

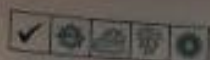
ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN



KHAI THÁC DỮ LIỆU

BÀI TẬP LAP 3

NGUYỄN THỊ NGỌC HÀ: 17520421



Date: _____

E A B D E C A D C E B A B D C →
 10 20 30 40 50 60 70 80 90

min frequent = 30%, min conf = 60%

width = 25 và step = 5

a. Đây phổ biến song song bằng thuật toán ~~WINEPT~~ WINEPT

- Thuật toán này sử dụng của số với độ rộng cố định để trượt qua chuỗi sk...

$$N = \frac{T_c - T_s + \text{width}}{\text{step}} - 1 =$$

$$\Rightarrow N = \frac{100 - 5 + 25}{5} - 1 = 23$$

Các của số

U_1 : -, -, -, -, E

U_{13} : D, A, -, C, E, B

U_2 : -, -, -, E, A

U_{14} : -, C, E, B, A

U_3 : -, -, E, A, B

U_{15} : C, E, B, A, -

U_4 : -, E, A, B, D

U_{16} : E, B, A, -, B

U_5 : E, A, B, D, E

U_{17} : B, A, -, B, -

U_6 : A, B, D, E, C

U_{18} : A, -, B, -, D

U_7 : B, D, E, C, -

U_{19} : -, B, -, D, C

U_8 : D, E, C, -, A

U_{20} : B, -, D, C, -

U_9 : E, C, -, A, D

U_{21} : -, D, C, -, -

U_{10} : C, -, A, D, -

U_{22} : D, C, -, -, -

U_{11} : -, A, D, -, C

U_{23} : C, -, -, -, -

U_{12} : A, D, -, C, E



Tập các episode song song ứng viên 1 phần tử

A: 15 C: 15 E: 15

B: 13 D: 15

Với số của số là 23 và min-frequent = 30%,
ta có các dãy song song 1 phần tử phổ biến sau:

$$L_1 = \{ \{A\}, \{B\}, \{C\}, \{D\}, \{E\} \}$$

• Tập ứng viên có 2 phần tử là:

$$C_2 = L_1 \times L_1 = \{ AB: 9, AC: 8, AD: 9, AE: 11$$

$$BC: 7, BD: 8, BE: 9, CD: 12, CE: 8, DE: 8 \}$$

⇒ Các dãy phổ biến 2 phần tử thỏa min-frequent:

$$L_2 = \{ \{AB\}, \{AC\}, \{AD\}, \{AE\}, \{BC\}, \{BD\}, \\ \{BE\}, \{CD\}, \{CE\}, \{DE\} \}$$

• Tập ứng viên có 3 phần tử là:

$$C_3 = L_2 \times L_2 = \{ ABC: 3, ABD: 4, ABE: 7, ACD: 6,$$

$$ACE: 6, ADE: 6, BCD: 5, BCE: 5, BDE: 5,$$

$$CDE: 6 \}$$

⇒ Các dãy phổ biến 3 phần tử thỏa min-frequent:

$$L_3 = \{ \{ABE\} \}$$

• Tập ứng viên có 4 phần tử:

$$C_4 = L_3 \times L_3 = \emptyset$$

Đến đây thuật toán tìm dãy phổ biến dừng lại và kết quả là các dãy L_1, L_2, L_3 thu được.



Date:

b. Dãy phổ biến tuân từ bằng thuật toán WINEPI

• Sử dụng dữ liệu của sở trượt và dãy phổ biến

L_1 đã tìm ở câu a.

• Tập các ứng viên có 2 phôi là:

$$C_2 = L_1 \bowtie L_1 = \{ AB:7, BA:4, AC:3, CA:5, AD:8, DA:1, AE:3, EA:9, BC:4, CB:3, BD:7, DB:1, BE:3, EB:7, CD:2, DC:10, CE:4, EC:4, DE:6, ED:3 \}$$

\Rightarrow Các dãy phổ biến 2 phôi thỏa min-frequent:

$$L_2 = \{ \{A, B\}, \{A, D\}, \{E, A\}, \{B, D\}, \{E, B\}, \{D, C\} \}$$

• Tập các ứng viên có 3 phôi là:

$$C_3 = L_2 \bowtie L_2 = \{ ABD:4, EAB:3 \}$$

\Rightarrow Không có dãy phổ biến 3 phôi thỏa min-frequent

• Đến đây thuật toán tìm dãy phổ biến dừng lại và kết quả là các dãy L_1, L_2 thu được

c. Các luật WI NEPI song song dựa trên dãy phổ biến tìm được ở câu a là:

Xét tập phổ biến ABE, tính confidence của các luật tạo ra:

$$A \Rightarrow ABE \text{ có } \text{conf} = \frac{\text{fre}(ABE)}{\text{fre}(A)} = \frac{7}{15} = 0,46 < \text{min-conf} = 0,6 \quad (1)$$

$$B \Rightarrow ABE \text{ có } \text{conf} = \frac{\text{fre}(ABE)}{\text{fre}(B)} = \frac{7}{13} = 0,54 < \text{min-conf} = 0,6 \quad (2)$$

$$E \Rightarrow ABE \text{ có } \text{conf} = \frac{\text{fre}(ABE)}{\text{fre}(E)} = \frac{7}{15} = 0,46 < \text{min-conf} = 0,6 \quad (3)$$



Date: _____

$$(4) AB \Rightarrow ABE \text{ có } \text{conf} = \frac{f_e(ABE)}{f_e(AB)} = \frac{7}{9} = 0,77 > \text{min-conf} = 0,6$$

$$(5) BE \Rightarrow ABE \text{ có } \text{conf} = \frac{f_e(ABE)}{f_e(BE)} = \frac{7}{9} = 0,77 > \text{min-conf} = 0,6$$

$$(6) AE \Rightarrow ABE \text{ có } \text{conf} = \frac{f_e(ABE)}{f_e(AE)} = \frac{7}{11} = 0,64 > \text{min-conf} = 0,6$$

Như vậy, ta nhận các luật (4), (5), (6)

d. Các luật WINEPI tuân hệ dựa trên dãy phổ biến tìm được

$$(1) A \Rightarrow AB \text{ có } \text{conf} [0,30, 0,46]$$

$$(2) A \Rightarrow AD [0,35, 0,53]$$

$$(3) A \Rightarrow EA [0,39, 0,60]$$

$$(4) B \Rightarrow AB [\cancel{0,37, 0,47}] [0,3, 0,54]$$

$$(5) D \Rightarrow AD [0,35, 0,53]$$

$$(6) E \Rightarrow EA [0,39, 0,60]$$

$$(7) B \Rightarrow BD [0,3, 0,54]$$

$$(8) D \Rightarrow BD [0,3, 0,47]$$

$$(9) E \Rightarrow EB [0,3, 0,46]$$

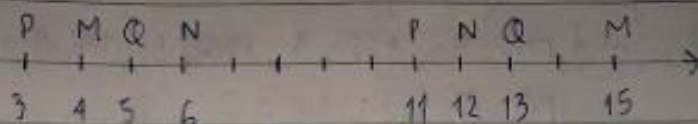
$$(10) B \Rightarrow EB [0,3, 0,54]$$

$$(11) D \Rightarrow DC [0,43, 0,66]$$

$$(12) C \Rightarrow DC [0,43, 0,66]$$

Vậy ta nhận các luật (3), (6), (11), (12)

2.



min-frequent = 20%, min-conf = 50%, width = 4, skip = 1
 a. Các dãy phổ biến song song bằng thuật toán WINEPT

$$N = \frac{T_e - T_s + \text{width}}{\text{step}} - 1 = \frac{16 - 3 + 4}{1} - 1 = 16$$

Các cửa sổ:

$U_1: _, _, _, P$

$U_9: _, _, _, P$

$U_2: _, _, P, M$

$U_{10}: _, _, P, N$

$U_3: _, P, M, Q$

$U_{11}: _, P, N, Q$

$U_4: P, M, Q, N$

$U_{12}: P, N, Q, _$

$U_5: M, Q, N, _$

$U_{13}: N, Q, _, M$

$U_6: Q, N, _, _$

$U_{14}: Q, _, M, _$

$U_7: N, _, _, _$

$U_{15}: _, M, _, _$

$U_8: _, _, _, _$

$U_{16}: M, _, _, _$

• Tập các episode song song 1phủ

$M: 8 \quad N: 8 \quad P: 8 \quad Q: 8$

⇒ Với 16 cửa sổ và min-frequent = 20%, ta tìm ra các dãy song song 1phủ phổ biến sau:

$$L_1 = \{ \{M\}, \{N\}, \{P\}, \{Q\} \}$$



Date:

- Tập các ứng viên có 2 chữ là:

$$C_2 = L_1 \bowtie L_1 = \{ MN:3, MP:3, MQ:5, NP:4, NQ:6, PQ:4 \}$$

→ Các dãy phổ biến 2 chữ thỏa min-frequent:

$$L_2 = \{ \{M, Q\}, \{N, P\}, \{N, Q\}, \{P, Q\} \}$$

- Tập các ứng viên có 3 chữ là:

$$C_3 = L_2 \bowtie L_2 = \{ NPQ:3 \}$$

⇒ Không có dãy phổ biến 3 chữ thỏa min-frequent.

- Đến đây, thuật toán tìm dãy phổ biến dừng lại và kết quả là các dãy L_1, L_2 thu được.

b. Các dãy phổ biến tuân từ bằng thuật toán WINEPI

- Sử dụng dữ liệu của số trượt và dãy phổ biến L_1 đã làm ở câu A.

- Tập các ứng viên có 2 chữ là:

$$C_2 = L_1 \bowtie L_1 = \{ MN:2, NM:0, MP:0, PM:3, MQ:3, QM:2, NP:0, PN:4, NQ:3, QN:3, PQ:4, QP:0 \}$$

→ Các dãy phổ biến 2 chữ thỏa min-frequent:

$$L_2 = \{ \{P, N\}, \{P, Q\} \}$$

- Tập các ứng viên có 3 chữ

$$C_3 = L_2 \bowtie L_2 = \emptyset$$

→ Không có dãy phổ biến 3 chữ

- Đến đây thuật toán tìm dãy phổ biến dừng lại và kết quả là các dãy L_1, L_2 thu được.



Date: _____

c. Các luật WINEPI song song dựa trên dãy phối biến tìm được ở câu a.

$$M \Rightarrow MQ \quad [0,31; 0,63] \quad (1)$$

$$Q \Rightarrow MQ \quad [0,31; 0,63] \quad (2)$$

$$N \Rightarrow NP \quad [0,25; 0,5] \quad (3)$$

$$P \Rightarrow NP \quad [0,25; 0,5] \quad (4)$$

$$N \Rightarrow NQ \quad [0,38; 0,75] \quad (5)$$

$$Q \Rightarrow NQ \quad [0,38; 0,75] \quad (6)$$

$$P \Rightarrow PQ \quad [0,25; 0,5] \quad (7)$$

$$Q \Rightarrow PQ \quad [0,25; 0,5] \quad (8)$$

Với $\text{min_conf} = 50\%$, ta nhận tất cả các luật nêu trên.

d. Các luật WINEPI tuân từ dựa trên dãy phối biến câu b.

$$(1) P \Rightarrow PN \quad [0,25; 0,5] \quad (3) P \Rightarrow PQ \quad [0,25; 0,5]$$

$$(2) N \Rightarrow PN \quad [0,25; 0,5] \quad (4) Q \Rightarrow PQ \quad [0,25; 0,5]$$

Với $\text{min_conf} = 50\%$, ta nhận tất cả các luật nêu trên.

3. Cho dữ liệu mua văn phòng phẩm sau.

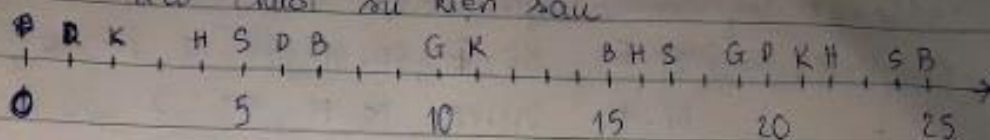
| Tên | Thời điểm mua |
|--------------|---------------|
| Bút bi | 7, 15, 25 |
| Bìa sổ mi | 5, 17, 24 |
| Hồ dán | 4, 16, 22 |
| Đĩa CD | 1, 6, 20 |
| Giấy ghi chú | 10, 19 |
| Kẹp giấy | 2, 11, 21 |

| | | |
|------|---------------|--------------------|
| Gọi: | Bút bi = B | Với: |
| | Bìa sổ mi = S | min frequent = 50% |
| | Hồ dán = H | min conf = 70% |
| | Đĩa CD = D | width = 6 |
| | Giấy ghi = G | step = 1 |
| | Kẹp giấy = K | |



Date: _____

Tạo chuỗi sự kiện sau



$$N = \frac{T_e - T_s + \text{width}}{\text{Step}} - 1 = 30$$

| | | |
|--|--------------------------------------|--------------------------------------|
| U_1 : -, -, -, -, -, D | U_{11} : D, B, -, -, G, K | U_{21} : H, S, -, -, G, D, K |
| U_2 : -, -, -, -, -, D, K | U_{12} : B, -, -, -, G, K, - | U_{22} : S, -, -, -, G, D, K, H |
| U_3 : -, -, -, -, -, D, K, - | U_{13} : -, -, -, -, G, K, - | U_{23} : -, -, -, -, G, D, K, H, - |
| U_4 : -, -, -, -, -, D, K, -, H | U_{14} : -, -, -, -, G, K, -, - | U_{24} : G, D, K, H, -, S |
| U_5 : -, -, -, -, -, D, K, -, H, S | U_{15} : G, K, -, -, -, B | U_{25} : D, K, H, -, S, B |
| U_6 : D, K, -, -, H, S, D | U_{16} : K, -, -, -, B, H | U_{26} : K, H, -, -, S, B, - |
| U_7 : K, -, -, H, S, D, B | U_{17} : -, -, -, -, B, H, S | U_{27} : H, -, -, -, S, B, -, - |
| U_8 : -, -, H, S, D, B, - | U_{18} : -, -, -, -, B, H, S, - | U_{28} : -, -, -, -, S, B, -, -, - |
| U_9 : H, S, D, B, -, - | U_{19} : -, -, -, -, B, H, S, -, G | U_{29} : S, B, -, -, -, -, - |
| U_{10} : S , D , B, -, -, -, G | U_{20} : B, H, S, -, -, G, D | U_{30} : B, -, -, -, -, -, - |

• Tập các episodes song song 1phần:

$$C_1 = \{ B: 18, D: 17, K: 18, H: 18, S: 18, G: 12 \}$$

Với 50 của số và nín - frequent = 30%, ta tìm ra các

dãy phổ biến song song 1phần sau:

$$L_1 = \{ \{B\}, \{D\}, \{K\}, \{H\}, \{S\}, \{G\} \}$$



Date: _____

- Tập các ứng viên có 2ptử là:

$$C_2 = L_1 \times L_1 = \{ BD:5, BK:7, BH:11, BS:12, BG:6, DK:12, DH:12, DS:11, DG:7, SG:6, KH:11, KS:8, KG:9, HS:15, HG:6 \}$$

- ⇒ Các dãy phổ biến 2ptử thỏa min-frequent:

$$L_2 = \{ \{B,H\}, \{B,S\}, \{D,K\}, \{D,H\}, \{D,S\}, \{K,H\}, \{K,G\}, \{H,S\} \}$$

- Tập các ứng viên có 3ptử là:

$$C_3 = L_2 \times L_2 = \{ DKH:9, BHS:10, DHS:10 \}$$

- ⇒ Các dãy phổ biến 3ptử thỏa min-frequent:

$$L_3 = \{ \{D,K,H\}, \{B,H,S\}, \{D,H,S\} \}$$

- Tập các ứng viên có 4ptử là

$$C_4 = L_3 \times L_3 = \emptyset$$

- ⇒ Không có dãy phổ biến 4ptử

- Đến đây thuật toán tìm dãy phổ biến dừng lại và kết quả là L_1, L_2, L_3 . Thu được.

b. Các dãy phổ biến tuân hệ bằng thuật toán WINEPI

- Sử dụng dữ liệu của số trượt và dãy phổ biến L_1 đã làm ở câu a

- Tập các ứng viên có 2ptử là:

$$C_2 = L_1 \times L_1 = \{ BD:1, DB:6, BK:2, KB:5, BH:5, HB:6, BS:4, SB:9, BG:5, GB:1, DK:11, KD:2, PH:7, HP:6, DS:4, SD:8, DG:2, GD:5, SG:5, GS:1, KH:10, HK:1, KS:6, SK:2, KG:0, GK:9, HS:14, SH:1, HG:3, GH:3 \}$$



Date:

⇒ Các dãy phổ biến & phổ biến min-afrequent

$$L_2 = \{ \{S, B\}, \{P, K\}, \{K, H\}, \{G, K\}, \{H, S\} \}$$

• Tập các ứng viên 3phần

$$C_3 = L_2 \times L_2 = \emptyset$$

⇒ Không có dãy phổ biến 3phần

• Đến đây thuật toán tìm dãy phổ biến dừng lại và kết quả là L_1, L_2 thu được.

c. Các luật WI NEPI song song dựa trên dãy phổ biến tìm được ở câu a.

$$(1) D \Rightarrow DKH \quad [0,3, 0,53] < \text{min-conf} = 0,7$$

$$(2) K \Rightarrow DKH \quad [0,3, 0,5] < \text{min-conf} = 0,7$$

$$(3) H \Rightarrow DKH \quad [0,3, 0,5] < \text{min-conf} = 0,7$$

$$(4) KH \Rightarrow DKH \quad [0,3, 0,82] > 0,7$$

$$(5) PH \Rightarrow PKH \quad [0,3, 0,75] > 0,7$$

$$(6) KD \Rightarrow DKH \quad [0,3, 0,75] > 0,7$$

Với $\text{min-conf} = 70\%$, ta nhận được các luật (1), (5), (6)

Lâm tương tự với tập phổ biến DHS và BHS, ta được

$$DH \Rightarrow DHS \quad [0,33, 0,83] > 0,7$$

$$DS \Rightarrow DHS \quad [0,33, 0,9] > 0,7$$

$$BH \Rightarrow BHS \quad [0,33, 0,9] > 0,7$$

$$BS \Rightarrow BHS \quad [0,33, 0,83] > 0,7$$

...



Date: _____

d. Các luật ~~phổ biến~~ WINEPI tuân từ thỏa min-conf dựa trên đây phổ biến ở câu b.

$$(1) S \Rightarrow SB [0,3, 0,5] \quad (6) K \Rightarrow KH [0,33, 0,56]$$

$$(2) B \Rightarrow SB [0,3, 0,5] \quad (7) H \Rightarrow KH [0,33, 0,56]$$

$$(3) D \Rightarrow DK [0,32, 0,65] \quad (8) G \Rightarrow GK [0,3, 0,5]$$

$$(4) K \Rightarrow DK [0,32, 0,61] \quad (9) K \Rightarrow GK [0,3, 0,5]$$

$$(5) H \Rightarrow HS [0,47, 0,78] \quad (10) S \Rightarrow HS [0,47, 0,78]$$

Với min-conf = 70%, ta nhận được các luật (5), (10).