

SEKVENČNÉ PRAVIDLÁ

Sekvenčné pravidlá sú odvodené od asociačných, rozdiely preto nie sú príliš veľké. K - sekvencia je sekvencia dĺžky k , t. j. obsahuje k stránok. Frekventovaná k - sekvencia je obdobou frekventovanej k - položkovej množiny, resp. kombinácie. Zvyčajne, frekventovaná 1 - položková množina je rovnaká ako frekventovaná 1 - sekvencia.

Na nasledujúcom príklade naznačíme rozdiely medzi algoritmom Apriori a AprioriAll, kde algoritmus Apriori slúži na hľadanie asociačných pravidiel a AprioriAll na hľadanie sekvenčných pravidiel:

$$D = \{T1 = \{U1, \langle A, B, C \rangle\}, T2 = \{U2, \langle A, C \rangle\}, T3 = \{U1, \langle B, C, E \rangle\}, T4 = \{U3, \langle A, C, D, C, E \rangle\}\},$$

kde D je databáza transakcií s časovou nálepkou, A, B, C, D, E sú webové časti a $U1, U2, U3$ sú používatelia. Každá transakcia je identifikovaná používateľom. Predpokladajme, že minimálna podpora (*min support*) je 30%.

V tomto príklade má používateľ $U1$ aktuálne dve transakcie. Pri hľadaní sekvenčných pravidiel uvažujeme jeho sekvenciu ako aktuálne spojenie webových častí v transakciách $T1$ a $T3$, t. j. sekvencia môže pozostávať z viacerých transakcií, pričom sa nevyžadujú súvislé prístupy na stránky. Rovnako podpora sekvencie je určená nie percentom transakcií, ale percentom používateľov, ktorí majú danú sekvenciu. Sekvencia je veľká/frekventovaná ak sa prinajmenšom nachádza v jednej sekvencii identifikovanej používateľom a spĺňa podmienku minimálnej podpory. Prvým krokom je zoradenie transakcií podľa používateľa s časovou nálepkou každej ním navštívenej stránky, zostávajúce kroky sú podobné ako v algoritme Apriori. Po zoradení sme získali aktuálne sekvencie identifikované používateľom, ktoré predstavujú kompletne odkazy od jedného používateľa:

$$D = \{T1 = \{U1, \langle A, B, C \rangle\}, T3 = \{U1, \langle B, C, E \rangle\}, T2 = \{U2, \langle A, C \rangle\}, T4 = \{U3, \langle A, C, D, C, E \rangle\}\}.$$

Podobne ako v algoritme Apriori začíname generovaním množiny kandidátov dĺžky 1: $C1 = \{\langle A \rangle, \langle B \rangle, \langle C \rangle, \langle D \rangle, \langle E \rangle\}$, z nej určíme množinu $L1 = \{\langle A \rangle, \langle B \rangle, \langle C \rangle, \langle D \rangle, \langle E \rangle\}$ jednoprvkových sekvencií takých, kde každá stránka je odkazovaná prinajmenšom jedným používateľom.

Následne množiny kandidátov $C2$ vygenerujeme z $L1$ tzv. úplným spájaním, t. j. zohľadňujeme, že používateľ webu prehľadáva stránky dopredu alebo dozadu. Z tohto dôvodu algoritmus Apriori nie je vhodný na objavovanie znalostí z webových prístupov (web log mining), na rozdiel od algoritmu AprioriAll, ktorý spomínanú skutočnosť zohľadňuje. Z množiny kandidátov dĺžky 2: $C2 = \{\langle A, B \rangle, \langle A, C \rangle, \langle A, D \rangle, \langle A, E \rangle, \langle B, A \rangle, \langle B, C \rangle,$

$\langle B, D \rangle, \langle B, E \rangle, \langle C, A \rangle, \langle C, B \rangle, \langle C, D \rangle, \langle C, E \rangle, \langle D, A \rangle, \langle D, B \rangle, \langle D, C \rangle, \langle D, E \rangle,$
 $\langle E, A \rangle, \langle E, B \rangle, \langle E, C \rangle, \langle E, D \rangle\}$, určíme množinu $L2 = \{\langle A, B \rangle, \langle A, C \rangle, \langle A, D \rangle, \langle A,$
 $E \rangle, \langle B, C \rangle, \langle B, E \rangle, \langle C, B \rangle, \langle C, D \rangle, \langle C, E \rangle, \langle D, C \rangle, \langle D, E \rangle\}$ dvojprvkových
sekvencií takých, kde sa každá sekvencia prinajmenšom nachádza v jednej sekvencii
identifikovanej používateľom. Analogicky postupujeme ďalej.