

Додатне технике правила придруживања

Ненад Митић

Математички факултет
nenad@matf.bg.ac.rs

Непрекидни и категорички атрибути

- Претходно приказани алгоритми су применљиви на податке представљене у облику асиметричних бинарних атрибута
- Шта ако су подаци представљени категоричким, непрекидним или симетричним бинарним атрибутима
- Потребне су модификације алгоритама

Индекс	Име	Презиме	Пол	Назив програма	Назив предмета	Статус	Датум полагања	Оцена
20150400	Ања	Јевтић	ж	Информатика	Алгебра 1	о	2016-02-19	10
20150220	Урош	Рачић	м	Математика	Геометрија 2	о	2016-02-03	10
20150220	Урош	Рачић	м	Математика	Програмирање 1	о	2016-02-02	10
20150220	Урош	Рачић	м	Математика	Страни језик	о	2016-01-30	10
...

Непрекидни и категорички атрибути

- Поступак са категоричким атрибутима
 - Трансформисати категоричке атрибуте у асиметричне бинарне променљиве
 - Могући проблеми
 - Шта ако атрибут има велики број могућих вредности (нпр. атрибут држава има више од 200 могућих вредности)
 - Велики број вредности атрибута може да има малу подршку
 - Могуће решење: агрегирати вредности таквих атрибута

Непрекидни и категорички атрибути

- Шта ако је расподела вредности атрибута јако одскаче на једну страну?
- Нпр. 95% студената има за оцену 10 вредност предмета "Страни језик"
- Могуће решење: изbrisati ставке са високом учесталошћу
- Поступак са непрекидним атрибутима?
- Различите врсте правила

[STATUS=0]+[OCENA=6]+[OCENA=7]+[DATPOLAGANJA >= 736950 AND < 737450]

→ [NAZIV_PREDMETA=Програмирање 2]

Поступак са непрекидним атрибутима

Различите методе

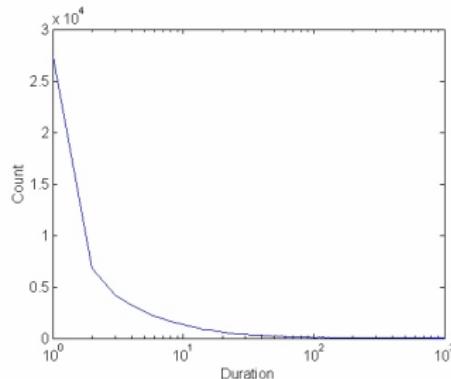
- Засноване на дискретизацији
- Статистички засноване методе
- Засноване на не-дискретизацији
 - *minApriori*

Методе засноване на дискретизацији

Дискретизација

- Без надзора

- Делови исте ширине
- Делови исте дубине
- Кластеровање



- Под надзором

На основу вредности атрибута A

Класа	A ₁	A ₂	A ₃	A ₄	A ₅	A ₆	A ₇	A ₈	A ₉
Аномалија	0	0	20	10	20	0	0	0	0
Нормалних	150	100	0	0	0	100	100	150	100

{ deo1 } { deo2 } { deo3 }

Проблеми код дискретизације

- Величина дискретизационих интервала утице на подршку и поузданост
- Ако је интервал јако мали подршка може да буде недовољна
- Ако је интервал јако велики поузданост може да буде недовољна
- Решење: користити све могуће интервале (ефикасност?)

Проблеми код дискретизације

- Време извршавања
 - Ако интервал садржи p вредности, у просеку постоји $O(p^2)$ могућих уређења
- Велики број правила
- Илустрација проблема: правила придрживања над табелом *polozeni_ispit*

Статистички засноване методе

- Primer: [status=o]+[datpolaganja] >= 736950 and < 737450] →
[оценка]: $\mu = 7.5$
- Последично правило - непрекидна променљиве
карактерисана статистиком (средина, медијана, stddev, ...)
- Приступ
 - Издвајање циљног атрибута из остатка података
 - Применити постојеће генерисање честог скупа ставки на
остатак података
 - За сваку честу ставку израчунати описну статистику за
одговарајућу циљну променљиву
 - Чест скуп ставки постаје правило за укључивање циљног
атрибута као последичног правила
 - Применити статистичке тестове ради одређивања
интересантности правила

Статистички засноване методе

- Како одредити да ли је правило придруживања интересантно?
- Поредити статистику дела популације покривену правилом у односу на део популације који није покривен правилом:
 $A \rightarrow B : \mu$ prema $\bar{A} \rightarrow B : \mu'$
- Статистичко тестирање хипотеза:
 - Нулта хипотеза: $H_0 : \mu' = \mu + \Delta$
 - Алтернативна хипотеза: $H_1 : \mu' > \mu + \Delta$

Статистички засноване методе

Да би се одредило која је хипотеза важећа рачуна се Z статистика

$$Z = \frac{\mu' - \mu - \Delta}{\sqrt{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)}}$$

где је

- n_1 број трансакција које подржавају А, n_2 број трансакција које не подржавају А
- s_1 стандардна девијација циљног атрибута у трансакцијама које подржавају А
- s_2 стандардна девијација циљног атрибута у трансакцијама које не подржавају А

Према нултој хипотези Z има средину 0 и варијансу 1. Ако је $Z > Z_0$ где је Z_0 критична вредност за одређен ниво поузданости, тада се нулта хипотеза одбацује

Статистички засноване методе

Пример: нека је правило

[Назив_предмета='Анализа 1']+ [ознакарока='Јун'] → [поени]: $\mu = 48$

подржано од стране 50 студената, и нека је стандардна девијација броја њихових поена 3.5. Са друге стране, за комплементаран скуп од 200 студената који не подржавају ово правило средња вредност броја поена је 55, а стандардна девијација је 6.5. Нека је правило интересантно ако је разлика између μ' и μ већа од 5 поена ($\Delta = 5$). Вредност Z је

$$Z = \frac{55 - 48 - 5}{\sqrt{\left(\frac{3.5^2}{50} + \frac{6.5^2}{200}\right)}} = \frac{2}{\sqrt{(0.245 + 0.21125)}} = \frac{2}{\sqrt{0.45625}} = \frac{2}{0.675462804} = 2,960932842$$

За 1-страни тест са поузданошћу од 95% критична вредност за одбацивање нулте хипотезе 1.64. Како је $Z > 1.64$ нулта хипотеза се одбацује ==> правило је интересантно

Методе засноване на не-дискретизацији

- Постоје случајеви када је интереснатније наћи везе између непрекидних атрибута него између њихових дискретних интервала
- Пример: на основу табеле појављивања речи у тексту може се закључити да W_1 и W_2 имају тенденцију да се појављују заједно у истом документу

	W_1	W_2	W_3	W_4	W_5
D_1	2	2	0	0	1
D_2	0	0	1	2	2
D_3	2	3	0	0	0
D_4	0	0	1	0	1
D_5	1	1	1	0	2

Методе засноване на не-дискретизацији

- Подаци садрже само непрекидне атрибуте истог "типа"
- Пример: фреквенција речи у неком документу

	W_1	W_2	W_3	W_4	W_5
D_1	2	2	0	0	1
D_2	0	0	1	2	2
D_3	2	3	0	0	0
D_4	0	0	1	0	1
D_5	1	1	1	0	2

- Могуће решење:
 - Конвертовати садржај у 0/1 матрицу где је 1 нормализована вредност која прелази одређен праг и применити постојеће алгоритме (губи се информација о фреквенцији речи)
 - Дискретизација је често неприменљива пошто корисници желе везе између речи а не између броја појављивања речи

Min-Apriori

- Како одредити подршку за реч?
 - Ако се саберу фреквенције подршка ће бити већа од укупног броја докумената!
 - Нормализује се вектор речи, нпр. употребом L_1 норме
 - Свака реч има подршку једнаку 1.0

	W_1	W_2	W_3	W_4	W_5			W_1	W_2	W_3	W_4	W_5
D_1	2	2	0	0	1	→	D_1	0.40	0.33	0.00	0.00	0.17
D_2	0	0	1	2	2	→	D_2	0.00	0.00	0.33	1.00	0.33
D_3	2	3	0	0	0	→	D_3	0.40	0.50	0.00	0.00	0.00
D_4	0	0	1	0	1	→	D_4	0.00	0.00	0.33	0.00	0.17
D_5	1	1	1	0	2	→	D_5	0.20	0.17	0.33	0.00	0.33

Min-Apriori

- Ставка је скуп речи
- Подршка представља меру колико су речи придружене једна другој
- Подршка скупа речи C у скупу докумената T се рачуна као

$$sup(C) = \sum_{i \in T} \min_{j \in C} D(i,j)$$

Пример: $sup(W_1, W_2, W_3) = 0 + 0 + 0 + 0 + 0.17 = 0.17$

	W_1	W_2	W_3	W_4	W_5
D_1	0.40	0.33	0.00	0.00	0.17
D_2	0.00	0.00	0.33	1.00	0.33
D_3	0.40	0.50	0.00	0.00	0.00
D_4	0.00	0.00	0.33	0.00	0.17
D_5	0.20	0.17	0.33	0.00	0.33

Min-Apriori

Ова мера подршке се назива *Min-Apriori* и има особине да подршка

- монтоно расте како расте нормализована фреквенција речи
- монтоно расте како расте број документата који садрже реч
- монотоно опада како расте број речи у скупу ставки - анти-монотоност

Пример

- $sup(W_1) = 0.4 + 0 + 0.4 + 0 + 0.2 = 1.0$
- $sup(W_1, W_2) = 0.33 + 0 + 0.4 + 0 + 0.17 = 0.9$
- $sup(W_1, W_2, W_3) = 0 + 0 + 0 + 0 + 0.17 = 0.17$

	W_1	W_2	W_3	W_4	W_5
D_1	0.40	0.33	0.00	0.00	0.17
D_2	0.00	0.00	0.33	1.00	0.33
D_3	0.40	0.50	0.00	0.00	0.00
D_4	0.00	0.00	0.33	0.00	0.17
D_5	0.20	0.17	0.33	0.00	0.33

Правила придрживања између више нивоа



Правила придрживања између више нивоа

- Зашто се укључује хијерархија концепата?
- Правила на нижим нивоима можда немају довољну подршку да се јаве у честим скуповима података
- Правила на нижим нивоима су превише специфична. Нпр. обрано млеко → бели хлеб, пуномасно млеко → црни хлеб, кефир → цри хлеб, ... су индикатори правила придрживања између млека и хлеба
- Ако се обилази дрво хијерархије концепата тада важи
 - 1 Ако је X родитељ ставка за X_1 и X_2 тада важи $\sigma(X) \leq \sigma(X_1) + \sigma(X_2)$
 - 2 Ако је $\sigma(X_1 \cup Y_1) \geq \text{minsup}$ и X је родитељ од X_1 и Y је родитељ од Y_1 тада $\sigma(X \cup Y_1) \geq \text{minsup}$, $\sigma(X_1 \cup Y) \geq \text{minsup}$, и $\sigma(X \cup Y) \geq \text{minsup}$
 - 3 Ако $\text{conf}(X_1 \rightarrow Y_1) \geq \text{minconf}$ тада $\text{conf}(X_1 \rightarrow Y) \geq \text{minconf}$

Правила придрживања између више нивоа

Приступ 1:

- Проширити текуће правило придрживања проширивањем сваке трансакције са ставком са вишег нивоа
 - Оригинална трансакције: {обрано млеко, бели хлеб}
 - Проширена трансакција: {обрано млеко, бели хлеб, млеко, хлеб, храна}
- Проблеми:
 - Ставке које се налазе на вишем нивоу имају много виши ниво подршке
 - ако је подршка јако мала велики број честих образца укључује ставке са вишег нивоа
 - Повећава се димензионалност података

Правила придрживања између више нивоа

Приступ 2:

- Формирати честе обрасце прво на највишем нивоу
- Затим формирати честе обрасце на наредним највишим нивоима, итд.

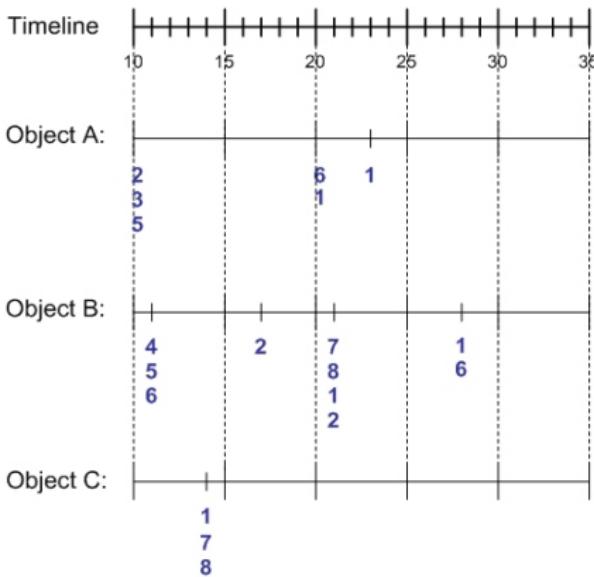
Проблеми:

- У/И захтеви се драматично повећавају због потребе вишеструког пролажења кроз податке
- Могу да се изгубе неке од потенцијално занимљивих веза између образца на различитим нивоима

Низ података

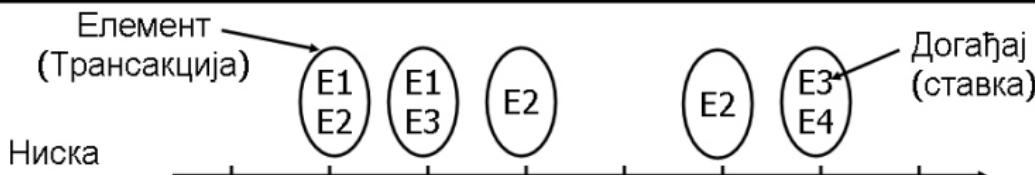
Низ у бази:

Object	Timestamp	Events
A	10	2, 3, 5
A	20	6, 1
A	23	1
B	11	4, 5, 6
B	17	2
B	21	7, 8, 1, 2
B	28	1, 6
C	14	1, 8, 7



Примери низа података

Niz u bazi	Niz podataka	Element (Transakcija)	Događaj (stavka)
Kupac	Istorija kupovanja datog kupca	Skup stavki kupljen od strane kupca u trenutku t	Knjige, tekući proizvodi, CDovi, itd.
Veb podaci	Pregledanje aktivnosti pojedinačnog posetioča Veba	Skup datoteka pregledan od strane poseocioca Veba posle pritiska na tipku miša	Glavna strana, strana sa indeksima, kontakt informacije, itd.
Podaci o događajima	Istorija događaja formirana pomoću datog senzora	Događaju uočeni senzorom u trenutku t	Tipovi alarma koje je formirao senzor
Niske genoma	DNK niske pojedinačnih vrsta	Elementi DNK niski	Baze A,T,G,C



Формална дефиниција ниске

- Ниска је уређена листа елемената (трансакција)
 $S = \langle e_1 e_2 e_3 \dots \rangle$
- Сваки елемент садржи скуп догађаја (ставки) $e_i = \{i_1 i_2 \dots i_k\}$
- Сваком елементу се додељује одређено време или место
- Број елемената у нисци s одређује дужину ниске $|s|$
- k -ниска је ниска која садржи k догађаја (ставки)
- Пример ниске: \langle Анализа1, Анализа 2, Анализа 3 \rangle

Формална дефиниција подниске

Дефиниција: Ниска $\langle a_1 a_2 \dots a_n \rangle$ је садржана у нисци $\langle b_1 b_2 \dots b_m \rangle$ ($m \geq n$) ако постоје цели бројеви $i_1 < i_2 < \dots < i_n$ такви да важи $a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots, a_n \subseteq b_{i_n}$

Ниска података	Подниска	Садржи
$\langle \{2,4\} \{3,5,6\} \{8\} \rangle$	$\langle \{2\} \{3,5\} \rangle$	Да
$\langle \{1,2\} \{3,4\} \rangle$	$\langle \{1\} \{2\} \rangle$	Не
$\langle \{2,4\} \{2,4\} \{2,5\} \rangle$	$\langle \{2\} \{4\} \rangle$	Да
$\langle \{2,4\} \{2,5\} \{4,5\} \rangle$	$\langle \{2\} \{4\} \{5\} \rangle$	Не
$\langle \{2,4\} \{2,5\} \{4,5\} \rangle$	$\langle \{2\} \{5\} \{5\} \rangle$	Да
$\langle \{2,4\} \{2,5\} \{4,5\} \rangle$	$\langle \{2,4,5\} \rangle$	Не

Дефиниција истраживања секвенцијалних образца

- Задато
 - база са нискама
 - кориснички дефинисана најмања подршка minsup
- Подршка ниске w је количник броја ниски података које садрже w у односу на укупан број ниски
- Секвенцијални образац је честа подниска (т.ј. подниска чија је подршка $\geq \text{minsup}$)
- Циљ: Наћи све подниске које имају подршку $\geq \text{minsup}$
- Из дате ниске дужине n може да се изведе $\binom{n}{k}$ k -подниски

Истраживање секвенцијалних образца - пример

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

$Minsup = 50\%$

Примери честих подниски:

- | | |
|-----------------------------------|----------|
| $\langle \{1,2\} \rangle$ | $s=60\%$ |
| $\langle \{2,3\} \rangle$ | $s=60\%$ |
| $\langle \{2,4\} \rangle$ | $s=80\%$ |
| $\langle \{3\} \{5\} \rangle$ | $s=80\%$ |
| $\langle \{1\} \{2\} \rangle$ | $s=80\%$ |
| $\langle \{2\} \{2\} \rangle$ | $s=60\%$ |
| $\langle \{1\} \{2,3\} \rangle$ | $s=60\%$ |
| $\langle \{2\} \{2,3\} \rangle$ | $s=60\%$ |
| $\langle \{1,2\} \{2,3\} \rangle$ | $s=60\%$ |

Секвенцијални обрасци / потрошачка корпа

Секвенце

Купац	Датум	Ставке
A	10	2, 3, 5
A	20	1, 6
A	23	1
B	11	4, 5, 6
B	17	2
B	21	1, 2, 7, 8
B	28	1, 6
C	14	1, 7, 8

Подаци из потрошачке корпе

Догађај
2, 3, 5
1, 6
1
4, 5, 6
2
1, 2, 7, 8
1, 6
1, 7, 8

$$\{2\} \rightarrow \{1\} \quad conf(\{2\} \rightarrow \{1\}) = \frac{\sigma(\{2\}\{1\})}{\sigma(\{2\})}$$

$$(1,8) \rightarrow (7) \quad conf(1,8) \rightarrow (7)) = \frac{\sigma(1,7,8)}{\sigma(\{1,8\})}$$

Издвајање секвенцијалних образца

- Дато је н догађаја: i_1, i_2, \dots, i_n
- Кандидатске 1-подниске:
 $\langle \{i_1\} \rangle, \langle \{i_2\} \rangle, \langle \{i_3\} \rangle, \dots, \langle \{i_n\} \rangle$
- Кандидатске 2-подниске:
 $\langle \{i_1, i_2\} \rangle, \langle \{i_1, i_3\} \rangle, \dots, \langle \{i_1\}\{i_1\} \rangle,$
 $\langle \{i_1\}\{i_2\} \rangle, \dots, \langle \{i_{n-1}\}\{i_n\} \rangle$
- Кандидатске 3-подниске:
 $\langle \{i_1, i_2, i_3\} \rangle, \langle \{i_1, i_2, i_4\} \rangle, \dots, \langle \{i_1, i_2\}\{i_1\} \rangle, \langle \{i_1, i_2\}\{i_2\} \rangle$
 $, \dots, \langle \{i_1\}\{i_1, i_2\} \rangle, \langle \{i_1\}\{i_1, i_3\} \rangle,$
 $\dots, \langle \{i_1\}\{i_1\}\{i_1\} \rangle, \langle \{i_1\}\{i_1\}\{i_2\} \rangle, \dots$

Издвајање секвенцијалних образца

- Нека су дата два догађаја: a и b
- Кандидатске 1-подниске:
 $\langle \{a\} \rangle, \langle \{b\} \rangle$
- Кандидатске 2-подниске:
 $\langle \{a\}\{a\} \rangle, \langle \{a\}\{b\} \rangle, \langle \{b\}\{a\} \rangle, \langle \{b\}\{b\} \rangle, \langle \{a,b\} \rangle$
- Кандидатске 3-подниске:
 $\langle \{a\}\{a\}\{a\} \rangle, \langle \{a\}\{a\}\{b\} \rangle, \langle \{a\}\{b\}\{a\} \rangle, \langle \{a\}\{b\}\{b\} \rangle,$
 $\langle \{b\}\{b\}\{b\} \rangle, \langle \{b\}\{b\}\{a\} \rangle, \langle \{b\}\{a\}\{b\} \rangle, \langle \{b\}\{a\}\{a\} \rangle$
 $\langle \{a,b\}\{a\} \rangle, \langle \{a,b\}\{b\} \rangle, \langle \{a\}\{a,b\} \rangle, \langle \{b\}\{a,b\} \rangle$

Формирање секвенцијалних образца

GSP (Generalized Sequential Patterns: Srikant & Agrawal 1996.) алгоритам

- Корак 1: Направити први пролаз кроз базу ниски D ради добијања свих 1-елемент честих подниски
- Корак 2: Понављати поступак све док има нових честих подниски
 - Формирање кандидата: спојити парове честих подниски нађених у $(k-1)$ -овом пролазу ради формирања кандидатских ниски које садрже к догађаја
 - Поткресивање списка кандидата: поткресати скуп кандидатских k -ниски које садрже ретке $(k-1)$ подниске
 - Израчунавање подршке: направити нови пролаз кроз базу ниски D ради налажења подршке за преостале кандидатске ниске
 - Уклањање кандидата: уклонити кандидатске k -ниске чија је подршка мања од minsup

Формирање кандидата

Основни случај ($\kappa=2$)

- Спајањем две честе 1-ниске $\langle \{i_1\} \rangle$ и $\langle \{i_2\} \rangle$ се формирају две кандидатске 2-ниске: $\langle \{i_1\}\{i_2\} \rangle$ и $\langle \{i_1 i_2\} \rangle$

Општи случај ($\kappa>2$)

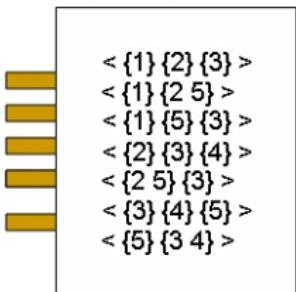
- Честа ($\kappa-1$)-ниска w_1 се спаја са другом честом ($\kappa-1$)-ниском w_2 и формира кандидатска κ -ниска ако је подниска добијена уклањањем првог догађаја из w_1 иста као и подниска добијена уклањањем последњег догађаја из w_2
 - Резултујућа кандидатска ниска је добијена проширењем ниске w_1 последњим догађајем из ниске w_2 . Ако последња два догађаја из w_2 припадају истом елементу, тада последњи догађај из w_2 постаје део последњег елемента у w_1
 - У супротном, последњи догађај из w_2 постаје посебан елемент додат на крај w_1

Примери формирање кандидата

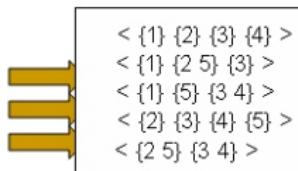
- Спајањем ниски $w_1 = <\{123\}\{46\}>$ и $w_2 = <\{23\}\{46\}\{5\}>$ се формира кандидатска ниска $<\{123\}\{46\}\{5\}>$ јер последњи елемент из w_2 (5) има само један догађај
- Спајањем ниски $w_1 = <\{1\}\{23\}\{4\}>$ и $w_2 = <\{23\}\{45\}>$ се формира кандидатска ниска $<\{1\}\{23\}\{45\}>$ јер последња два догађаја из w_2 (4 и 5) припадају истом елементу
- Спајањем ниски $w_1 = <\{1\}\{23\}\{4\}>$ и $w_2 = <\{23\}\{4\}\{5\}>$ се формира кандидатска ниска $<\{1\}\{23\}\{4\}\{5\}>$ јер последња два догађаја из w_2 (4 и 5) не припадају истом елементу
- Не могу да се споје ниске $w_1 = <\{1\}\{26\}\{4\}>$ и $w_2 = <\{1\}\{2\}\{45\}>$ да би се добила кандидатска ниска $<\{1\}\{26\}\{45\}>$ јер би, у случају да је цео поступак коректан, ниска добила спајањем ниске w_1 са ниском $<\{26\}\{45\}>$

Примери формирање кандидата

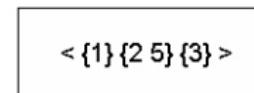
Честе 3-ниске



Формирање кандидата



Поткресивање Кандидата



Налажење секвенцијалних образца - алгоритам

Алгоритам за налажење секвенцијалних образца - верзија слична Apriori

k=1

$F_k = \{i \mid i \in I \wedge \frac{\sigma(\{i\})}{N} \geq minsup\}$ { Наци све 1-подниске }

repeat

k=k+1

$c_k = \text{apriori_generisan}(F_{k-1})$ { Formirati kandidatske k-podnische }

 for svaka_sekvencia_podataka $t \in T$ do

$C_t = \text{podniska}(C_k, t)$ { Наци све кандидате из t }

 for svaka_kandidatska_k-podniska $c \in C_t$ do

$\sigma(c) = \sigma(c) + 1$ { Повећање подрске }

 end for

 end for

$F_k = \{c \mid c \in C_k \wedge \frac{\sigma(\{c\})}{N} \geq minsup\}$ { izdvajanje cestih k-podsniski }

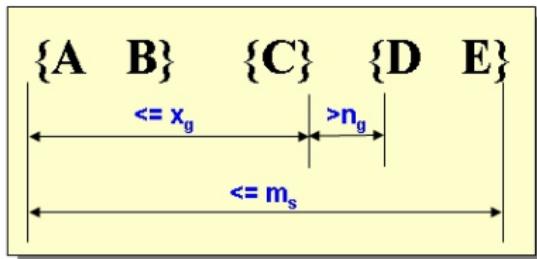
until $F_k = \emptyset$

Rezultat = $\bigcup F_k$

Временска ограничења

- Као један од услова да би ниска била честа може се поставити временско ограничење у коме се ниска појављује
- Ограниченије може да укључи најмању и највећу вредност временског интервала између два појављивања ниске
- Интервал може да се односи на разлику између појављивања прве и последње ставке у комплетној секвенци или на најмању/највећу разлику између појављивања две ниске, или на временски прозор који представља разлику између појављивања прве/последње ставке у појединачној нисци

Временска ограничења



x_g : максимални јаз (max-gap)

n_g : минимални јаз (min-gap)

m_s : максимални размак

$$x_g = 2, n_g = 0, m_s = 4$$

Ниска података	Подниска	Садржи?
$\langle \{2,4\} \{3,5,6\} \{4,7\} \{4,5\} \{8\} \rangle$	$\langle \{6\} \{5\} \rangle$	Да
$\langle \{1\} \{2\} \{3\} \{4\} \{5\} \rangle$	$\langle \{1\} \{4\} \rangle$	Не
$\langle \{1\} \{2,3\} \{3,4\} \{4,5\} \rangle$	$\langle \{2\} \{3\} \{5\} \rangle$	Да
$\langle \{1,2\} \{3\} \{2,3\} \{3,4\} \{2,4\} \{4,5\} \rangle$	$\langle \{1,2\} \{5\} \rangle$	Не

Секвенцијални обрасци са временским ограничењима

- Приступ 1

- Истраживати секвенцијалне обрасце без временских ограничења
- Додатно обрадити откривене обрасце

- Приступ 2

- Модификовати претходне алгоритме да директно поткресују кандидате који крше временска ограничења
- Да ли још увек важи Априори принцип?

Априори принцип за низ података

Object	Timestamp	Events
A	1	1,2,4
A	2	2,3
A	3	5
B	1	1,2
B	2	2,3,4
C	1	1, 2
C	2	2,3,4
C	3	2,4,5
D	1	2
D	2	3, 4
D	3	4, 5
E	1	1, 3
E	2	2, 4, 5

Претпоставка:

$$x_g = 1 \text{ (max-gap)}$$

$$n_g = 0 \text{ (min-gap)}$$

$$m_s = 5 \text{ (maximum span)}$$

$$minsup = 60\%$$

$\langle\{2\} \{5\}\rangle$ подршка = 40%

али

$\langle\{2\} \{3\} \{5\}\rangle$ подршка = 60%

Проблем постоји због ограничења максималног јаза (*max-gap*)

Проблем се не јавља ако је максимални јаз бесконачан

Непрекидне подниске

Важење Априори принципа се обезбеђује увођењем концепта **непрекидних подниски**

Ниска s је **непрекидна подниска** од $w = \langle e_1 \ e_2 \ \dots \ e_k \rangle$ ако важи

- s је добијено из w брисањем догађаја или из e_1 или из e_k , или
- s је добијено из w брисањем догађаја из неког елемента $e_i \in w$ који садржи најмање два догађаја, или
- s је непрекидна подниска од t и t је непрекидна подниска од w (рекурзивна дефиниција)

Ниска података s	Образац t	t непрекидна подниска s
$\langle \{1\} \ \{2,3\} \rangle$	$\langle \{1\} \ \{2\} \rangle$	Да
$\langle \{1,2\} \ \{2\} \ \{3\} \rangle$	$\langle \{1\} \ \{2\} \rangle$	Да
$\langle \{3,4\} \ \{1,2\} \ \{2,3\} \ \{4\} \rangle$	$\langle \{1\} \ \{2\} \rangle$	Да
$\langle \{1\} \ \{3\} \ \{2\} \rangle$	$\langle \{1\} \ \{2\} \rangle$	Не
$\langle \{1,2\} \ \{1\} \ \{3\} \ \{2\} \rangle$	$\langle \{1\} \ \{2\} \rangle$	Не

Модификовани Априори принцип

Дефиниција (модификовани Априори принцип): Ако је k -ниска честа тада су и све њене непрекидне $k - 1$ ниске честе

Применом модификованог Априори принципа на истраживање секвенцијалних образца разматрају се и поткресују кандидатске ниске

- Без ограничења на величину максималног јаза
 - Разматрају се све $(k-1)$ подниске и кандидатска k -ниска се поткресује ако је најмање једна од њених $(k-1)$ -подниски ретка
- Са ограничењем на величину максималног јаза
 - Разматрају се само непрекидне подниске и кандидатска k -ниска се поткресује ако је најмање једна непрекидна $(k-1)$ -подниска ретка

Ограничења величине прозора

Додатно ограничење - величина прозора којим се дефинише највећи дозвољени временски размак између првог и последњег појављивања догађаја у елементима секвенцијалног обрасца. Прозор величине 0 означава да се сви догађаји у истом елементу дешавају истовремено

`ws=2, mingap=0, maxgap=3, maxspan= ∞`

Ниска података	Образац	Садржи?
<code><{1,3} {3,4} {4} {5} {6,7} {8} ></code>	<code><{3,4} {5}></code>	Да
<code><{1,3} {3,4} {4} {5} {6,7} {8} ></code>	<code><{4,6} {8}></code>	Да
<code><{1,3} {3,4} {4} {5} {6,7} {8} ></code>	<code><{3,4,6} {8}></code>	Не
<code><{1,3} {3,4} {4} {5} {6,7} {8} ></code>	<code><{1,3,4} {6,7,8}></code>	Не

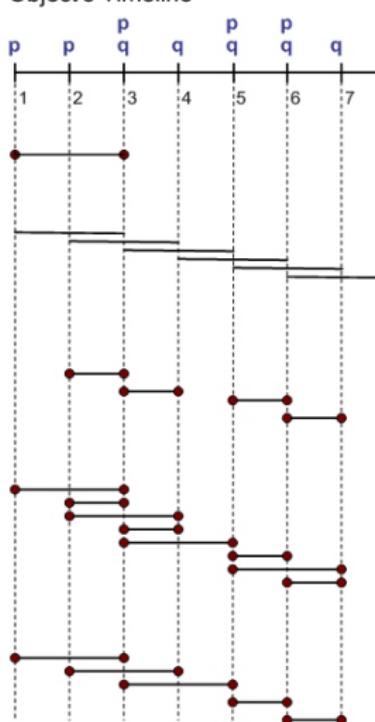
Алтернативне шеме преbroјавања ставки

Пребројавање колико пута се подниска садржи у ниски

- COBJ - једно појављивање по објекту
- CWIN - једно појављивање по помичном прозору величине maxspan
- CMINWIN - број најмањих прозора у којима се ставка појављује
- CDIST_O - различита појављивања са могућношћу временског преклапања
- CDIST - различита појављивања без временског преклапања

Алтернативне шеме преbroјавања ставки

Object's Timeline



Sequence: (p) (q)

Method	Support Count
--------	---------------

COBJ 1

CWIN 6

CMINWIN 4

CDIST_O 8

CDIST 5

Assume:

 $x_g = 2$ (max-gap) $n_g = 0$ (min-gap) $ws = 0$ (window size) $m_s = 2$ (maximum span)

Ретки обрасци

Дефиниција (ретки обрасци): Редак образац је скуп ставки или правило чија је подршка мања од задатог прага $minsup$

Примери:

- Истовремена продаја DVD и VCR снимача је релативно ретка. Ове ставке су негативно корелисане и представљају конкурентске ставке
- Истовремено уписивање изборних предмета Анализа 4, Алгебра 2 и Дискретне структуре 3 на 4. години студија (И смера) је ретко. Ови предмети су међусобно конкурентни и њихова заједничка појава је негативно корелисана
- Ако је {Пожар = Да} честа ставка али {Пожар = Да, Аларм = Не} је ретка, ово правило је важна ретка ставка јер означава грешку у алармном систему

Негативни обрасци

Нека је $I = \{i_1, i_2, \dots, i_d\}$ скуп ставки. Негативна ставка $\overline{i_k}$ означава одсуство ставке i_k из дате трансакције.

Дефиниција (негативан скуп ставки): Негативан скуп ставки X је скуп ставки који има следеће особине

- $X = A \cup \overline{B}$ где је A скуп позитивних ставки, \overline{B} је скуп негативних ставки, $|\overline{B}| \geq 1$
- $sup(X) \geq minsup$

Дефиниција (негативно правило придрживања): Негативно правило придрживања је правило придрживања које има следеће особине

- правило је издвојено из негативног скупа ставки
- подршка правила је $\geq minsup$
- поузданост правила је $\geq minconf$

Негативно корелисани обрасци

Дефиниција (негативно корелисани обрасци): Скуп ставки $X = \{x_1, x_2, \dots, x_k\}$ је негативно корелисан ако

$$sup(X) < \prod_{j=1}^k sup(x_j) = sup(x_1) \times sup(x_2) \times \dots \times sup(x_k)$$

где $sup(x_j)$ означава подршку ставке x_j

Дефиниција (негативно корелисано правило придрживања, делимични услов): Правило придрживања $X \longrightarrow Y$ је негативно корелисано ако

$$sup(X \cup Y) < sup(X) \times sup(Y)$$

где су X и Y дисјунктни скупови ставки

Негативно корелисани обрасци

Дефиниција (негативно корелисано правило придрживања, пуни услов): Правило придрживања $X \rightarrow Y$ је негативно корелисано ако

$$sup(X \cup Y) < \prod_j sup(x_j) \times \prod_j sup(y_j)$$

где су X и Y дисјунктни скупови ставки

Како су ставке у X и Y међусобно позитивно корелисане, практичније је користити делимичан уместо пуног услова негативне корелисаности правила придрживања. На пример, иако је правило

{наочаре, марамице за брисање стакала} \rightarrow {контактна сочива, течност за сочива}

негативно корелисано, ставке на левој и десној страни су међусобно корелисане и у случају примене пуног услова корелисаности ово правило се не би јавило

Негативно корелисани обрасци

Како је

$$\begin{aligned} & \text{sup}(X \cup Y) - \text{sup}(X) \times \text{sup}(Y) \\ &= \text{sup}(X \cup Y) - [\text{sup}(X \cup Y) + \text{sup}(X \cup \bar{Y})] \times [\text{sup}(X \cup Y) + \text{sup}(\bar{X} \cup Y)] \\ &= \text{sup}(X \cup Y) \times [1 - \text{sup}(X \cup Y) - \text{sup}(X \cup \bar{Y}) - \text{sup}(\bar{X} \cup Y)] - \text{sup}(X \cup \bar{Y}) \times \text{sup}(\bar{X} \cup Y) \\ &= \text{sup}(X \cup Y) \times \text{sup}(\bar{X} \cup \bar{Y}) - \text{sup}(X \cup \bar{Y}) \times \text{sup}(\bar{X} \cup Y) \end{aligned}$$

Услов негативне корелације може да се запише као

$$\text{sup}(X \cup Y) \times \text{sup}(\bar{X} \cup \bar{Y}) < \text{sup}(X \cup \bar{Y}) \times \text{sup}(\bar{X} \cup Y)$$

Поређење образца

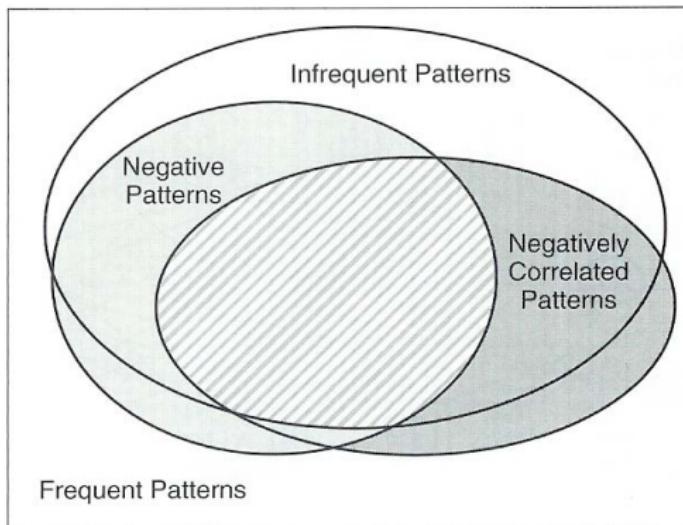
- Многи ретки обрасци имају одговарајуће негативне обрасце
- Многи негативно корелисани обрасци такође имају одговарајуће негативне обрасце
- Што је мања подршка за $X \cup Y$ образац је више негативно корелисан
- Ретки негативно корелисани обрасци су интересантнији од честих негативно корелисаних образаца

Табела контингената за правило придрживања $X \rightarrow Y$

	Y	\bar{Y}	
X	$sup(X \cup Y)$	$sup(X \cup \bar{Y})$	$sup(X)$
\bar{X}	$sup(\bar{X} \cup Y)$	$sup(\bar{X} \cup \bar{Y})$	$sup(\bar{X})$
	$sup(Y)$	$sup(\bar{Y})$	1

Ако су X и Y негативно корелисани тада $X \cup \bar{Y}$, $\bar{X} \cup Y$ или оба морају да имају релативно високу подршку

Порећење



Технике за истраживање негативних образца

Техника заснована на симетричним бинарним променљивим

TID	Items
1	{A,B}
2	{A,B,C}
3	{C}
4	{B,C}
5	{B,D}



Original Transactions

TID	A	\bar{A}	B	\bar{B}	C	\bar{C}	D	\bar{D}
1	1	0	1	0	0	1	0	1
2	1	0	1	0	1	0	0	1
3	0	1	0	1	1	0	0	1
4	0	1	1	0	1	0	0	1
5	0	1	1	0	0	1	1	0

Transactions with Negative Items

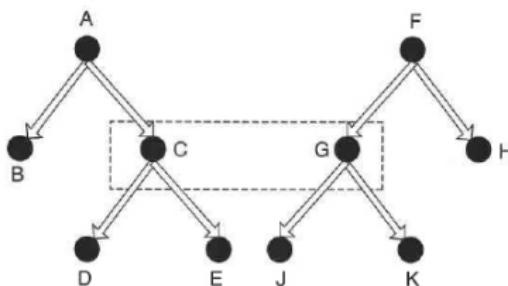
Проблеми

- Велики број ставки (уместо решетке са 2^d ставки добија се решетка са 2^{2d})
- Поткресивање на основу подршке не може да се примени (за сваку ставку или x или \bar{x} имају подршку $\geq 50\%$)
- Дужина трансакције се повећава када се укључе негативне ставке (на бар укупан број ставки)

Технике засноване на очекиваној подршци

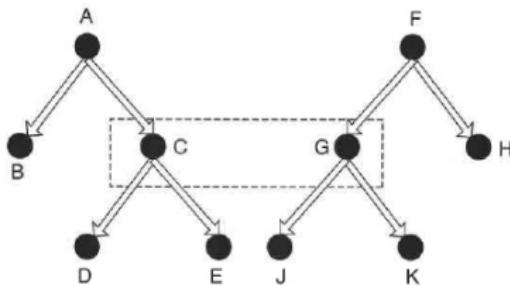
- Технике засноване на хијерархији концепата
- Технике засноване на индиректном придрживању

Технике засноване на хијерархији концепата



Ако је ставка $\{C, G\}$ честа, а ставка $\{D, J\}$ није, тада D и J формирају интересантан образац

Технике засноване на хијерархији концепата

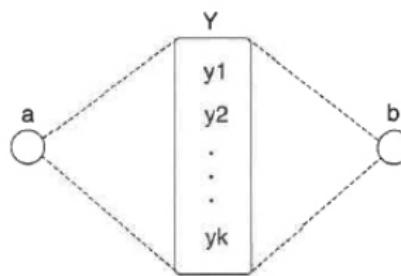


$$\varepsilon(s(E, J)) = s(C, G) \times \frac{s(E)}{s(C)} \times \frac{s(J)}{s(G)}$$

$$\varepsilon(s(C, J)) = s(C, G) \times \frac{s(J)}{s(G)}$$

$$\varepsilon(s(C, H)) = s(C, G) \times \frac{s(H)}{s(G)}$$

Технике засноване на индиректном
придруживању



Пар ставки a и b је индиректно придружен преко медијатора ('комшија' који је доволно близу) Y ако важи

- подршка за пар $\sup(\{a, b\}) < t_s$
 - $\exists Y \neq \emptyset$ тако да
 - $s(\{a\} \cup Y) \geq t_f$ и $s(\{b\} \cup Y) \geq t_f$
 - $d(\{a\} \cup Y) \geq t_d$ $d(\{b\} \cup Y) \geq t_d$ где је $d(X, Z)$ мера придрживања X и Z

Класификација помоћу правила придрживања

Добијање правила за класификацију

- Скуп правила са великим поузданошћу
- На десном делу правила само једна ставка
- Десни део правила представља класу која се предвиђа
- Ставке на левом делу правила чине атрибуте на основу којих се врши предвиђање