่าง 0 ถึง 360
  - (จ) เหรียญทอง เหรียญเงิน เหรียญทองแดง ของกีฬาโอลิมปิก
  - (ฉ) ความสูงจากระดับน้ำทะเล
  - (ช) จำนวนผู้ป่วยในโรงพยาบาลแห่งหนึ่ง
  - (ซ) เลข ISBN ของหนังสือ
  - (ณ) ความหนาแน่นของสารหน่วยเป็น กรัมต่อลูกบาศก์เซนติเมตร
  - (ญ) ระยะทางจากจุดศูนย์กลางของวิทยาเขตหน่วยเป็นเมตร
  - (ฎ) ชั้นยศของกองทัพบก
  - (ฏ) ความสามารถในการส่งผ่านคลื่นแสง: opaque translucent transparent
2. จงยกตัวอย่างสถานการณ์ที่ identification numbers (รหัสประจำตัว) น่าจะมีประโยชน์สำหรับการทำนาย
  3. ปริมาณใดต่อไปนี้ที่มีคุณสมบัติ spatial autocorrelation : daily rainfall หรือ daily temperature และทำไมจึงเป็นเช่นนั้น
  4. โปรแกรมเมอร์คนหนึ่งได้ออกแบบอัลกอริทึม k-nearest neighbors ดังนี้

---

**Algorithm 2.1** Algorithm for finding  $K$  nearest neighbors.

---

- 1: **for**  $i = 1$  to number of data objects **do**
  - 2:   Find the distances of the  $i^{th}$  object to all other objects.
  - 3:   Sort these distances in decreasing order.  
       (Keep track of which object is associated with each distance.)
  - 4:   **return** the objects associated with the first  $K$  distances of the sorted list
  - 5: **end for**
- 

- (ก) จงอภิปรายว่าจะมีปัญหอะไรเกิดขึ้นได้บ้างกับอัลกอริทึมนี้ ถ้าดาต้าเซตมีข้อมูลซ้ำ (duplicates)
  - (ข) จงเสนอวิธีการแก้ไขปัญหที่เกิดขึ้นจากการมีข้อมูลซ้ำซ้อนในดาต้าเซต
5. คำนวณค่า cosine, correlation, Jaccard และ Euclidean distance ของ ดาต้าอ็อบเจกต์  $x$  และ  $y$  ดังต่อไปนี้
    - (ก)  $x = (1, 1, 1, 1), y = (2, 2, 2, 2)$
    - (ข)  $x = (0, 1, 0, 1), y = (1, 0, 1, 0)$
    - (ค)  $x = (0, -1, 0, 1), y = (1, 0, -1, 0)$
  6. คำนวณค่า Mutual information ของดาต้าอ็อบเจกต์  $x$  และ  $y$  ดังต่อไปนี้
    - (ก)  $x = (-7, -2, 1, 0, 1, 2), y = (9, 4, 1, 0, 4, 1)$
    - (ข)  $x = (1, 1, 1, 1), y = (2, 2, 2, 2)$

### เอกสารอ้างอิง

[1] Pang-Ning Tan, Michael Steinbach, Anuj Karpatne, Vipin Kumar. "Introduction to Data Mining". Pearson, 2nd edition, 2018.

1. ก.) เชิงคุณภาพ ชนิด nominal แบบ binary  
 ข.) เชิงปริมาณ interval แบบ continuous  
 ค.) เชิงคุณภาพ ordinal แบบ discrete  
 ง.) เชิงปริมาณ ratio แบบ continuous  
 จ.) เชิงคุณภาพ ordinal แบบ discrete  
 ฉ.) เชิงปริมาณ interval แบบ continuous  
 ช.) เชิงปริมาณ interval แบบ discrete  
 ซ.) เชิงคุณภาพ nominal แบบ discrete  
 ฌ.) เชิงปริมาณ ratio แบบ continuous  
 ญ.) เชิงปริมาณ interval แบบ continuous  
 ฎ.) เชิงคุณภาพ ordinal แบบ discrete  
 ฏ.) เชิงปริมาณ ratio แบบ continuous

2. - ต้องการรู้ช่วงอายุ  
 - ต้องการรู้ช่วงปีจบ  
 - ต้องการรู้วิชาที่น่านจะเรียนในปัจจุบัน

3. daily temperature เพราะ ข้อมูลจากพื้นที่หนึ่งที่ได้นั้นจะมีค่าใกล้เคียงกลับพื้นที่ ๆ อยู่ใกล้ ๆ

4. ก. - มีประสิทธิภาพต่ำจากการที่ต้องนำตัวเข้ามาประมวลผลร่วมด้วย  
 - ได้เพื่อนบ้านที่เป็นตัวซ้ำกับตัวต้น  
 - หากมีตัวใกล้ ๆ เป็นตัวซ้ำอาจทำให้เพื่อนบ้านที่ได้มีแต่นั้น ๆ  
 ข. จัดรวมกลุ่มตัวซ้ำเป็นกลุ่มเดียวกันแล้วให้บันทึกกลุ่มนั้นเป็น data object ตัวหนึ่ง

5.)

ก. cosine

$$\langle x, y \rangle = 1 \cdot 2 + 1 \cdot 2 + 1 \cdot 2 + 1 \cdot 2 = 8$$

$$\|x\| = (1^2 + 1^2 + 1^2 + 1^2)^{1/2} = 2 \quad \|y\| = (2^2 + 2^2 + 2^2 + 2^2)^{1/2} = 4$$

$$\text{Cosine} = 8 / (2 \cdot 4) = 1$$

$$\text{Correlation mean}(x) = (1+1+1+1)/4 = 1 \quad \text{mean}(y) = (2+2+2+2)/4 = 2$$

$$S_x = (1/4 - 1 \cdot ((1-1)^2 + (1-1)^2 + (1-1)^2 + (1-1)^2))^{1/2} = 0 \quad S_y = (1/4 - 1 \cdot ((2-2)^2 + (2-2)^2 + (2-2)^2 + (2-2)^2))^{1/2} = 0$$

$$S_{xy} = (1/4 - 1 \cdot ((1-1)(2-2) + (1-1)(2-2) + (1-1)(2-2) + (1-1)(2-2)))^{1/2} = 0$$

$$\text{Correlation} = 0/0*0 = 0$$

**Jaccard**

$$f_{1,2} = 4$$

$$J = 0/4 = 0$$

**Euclidean distance**

$$d(x,y) = ((1-2)^2+(1-2)^2+(1-2)^2+(1-2)^2)^{1/2} = 2$$

**u.) cosine**

$$\langle x,y \rangle = 0*1 + 1*0 + 0*1 + 1*0 = 0$$

$$\|x\| = (0^2+1^2+0^2+1^2)^{1/2} = 2^{1/2} \quad \|y\| = (1^2+0^2+1^2+0^2)^{1/2} = 2^{1/2}$$

$$\text{Cosine} = 0/2^{1/2}*2^{1/2} = 0$$

$$\text{Correlation mean}(x) = (1+0+1+0)/4 = 1/2 \quad \text{mean}(y) = (0+1+0+1)/4 = 1/2$$

$$S_x = (1/4-1*((1-1/2)^2+(0-1/2)^2+(1-1/2)^2+(0-1/2)^2))^{1/2} = 1/3^{1/2} \quad S_y = (1/4-1*((0-1/2)^2+(1-1/2)^2+(0-1/2)^2+(1-1/2)^2))^{1/2} = 1/3^{1/2}$$

$$S_{xy} = (1/4-1*((1-1/2)(0-1/2)+(0-1/2)(1-1/2)+(1-1/2)(0-1/2)+(0-1/2)(1-1/2)))^{1/2} = (-1/3)^{1/2}$$

$$\text{Correlation} = (-1/3)^{1/2}/(1/3)^{1/2}*(1/3)^{1/2} = 1/(-1/3)^{1/2}$$

**Jaccard**

$$f_{1,0} = 2 \quad f_{0,1} = 2$$

$$J = 0/2+2 = 0$$

**Euclidean distance**

$$d(x,y) = ((1-0)^2+(0-1)^2+(1-0)^2+(0-1)^2)^{1/2} = 2$$

**u. cosine**

$$\langle x,y \rangle = 0*1 + -1*0 + 0*-1 + 1*0 = 0$$

$$\|x\| = (0^2+1^2+0^2+(-1)^2)^{1/2} = 2^{1/2} \quad \|y\| = (1^2+0^2+(-1)^2+0^2)^{1/2} = 2^{1/2}$$

$$\text{Cosine} = 0/2^{1/2}*2^{1/2} = 0$$

$$\text{Correlation mean}(x) = (0+1+0-1)/4 = 0 \quad \text{mean}(y) = (1+0-1+0)/4 = 0$$

$$S_x = (1/4-1*((0-0)^2+(1-0)^2+(0-0)^2+(-1-0)^2))^{1/2} = 2/3^{1/2} \quad S_y = (1/4-1*((1-0)^2+(0-0)^2+(-1-0)^2+(0-0)^2))^{1/2} = 1/3^{1/2}$$

$$S_{xy} = (1/4-1*((0-0)(1-0)+(1-0)(0-0)+(0-0)(-1-0)+(1-0)(0-0)))^{1/2} = 0$$

$$\text{Correlation} = 0/1/3^{1/2}*1/3^{1/2} = 0$$

**Jaccard**

$$F_{0,1} = 1 \quad f_{1,0} = 1 \quad f_{0,1} = 1 \quad f_{1,0} = 1$$

$$J = 0/(1+1+1+1) = 0$$

**Euclidean distance**

$$d(x,y) = ((0-1)^2+(1-0)^2+(0-(-1))^2+(-1-0)^2)^{1/2} = 2$$

6. n.

$x_j$	$P(X=x_j)$	$-P(X=x_j)\log_2 P(X=x_j)$
-7	1/6	0.43082
-2	1/6	0.43082
0	1/6	0.43082
1	2/6	0.52832
2	1/6	0.43082
$H(x)$		2.2516

$Y_k$	$P(Y=y_k)$	$-P(Y=y_k)\log_2 P(Y=y_k)$
0	1/6	0.43082
1	2/6	0.52832
4	2/6	0.52832
9	1/6	0.43082
$H(y)$		1.91828

$x_j$	$Y_k$	$P(X=x_j, Y=y_k)$	$-P(X=x_j, Y=y_k) \log_2 P(X=x_j, Y=y_k)$
-7	9	1/6	0.43082
-2	4	1/6	0.43082
1	1	1/6	0.43082
0	0	1/6	0.43082
1	4	1/6	0.43082
2	1	1/6	0.43082
$H(x,y)$			2.58492

$$I(x,y) = 2.2516 + 1.91828 - 2.58492 = 1.58496$$

7.

$x_j$	$P(X=x_j)$	$-P(X=x_j)\log_2 P(X=x_j)$
1	4/4	0
$H(x)$		0

$Y_k$	$P(Y=y_k)$	$-P(Y=y_k)\log_2 P(Y=y_k)$
2	4/4	0
$H(y)$		0

$x_j$	$Y_k$	$P(X=x_j, Y=y_k)$	$-P(X=x_j, Y=y_k) \log_2 P(X=x_j, Y=y_k)$
1	2	4/4	0
$H(x,y)$			0

$$I(x,y) = 0 + 0 - 0 = 0$$

# สรุป

## ประเภทของข้อมูล

- **ชุดข้อมูล** คือ กลุ่มของ data objects ประกอบด้วย Attributes ที่บอกลักษณะ

ตารางที่ 2.1. ชุดข้อมูลนิสิต (Student Information Data Set)

### Attributes

ลำดับที่	รหัสนิสิต	ชั้นปี	เกรดเฉลี่ยสะสม
1	1034261	2	2.75
2	1034262	3	3.24
3	1034263	2	3.51
4	1034265	1	2.99
5	1034266	3	3.12

Data objects

Attributes ของ Data objects จำนวน เปลี่ยน แปลง ได้ตลอด  
เช่น อุณหภูมิพื้นผิวโลก-ผิวดิน

## การแบ่งประเภทของ Attributes โดย Operation ของระบบจำนวน

ตารางที่ 2.2. ชนิดของแอททริบิวต์

ชนิดของแอททริบิวต์	คำอธิบาย	ตัวอย่าง	โอเปอเรชัน
Categorical (เชิงคุณภาพ)	Nominal คำของ Nominal attribute สามารถใช้ในการแยกแยะค่าด้วยโอเปอเรชันได้ ด้วยโอเปอเรชัน Distinctness ( $=$ , $\neq$ )	รหัสไปรษณีย์ รหัสพนักงาน สีตา เพศ	ฐานนิยม, entropy, contingency correlation, Chi-squared test
	Ordinal ก. binary ด. Discrete จ. มีคุณสมบัติและโอเปอเรชัน Distinctness เช่นเดียวกับ Nominal attributes และ คำของ Ordinal attribute สามารถใช้ในการเรียงลำดับค่าด้วยโอเปอเรชันได้ ด้วยโอเปอเรชัน Order ( $<$ , $>$ )	ความแข็งแรงของแร่ธาตุ, เกรด {A, B+, B, C+, C, D+, D, F},	มัธยฐาน, เปอร์เซนต์ไทล์, rank correlation, run tests, sign tests
Numeric (เชิงปริมาณ)	Interval มีคุณสมบัติและโอเปอเรชัน Distinctness และ Order เช่นเดียวกับ Nominal attributes และ Ordinal attributes นอกจากนี้ ความแตกต่างระหว่าง interval attributes สองค่า คำนวณได้ด้วยโอเปอเรชัน Addition ( $+$ , $-$ ) สามารถตีความได้ กล่าวคือ interval attributes จะมีหน่วยของการวัด	อุณหภูมิในหน่วย องศาเซลเซียส หรือ องศาฟาเรนไฮต์, วันที่ตามปฏิทิน	ค่าเฉลี่ย, ส่วนเบี่ยงเบนมาตรฐาน, Pearson's correlation, t-test, F-test
	Ratio ว. Continuous มีคุณสมบัติและโอเปอเรชัน Distinctness, Order, และ Interval เช่นเดียวกับ Nominal attributes, Ordinal attributes, และ Interval attributes นอกจากนี้ อัตราส่วนของ ratio attributes ซึ่งคำนวณได้โดยใช้โอเปอเรชัน Multiplication ( $\times$ , $/$ ) สามารถตีความได้	อุณหภูมิในหน่วยเคลวิน (Kelvin), อายุ, มวล, ความยาว, กระแสไฟฟ้า	ค่าเฉลี่ยเรขาคณิต, ค่าเฉลี่ยฮาร์โมนิก, เปอร์เซ็นต์ความผันแปร

# การแบ่งประเภท Attributes. ด้วยจำนวนของค่าที่เป็นไปได้

1. Discrete Attribute ที่มีค่าเป็นไปได้น้อยกว่าหนึ่งร้อยค่า แต่สามารถนับแยกได้ **มักใช้กับ Attributes เชิงคุณภาพ**
2. Binary Attribute มีเพียง 2 ค่า 0/1, จริง/เท็จ, ใช่/ไม่ใช่, ชาย/หญิง
3. Continuous Attribute มีค่าเป็นจำนวนจริง **มักใช้กับ Attributes เชิงปริมาณ**

## ชนิดของ Data sets มี 3 ประเภท

### 1. ข้อมูล Record data

ข้อมูลที่ใช้ในการคำนวณเชิงข้อมูล มักอยู่ในรูปแบบ

ของ Record data ซึ่งประกอบด้วย Attributes หรือ field ข้อมูลแต่ละค่า

Tid	Refund	Marital Status	Taxable Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

(a) Record data.

TID	ITEMS
1	Bread, Soda, Milk
2	Beer, Bread
3	Beer, Soda, Diapers, Milk
4	Beer, Bread, Diapers, Milk
5	Soda, Diapers, Milk

(b) Transaction data.

Projection of Color	Projection of Flavor	Distance	Load	Thickness
10.23	5.27	15.22	27	1.2
12.65	6.25	16.22	22	1.1
13.54	7.23	17.34	23	1.2
14.27	8.43	18.45	25	0.9

(c) Data matrix.

	team	coach	play	ball	score	game	win	lost	timeout	season
Document 1	3	0	5	0	2	6	0	2	0	2
Document 2	0	7	0	2	1	0	0	3	0	0
Document 3	0	1	0	0	1	2	2	0	3	0

(d) Document-term matrix.

### 2. ข้อมูลกราฟ

กราฟเป็นโครงสร้างข้อมูลสำหรับแสดง

ความสัมพันธ์ระหว่าง Data object เรามักใช้กับ

กราฟ 2 ชนิด

#### 1. แสดงความสัมพันธ์ระหว่าง Data object

**Useful Links:**

- Bibliography
- Other Useful Web sites
  - ACM SIGKDD
  - KDDmag
  - The Data Mine

**Knowledge Discovery and Data Mining Bibliography**  
(Get updated frequently, so visit often!)

- Books
- General Data Mining

**Book References in Data Mining and Knowledge Discovery**

Usama Fayyad, Gregory Piatetsky-Shapiro, Padhraic Smyth, and Ramasamy Srikant, "Advances in Knowledge Discovery and Data Mining", AAAI Press/The MIT Press, 1996.

J. Ross Quinlan, "C4.5: Programs for Machine Learning", Morgan Kaufmann Publishers, 1993.

Michael Berry and Gordon Linoff, "Data Mining Techniques (For Marketing, Sales, and Customer Support)", John Wiley & Sons, 1997.

Usama Fayyad, "Mining Databases: Towards Algorithms for Knowledge Discovery", Bulletin of the IEEE Computer Society Technical Committee on data Engineering, vol. 21, no. 1, March 1998.

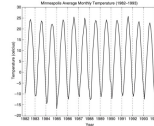
Christopher Mathews, Philip Chan, and Gregory Piatetsky-Shapiro, "Systems for Knowledge Discovery in Databases", IEEE Transactions on Knowledge and Data Engineering, 5(6):903-913, December 1993.

Time	Customer	Items Purchased
T1	C1	A, B
T2	C3	A, C
T2	C1	C, D
T3	C2	A, D
T4	C2	E
T5	C1	A, E

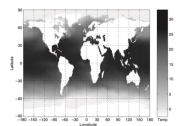
(a) Sequential transaction data.

```
GTTCCGCGCTTCAGCCCGCGCCGCGCGAGAGGCGCCCGCTGGCGCGCGGGGAGGCGGGGCGCCCGAGCCCAACCGAGTCCGACCAAGTGCCCGCTTCGCTCGGCTAGACCTGAGCTCATTAGGGGCGACGCGAGGCCAAGTAGAACACCGGAGCGCTGGGCTGCTGTCGACACAGGG
```

(b) Genomic sequence data.



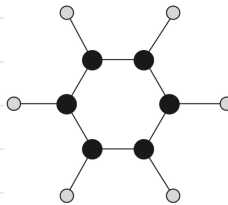
(c) Temperature time series.



(d) Spatial temperature data.

(a) Linked web pages.

### 2. Data object มีโครงสร้างแบบกราฟ



(b) Benzene molecule.

# คุณภาพข้อมูล

- การได้มาซึ่งคุณภาพข้อมูล เทคนิคการวัดการวัดคุณภาพของข้อมูล แบ่งได้ 2 กลุ่ม
1. เทคนิคด้านวิธีการรวบรวมและการนำไปใช้ เช่น คุณภาพข้อมูล
  2. การใช้เทคนิคการคำนวณเชิงข้อมูล ที่คำนวณต่อข้อมูลคุณภาพต่าง

## การวัดคุณภาพจากการวัดและการเก็บข้อมูล

- ปัญหาคุณภาพข้อมูลเกิดได้จากกระบวนการผลิตของมนุษย์ ซึ่งจำกัดของอุปกรณ์ เช่น ข้อมูลรบกวน, ข้อมูลที่ซ้ำ, การแปลอรรถ, การที่ข้อมูล และ การไม่แม่นยำ
- การวัดคุณภาพจากการวัด
  - การวัดคุณภาพจากการเก็บข้อมูล

## ปัญหาคุณภาพที่เกี่ยวกับองค์ประกอบข้อมูล

- ข้อมูลที่มีคุณภาพสูง คือข้อมูลที่นำมาใช้ในการประยุกต์ใช้งาน คุณภาพข้อมูลตามมุมมองของการใช้
- แต่ปัญหาก็มี 3 อย่าง
1. Timeliness เวลาในการเก็บ
  2. Relevance ข้อมูลทุกชั้นที่จำเป็นต่อการวิเคราะห์
  3. Knowledge about the Data คุณภาพเอกสาร

## ประเมิน Present

- ข้อ 8 คะแนนรวม มีเนื้อหาเพิ่มเติม ผู้ศึกษา อธิบายนำฟัง
- ข้อ 2 คะแนนรวม ผู้ศึกษา อธิบายนำฟัง