

Додатне технике правила придруживања

Ненад Митић

Математички факултет
`nenad@matf.bg.ac.rs`

Статистички засноване методе

- Primer: `[status=0]+[datpolaganja >= 736950 and < 737450]` →
[ocena]: $\mu = 7.5$
- 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 правилном у односу на део популације који није покривен правилном:
 $A \longrightarrow B : \mu$ према $\bar{A} \longrightarrow B : \mu'$
- 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

Методe заснованe 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

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

Пример: $\text{sup}(W_1, W_2, W_3) = 0 + 0 + 0 + 0 + 0.17 = 0.17$

	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

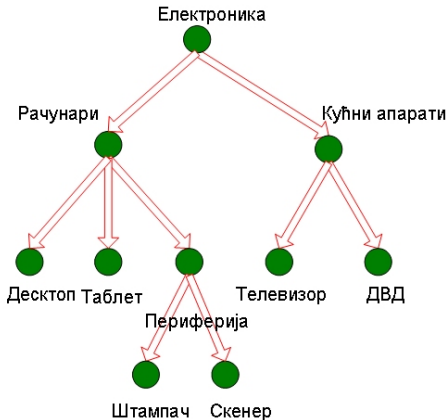
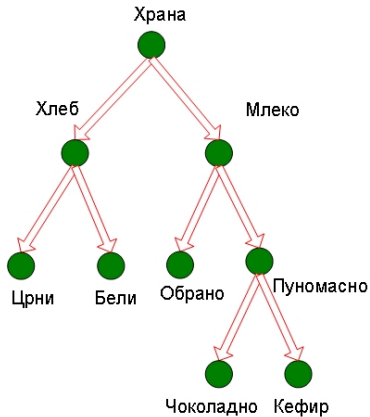
- монтоно расте како расте нормализована фреквенција речи
- монтоно расте како расте број документата који садрже реч
- монотонно опада како расте број речи у скупу ставки -
АНТИ-МОНОТОНОСТ

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 придруживања између више нивоа



Правила придруживања 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 \rightarrow бели хлеб, пуномасно млеко \rightarrow црни хлеб, кефир \rightarrow цри хлеб, ... су индикатори правила придруживања 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 $\text{conf}(X_1 \rightarrow Y) \geq \text{minconf}$

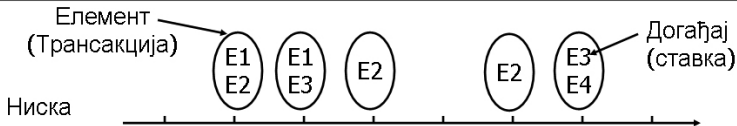
Правила придруживања између више нивоа

Приступ 1:

- Проширити текуће правило придруживања 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 образаца укључује ставке са вишег нивоа
 - Повећава се димензионалност података

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 posetioca Veba	Skup datoteka pregledan od strane posetioca 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: $\langle \text{Анализа 1, Анализа 2, Анализа 3} \rangle$

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 образаца

- 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 образаца - 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%

Primeri čestih podniski:

< {1,2} > s=60%
 < {2,3} > s=60%
 < {2,4} > s=80%
 < {3} {5} > s=80%
 < {1} {2} > s=80%
 < {2} {2} > s=60%
 < {1} {2,3} > s=60%
 < {2} {2,3} > s=60%
 < {1,2} {2,3} > s=60%

Издавање секвенцијалних образаца

- Дато је n 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:
 $\langle \{i_1, i_2, i_3\} \rangle, \langle \{i_1, i_2, i_4\} \rangle, \dots, \langle \{i_1, i_2\}\{i_1\} \rangle, \langle \{i_1, i_2\}\{i_2\} \rangle$
 $, \dots, \langle \{i_1\}\{i_1, i_2\} \rangle, \langle \{i_1\}\{i_1, i_3\} \rangle,$
 $\dots, \langle \{i_1\}\{i_1\}\{i_1\} \rangle, \langle \{i_1\}\{i_1\}\{i_2\} \rangle, \dots$

Издавање секвенцијалних образаца

- 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 образаца

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 k 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 ($k=2$)

- 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 ($k>2$)

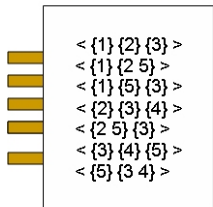
- Честа $(k-1)$ -ниска w_1 се спаја са другом честом $(k-1)$ -ниском w_2 и формира кандидатска k -ниска ако је подниска добијена уклањањем првог догађаја из w_1 иста као и подниска добијена уклањањем последњег догађаја из w_2
 - Резултујућа кандидатска ниска је добијена проширењем ниске w_1 последњим 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

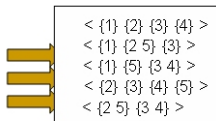
- Спајањем ниски $w_1 = \langle \{123\}\{46\} \rangle$ и $w_2 = \langle \{23\}\{46\}\{5\} \rangle$ се формира кандидатска ниска $\langle \{123\}\{46\}\{5\} \rangle$ јер последњи елемент из w_2 (5) има само један догађај
- Спајањем ниски $w_1 = \langle \{1\}\{23\}\{4\} \rangle$ и $w_2 = \langle \{23\}\{45\}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

Примери формирање кандидата

Честе 3-ниске



Формирање кандидата



Поткресивање
Кандидата

< {1} {2 5} {3} >

Налажење секвенцијалних образаца - алгоритам

Алгоритам за налажење секвенцијалних образаца - верзија слична Apriori

$k=1$

$F_k = \{i \mid i \in I \wedge \frac{\sigma(\{i\})}{N} \geq \text{minsup}\}$ {Naci sve ceste 1-podniske}

repeat

$k=k+1$

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

for svaka_sekvenca_podataka $t \in T$ **do**

$C_t = \text{podniska}(C_k, t)$ {Naci sve kandidate iz t}

for svaka_kandidatska k-podniska $c \in C_i$ **do**

$\sigma(c) = \sigma(c) + 1$ {Povecanje podrske}

end for

end for

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

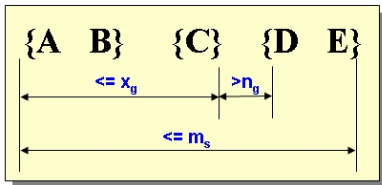
until $F_k = \emptyset$

Rezultat = $\cup F_k$

Временска ограничења

- Као један од услова да би ниска била честа може се поставити временско ограничење у коме се ниска појављује
- Ограничење 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



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

x_g : максимални јаз (max-gap)

n_g : минимални јаз (min-gap)

m_s : максимални размак

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:

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

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

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

$$\text{minsup} = 60\%$$

$\langle \{2\} \{5\} \rangle$ 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 превазилази увођењем концепта *непрекидних подниски*

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

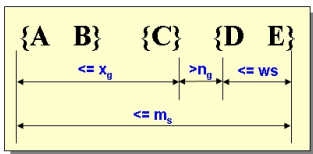
Дефиниција (модификовани Априори принцип): Ако је k -ниска честа тада су и све њене непрекидне $k - 1$ ниске честе

Применом модификованог Априори принципа на истраживање секвенцијалних образаца 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



x_g : максимални јаз (max-gap)

n_g : минимални јаз (min-gap)

ws: величина прозора

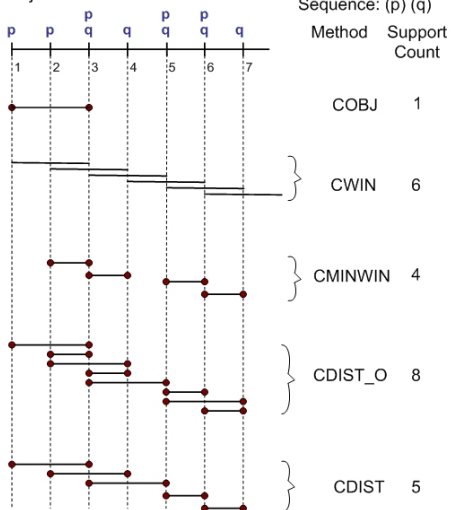
m_s : максимални размак

$$x_g = 2, n_g = 0, ws = 1, m_s = 5$$

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

Алтернативне шеме пребројавања ставки

Object's Timeline



Assume:

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