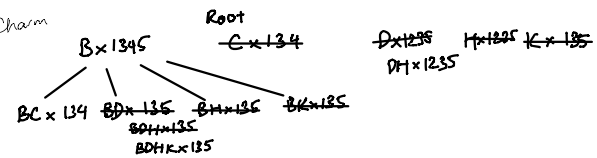


Charm



→ ≈ 3 transactions  
→  $\sigma(X) \geq 3$

Câu 1: Cho CSDL sau và minsupp= 60% và minconf= 100%

Đổi TID sang 1, 2, 3, 4, 5

TID	Items
10	D, H, C, A, B, K, M
20	E, H, D, G, P, I
30	B, C, D, G, H, K
40	E, A, C, B, P, I
50	K, B, M, F, H, D

Chuyển đổi bảng dl:

Item	Tidset	Item	Diffset
A	1, 4	B	2
B	1, 3, 4, 5	C	2, 5
C	1, 3, 4	D	4
D	1, 2, 3, 5	H	4
E	2, 4	K	2, 4
F	5		
G	2, 3		
H	1, 2, 3, 5		
I	2, 4		
K	1, 3, 5		
M	1, 5		
P	2, 4		

$$Do r = \frac{t(PXY)}{d(PXY)} \geq 1$$

nên sử dụng diffset sẽ có lợi hơn

- a) Liệt kê các tập phổ biến tối đại và tập phổ biến đóng thỏa mãn ngưỡng minsupp đã cho sử dụng thuật toán Apriori.
- b) Tìm các luật kết hợp có dạng sau và thỏa mãn ngưỡng minsupp, minconf đã cho sử dụng thuật toán Apriori
- item1 & item 2 -> item 3 & item 4 (về trái và phải của luật đều có 2 hạng mục)
  - D -> item (về phải có một hạng mục khác với hạng mục D)
- Yêu cầu trình bày chi tiết các bước (không chỉ liệt kê tập luật tìm được)

Câu 2: Cho tập dữ liệu gồm 7 điểm trong không gian 2 chiều : P1, P2, P3, P4, P5, P6, P7. Cho ma trận khoảng cách giữa các điểm như trong bảng 1.

- a) Hãy sử dụng lần lượt thuật toán AGNES với Single link và Complete link để gom nhóm (trình bày chi tiết các bước). Vẽ sơ đồ hình cây (dendogram) cho kết quả gom nhóm. (Sơ đồ hình cây phải vẽ rõ ràng để nhận biết được thứ tự và giá trị của vị trí các NHÓM gộp lại với nhau.)
- b) Dựa trên sơ đồ hình cây tương ứng (dùng Single Link/ Complete Link) xác định 3 nhóm thu được. So sánh kết quả.

Bảng 1 . Ma trận khoảng cách cho Câu 2

	P1	P2	P3	P4	P5	P6	P7
P1	0.00	0.27	0.23	0.56	0.17	0.40	0.14
P2	0.27	0.00	0.06	0.75	0.33	0.25	0.26
P3	0.23	0.06	0.00	0.59	0.28	0.24	0.22
P4	0.56	0.75	0.59	0.00	0.44	0.48	0.46
P5	0.17	0.33	0.28	0.44	0.00	0.37	0.09
P6	0.40	0.25	0.24	0.48	0.37	0.00	0.31
P7	0.14	0.26	0.22	0.46	0.09	0.31	0.00

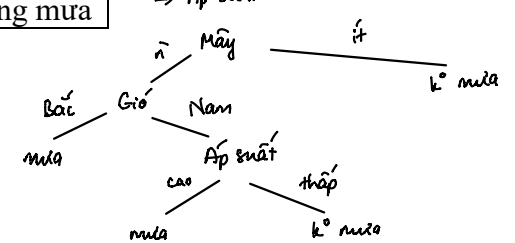
Câu 3: Sử dụng phương pháp cây quyết định để tìm các luật phân lớp từ bảng dữ liệu sau. Giả sử thuộc tính “kết quả” là thuộc tính phân lớp.

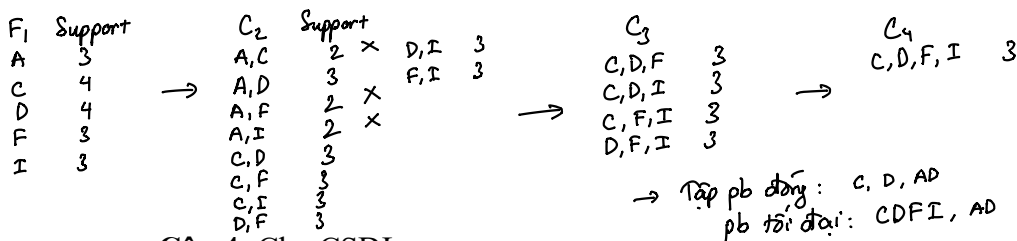
Đối tượng	Mây	Áp suất	Gió	Kết quả
1	ít	cao	Bắc	không mưa
2	nhiều	cao	Nam	mưa
3	nhiều	trung bình	Bắc	mưa
4	ít	thấp	Bắc	không mưa
5	nhiều	thấp	Bắc	mưa
6	nhiều	cao	Bắc	mưa
7	nhiều	thấp	Nam	không mưa
8	ít	cao	Nam	không mưa

$$Gain(S, Mây) = 1 - \frac{3}{8} \left( -\frac{3}{3} \log_2 \frac{3}{3} \right) - \frac{5}{8} \left( -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} \right) = 0.549$$

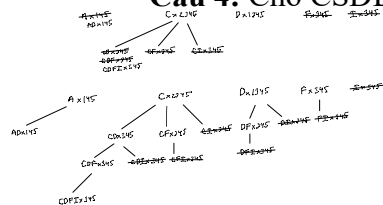
$$Gain(S, Gió) = 0.049$$
$$Gain(S, Áp suất) = 0.156$$

⇒ Mây  
 $Entropy(S, ít) = 0$   
 $Gain(S, nhiều, Áp suất) = 0.722 - 0.4 = 0.322$   
 $Gain(S, nhiều, Gió) = 0.722 - 0.4 = 0.322$   
⇒ Gió  
 $Gain(S, nhiều, Nam, Áp suất) = 1 - 0 = 1$   
⇒ Áp suất





Câu 4: Cho CSDL sau



TID	A	<del>B</del>	C	D	<del>E</del>	F	<del>G</del>	<del>H</del>	I
10	1			1			1	1	
20			1		1				
30		1	1	1		1			1
40	1		1	1	1	1	1		1
50	1		1	1		1		1	1

- a) Hãy sử dụng **một** trong hai thuật toán : **Apriori** hoặc **FP-Growth** để tìm **tất cả** các tập phổ biến thỏa mãn ngưỡng **minsupp=60%**. Liệt kê các tập phổ biến tối đại và tập bao phổ biến.
- b) Tìm các luật kết hợp được xây dựng từ tập phổ biến tối đại, thỏa mãn ngưỡng **minconf =80%**.

Câu 5: Cho CSDL sau :

Đối tượng	Mây	Áp suất	Gió	Kết quả
1	ít	cao	Bắc 0	không mưa
2	nhiều	cao	Nam	mưa
3	ít	thấp	Bắc	không mưa
4	nhiều	trung bình	Bắc	mưa
5	nhiều	thấp 0	Nam	không mưa
6	nhiều	thấp	Bắc	mưa
7	ít	cao	Nam	không mưa
8	nhiều	cao	Bắc	mưa

- L1: Nếu mây là ít thì kết quả là không mưa  
L2: Nếu áp suất thấp, gió nam thì kết quả là mưa  
L3: Nếu gió Bắc thì kết quả là không mưa  
L4: Nếu mây nhiều, áp suất cao thì kết quả là mưa  
L5: Nếu mây nhiều, áp suất thấp, gió Bắc thì mưa

- a) Sử dụng **thuật toán ILA** để tìm các luật phân lớp với cột **"Kết quả"** là thuộc tính phân lớp. Sử dụng **bộ luật phân lớp tìm được để xác định lớp cho các đối tượng mới** :

Đối tượng	Mây	Áp suất	Gió	Kết quả
9	ít	trung bình	Bắc	?
10	ít	thấp	Nam	?
11	nhiều	trung bình	Nam	?

ILA DT  
k° mưa k° mưa  
k° mưa k° mưa  
mưa X

- b) Sử dụng thuật toán **cây quyết định** để tìm các luật phân lớp với cột **"Kết quả"** là thuộc tính phân lớp. Sử dụng **bộ luật phân lớp tìm được để xác định lớp cho các đối tượng mới ở trên và so sánh kết quả với câu a)**. DT k° xác định lớp cho mẫu 11

Câu 6: Cho CSDL sau :

Đối tượng	Mây	Áp suất	Gió	Kết quả
1	ít	cao	Bắc	không mưa
2	nhiều	cao	Bắc	mưa
3	ít	thấp	Bắc	không mưa
4	nhiều	thấp	Bắc	mưa
5	nhiều	trung bình	Bắc	mưa
6	ít	cao	Nam	không mưa
7	nhiều	cao	Nam	mưa
8	nhiều	thấp	Nam	không mưa

Sử dụng thuật toán **Naïve Bayes** để xác định lớp cho mẫu mới sau:

$$\begin{aligned}
 X_1: S("k^{\circ} \text{ mưa}") &= P("k^{\circ} \text{ mưa}").R_1("ít", "k^{\circ} \text{ mưa}").R_2("thấp", "k^{\circ} \text{ mưa}").R_3("Nam", "k^{\circ} \text{ mưa}") \\
 &= \frac{4+1}{8+2} \cdot \frac{3+1}{4+2} \cdot \frac{2+1}{4+3} \cdot \frac{2+1}{4+2} = 0.071 \\
 S("có \text{ mưa}") &= P("có \text{ mưa}").R_1(\dots) \dots \\
 &= \frac{4+1}{8+2} \cdot \frac{0+1}{4+2} \cdot \frac{1+1}{4+3} \cdot \frac{1+1}{4+2} = 0.008
 \end{aligned}$$

$$X_2: S("k^o mua") = \frac{4+1}{8+2} \cdot \frac{3+1}{4+2} \cdot \frac{0+1}{4+3} \cdot \frac{2+1}{4+2} = 0.024$$

$$S("có mua") = \frac{4+1}{8+2} \cdot \frac{0+1}{4+2} \cdot \frac{1+1}{4+3} \cdot \frac{3+1}{4+2} = 0.016$$

$$X_3: S("k^o mua") = \frac{4+1}{8+2} \cdot \frac{1+1}{4+2} \cdot \frac{2+1}{4+3} \cdot \frac{2+1}{4+2} = 0.036$$

$$S("có mua") = \frac{4+1}{8+2} \cdot \frac{4+1}{4+2} \cdot \frac{2+1}{4+3} \cdot \frac{3+1}{4+2} = 0.119$$

$$X_4: S("k^o mua") = \frac{4+1}{8+2} \cdot \frac{3+1}{4+2} \cdot \frac{0+1}{4+3} \cdot \frac{2+1}{4+2} = 0.024$$

$$S("có mua") = \frac{4+1}{8+2} \cdot \frac{0+1}{4+2} \cdot \frac{1+1}{4+3} \cdot \frac{1+1}{4+2} = 0.008$$

Đối tượng	Mây	Áp suất	Gió	Kết quả
9	ít	thấp	Nam	?
10	ít	trung bình	Bắc	?
11	nhiều	cao	Bắc	?
12	nhiều	trung bình	Nam	?

*k<sup>o</sup> mua*

*k<sup>o</sup> mua*

*có mua*

*k<sup>o</sup> mua*

**Câu 7:** Cho bảng dữ liệu thống kê kết quả của một thuật toán phân lớp số khách hàng đến siêu thị có mua hay không mua sản phẩm trong 1 tháng:

*Lớp dự đoán*

<i>Lớp thực</i>	<i>Lớp</i>	Mua	Không mua
<i>Sự</i>	Mua	8986	1009
	Không mua	1358	2547

$$\text{confusion matrix} = \begin{bmatrix} 8986 & 1009 \\ 1358 & 2547 \end{bmatrix}$$

- Lập ma trận sai số (confusion matrix)  $\text{Accuracy} = \frac{8986 + 2547}{8986 + 1009 + 1358 + 2547} = 82.97\%$

- Tính các độ đo accuracy, error rate, sensitivity, specificity, precision

$$\text{Precision} = \frac{8986}{8986 + 1358} = 86.87\%$$

$$\text{Error rate} = 1 - \text{accuracy} = 17.03\%$$

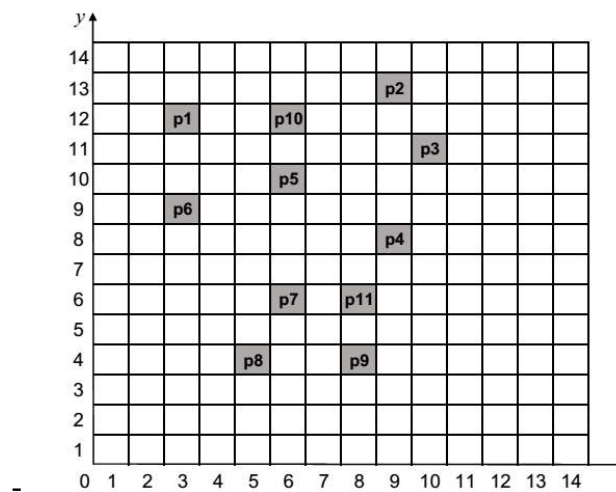
$$\text{Sensitivity} = 1 - \frac{1009}{1009 + 8986} = 89.9\%$$

$$\text{Specificity} = \frac{2547}{2547 + 1358} = 65.22\%$$

**Câu 8:** Cho các mẫu dữ liệu được phân bố trong không gian hai chiều Oxy như hình vẽ 1 (trang sau). Ví dụ: điểm P1 ở tọa độ (3,12). Giả sử người ta tiến hành gán nhãn cho mỗi điểm như sau:

*p1: xanh, p2: xanh, p3: đỏ, p4: xanh, p5: đỏ, p6: xanh, p7: đỏ, p8: đỏ, p9: xanh.*

Sử dụng thuật toán k-NN với khoảng cách Euclide để phân lớp 2 mẫu sau: p10, p11 với số lân cận  $k = 3$ . Thể hiện việc tính toán đầy đủ.



**Hình 1:** Phân bố các điểm dữ liệu trong không gian Oxy

Gợi ý: Công thức Euclide của 2 điểm A, B trong không gian Oxy:

$$AB = \sqrt{(x_A - x_B)^2 + (y_A - y_B)^2}$$

**Câu 9:** Cho tập dữ liệu gồm 12 giá trị như bên dưới (đã sắp xếp theo thứ tự tăng dần).

5, 10, 11, 13, 15, 35, 50, 55, 72, 92, 204, 215

a. Hãy áp dụng phương pháp chia giỏ để chia dữ liệu thành **3 giỏ** bằng hai phương pháp:

– Chia giỏ theo độ rộng  $[5, 75) \quad [75, 145) \quad [145, 215]$

– Chia giỏ theo độ sâu  $[5, 13] \quad [15, 55] \quad [72, 215]$

b. Áp dụng làm tròn bằng giá trị trung bình, giá trị trung vị và biên giỏ cho trường hợp

chia giỏ theo độ sâu. *Means:*  
Bin 1: 9.75, 9.75, ... Bin 3: 145.75, ... *Boundaries:*  
Bin 1: 5, 13, 13, 13 Bin 3: 72, 72, 215, 215  
Bin 2: 38.75, ... Bin 2: 15, 15, 55, 55

**Câu 10:** Cho tập dữ liệu gồm 8 điểm trong không gian 2 chiều: A1=(2,10), A2=(2,5), A3=(8,4), A4=(5,8), A5=(7,5), A6=(6,4), A7=(1,2), A8=(4,9).

Hãy sử dụng lần lượt thuật toán **DBSCAN** để gom nhóm với Eps = 2 và Minpts = 2.