

## Κανόνες Συσχέτισης II

Οι διαφάνειες στηρίζονται στο P.-N. Tan, M. Steinbach, V. Kumar, «Introduction to Data Mining», Addison Wesley, 2006



## Σύντομη Ανακεφαλαίωση



### Εισαγωγή

#### Market-Basket transactions (Το καλάθι της νοικοκυράς!)

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

δοσοληψία

- Προώθηση προϊόντων
- Τοποθέτηση προϊόντων στα ράφια
- Διαχείριση αποθεμάτων

Το πρόβλημα: Δεδομένου ενός συνόλου δοσοληψιών (transactions), βρες κανόνες που προβλέπουν την εμφάνιση στοιχείων (item) με βάση την εμφάνιση άλλων στοιχείων στις συναλλαγές

#### Παραδείγματα κανόνων συσχέτισης

{Diaper} → {Beer},  
{Milk, Bread} → {Eggs, Coke},  
{Beer, Bread} → {Milk}

Σημαίνει ότι εμφανίζονται μαζί, όχι ότι η εμφάνιση του ενός είναι η αιτία της εμφάνισης του άλλου (co-occurrence, not causality όχι έννοια χρόνου ή διάταξης)



### Ορισμοί

**στοιχειοσύνολο (itemset):** Ένα υποσύνολο του συνόλου των στοιχείων

**k-στοιχειοσύνολο (k-itemset):** ένα στοιχειοσύνολο με k στοιχεία

**support count (σ) ενός στοιχειοσυνόλου:** ο αριθμός εμφανίσεων του στοιχείου

**Υποστήριξη (Support (s)) ενός στοιχειοσυνόλου** Το ποσοστό των δοσοληψιών που περιέχουν ένα στοιχειοσύνολο

**Συχνό Στοιχειοσύνολο (Frequent Itemset)** Ένα στοιχειοσύνολο του οποίου η υποστήριξη είναι μεγαλύτερη ή ίση από κάποια τιμή κατωφλίου *minsup*

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke



#### Κανόνας Συσχέτισης (Association Rule)

Είναι μια έκφραση της μορφής  $X \rightarrow Y$ , όπου X και Y είναι στοιχειοσύνολα  $X \subseteq I, Y \subseteq I, X \cap Y = \emptyset$

Παράδειγμα: {Milk, Diaper} → {Beer}

#### Υποστήριξη Κανόνα Support (s)

Το ποσοστό των δοσοληψιών που περιέχουν και το X και το Y ( $X \cup Y$ )

#### Εμπιστοσύνη - Confidence (c)

Πόσες από τις δοσοληψίες (ποσοστό) που περιέχουν το X περιέχουν και το Y

#### Πρόβλημα

Εύρεση Κανόνων Συσχέτισης

Είσοδος: Ένα σύνολο από δοσοληψίες T  
Έξοδος: Όλοι οι κανόνες με  $support \geq minsup$   
 $confidence \geq minconf$

### Ορισμοί

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke



## Εξόρυξη Κανόνων Συσχέτισης

Χωρισμός του προβλήματος σε δύο υπο-προβλήματα:

- **Εύρεση όλων των συχνών στοιχειοσυνόλων (Frequent Itemset Generation)**

Εύρεση όλων των στοιχειοσυνόλων με υποστήριξη  $\geq minsup$

- **Δημιουργία Κανόνων (Rule Generation)**

Για κάθε στοιχειοσύνολο, δημιούργησε κανόνες με μεγάλη υποστήριξη, όπου κάθε κανόνες είναι μια δυαδική διαμέριση του συχνού στοιχειοσυνόλου



### Εύρεση Συχνών Στοιχειοσυνόλων

**Itemset Lattice - Πλέγμα Στοιχειοσυνόλων**

Για  $d$  στοιχεία,  $2^d$  πιθανά στοιχειοσύνολα

Εξόφλη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΕΥΧΕΤΙΣΤΗΣ II 7

### Εύρεση Συχνών Στοιχειοσυνόλων: Στρατηγική αρriori

**Αρχή Αρriori**  
Αν ένα στοιχειοσύνολο είναι συχνό, τότε όλα τα υποσύνολα του είναι συχνά

Εξόφλη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΕΥΧΕΤΙΣΤΗΣ II 8

### Στρατηγική αρriori

#### Γενικός Αλγόριθμος

Έστω  $k = 1$      $\#k$ : μήκος στοιχειοσυνόλου

Παρήγαγε τα συχνά **1-στοιχειοσύνολα**

**Repeat until** να μην παράγονται νέα συχνά στοιχειοσύνολα

1. Παρήγαγε υποψήφια ( $k+1$ )-στοιχειοσύνολα
2. Ψαλίδισε τα υποψήφια στοιχειοσύνολα που περιέχουν μη συχνά στοιχειοσύνολα μεγέθους  $k$
3. Υπολόγισε την υποστήριξη κάθε υποψήφιου ( $k+1$ )-στοιχειοσυνόλου διασχίζοντας τη βάση των δοσοληφιών
4. Σβήσε τα υποψήφια στοιχειοσύνολα που δεν είναι συχνά
5.  $k = k + 1$

Εξόφλη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΕΥΧΕΤΙΣΤΗΣ II 9

### Στρατηγική αρriori: Δημιουργία Στοιχειοσυνόλων

Για την παραγωγή υποψήφιων  $k$ -στοιχειοσυνόλων

$F_{k-1} \times F_1$

Επέκταση κάθε συχνού ( $k-1$ ) στοιχειοσυνόλου με άλλα συχνά στοιχεία

$F_{k-1} \times F_{k-1}$

Συγχώνευση δύο συχνών ( $k-1$ ) στοιχειοσυνόλου αν τα πρώτα  $k-2$  στοιχεία τους είναι τα ίδια

Παρατηρήσεις

- Για να αποφύγουμε τη δημιουργία του ίδιου στοιχειοσυνόλου, κρατάμε κάθε στοιχειοσύνολο (λεξικογραφικά) ταξινομημένο
- Είναι δυνατόν να γίνουν απλοί έλεγχοι αν τα παραγόμενα πιθανά στοιχειοσύνολα είναι συχνά ελέγχοντας αν τα υποσύνολα τους είναι συχνά και έτσι να αποφύγουμε να υπολογίσουμε την υποστήριξή τους

Εξόφλη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΕΥΧΕΤΙΣΤΗΣ II 10

### Στρατηγική αρriori: Υπολογισμός Υποστήριξης

Για κάθε νέο υποψήφιο  $k+1$ -στοιχειοσυνόλο, πρέπει να υπολογίσουμε την υποστήριξή του

- Για να μειώσουμε τον αριθμό των πράξεων, σε κάθε βήμα, αποθηκεύουμε τα υποψήφια  $k+1$ -στοιχειοσυνόλα σε ένα **δέντρο κατακερματισμού**
- Αντί να ταιριάζουμε κάθε δοσοληφία με κάθε υποψήφιο στοιχειοσυνόλο, κατακερματίζουμε τα στοιχειοσύνολα της δοσοληφίας και ενημερώνουμε μόνο τους αντίστοιχους κώδους του δέντρου κατακερματισμού των συχνών στοιχειοσυνόλων

Εξόφλη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΕΥΧΕΤΙΣΤΗΣ II 11

### Παραγωγή Κανόνων

Δοθέντος ενός συχνού στοιχειοσυνόλου  $L$ , βρες όλα τα μη κενά υποσύνολα  $f \subset L$  τέτοια ώστε ο κανόνας  $f \rightarrow L - f$  ικανοποιεί τον περιορισμό της ελάχιστης εμπιστοσύνης

Η εμπιστοσύνη για τους κανόνες που παράγονται από το ίδιο στοιχειοσύνολο έχει μια αντι-μονότονη ιδιότητα

Για παράδειγμα  $L = \{A, B, C, D\}$ :  $c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$

Η εμπιστοσύνη είναι αντι-μονότονη σε σχέση με των αριθμό των στοιχείων στο RHS του κανόνα

Εξόφλη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΕΥΧΕΤΙΣΤΗΣ II 12

### Παραγωγή Κανόνων για τον Αλγόριθμο αργίοι

Πιλέγμα Κανόνων για το Στοιχειοσύνολο {A, B, C, D}

Ψαλίδισμα με βάση την εμπιστοσύνη

Εστω κόμβος με μικρή εμπιστοσύνη

Pruned Rules

Για κάθε συχνό στοιχειοσύνολο, ξεκινάμε με έναν κανόνα που έχει μόνο  $k-1$  στοιχεία στο δεξιό μέρος του

Υπολογίζουμε την εμπιστοσύνη

Παράγουμε κανόνες με  $k+1$  στοιχεία και υπολογίζουμε την εμπιστοσύνη τους

Για τον υπολογισμό της εμπιστοσύνης δεν χρειάζεται να διαπεράσουμε τη βάση

Εξώφυλλο Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 13

### Αναπαράσταση Στοιχειοσυνόλων

Τα στοιχειοσύνολα που παράγονται είναι πολλά, κάποια ίσως περιττά - οδηγούν σε παραγωγή πολλών κανόνων

Τι να κρατήσουμε;

Αντιπροσωπευτικά συχνά στοιχειοσύνολα:

- Maximal συχνά
- Κλειστά συχνά

Εξώφυλλο Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 14

### Αναπαράσταση Στοιχειοσυνόλων

Ένα στοιχειοσύνολο είναι **maximal συχνό** αν κανένα από τα άμεσα υπερσυνόλά του δεν είναι συχνό

Προσφέρουν μια συνοπτική αναπαράσταση των συχνών στοιχειοσυνόλων: το μικρότερο σύνολο στοιχειοσυνόλων από το οποίο μπορούμε να πάρουμε όλα τα συχνά στοιχειοσύνολα - είναι τα υποσυνόλά τους

Πρόβλημα: Δεν προσφέρουν καμία πληροφορία για την υποστήριξη των υποσυνόλων τους

Συχνά

Μη συχνά

Εξώφυλλο Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 15

### Αναπαράσταση Στοιχειοσυνόλων

Ένα στοιχειοσύνολο είναι **κλειστό (closed)** αν κανένα από τα άμεσα υπερσυνόλά του δεν έχει την ίδια υποστήριξη με αυτό (δηλαδή, έχει μικρότερη)

Ένα στοιχειοσύνολο είναι **κλειστό συχνό στοιχειοσύνολο** αν είναι κλειστό και η υποστήριξη του είναι μικρότερη ή ίση με  $\text{minsup}$

Πάλι τα υποσύνολα τους μας δίνουν όλα τα συχνά υποσύνολα, τώρα όμως μπορούμε να υπολογίσουμε την υποστήριξη των υποσυνόλων τους

Πως: Η υποστήριξη ενός μη κλειστού στοιχειοσυνόλου πρέπει να είναι ίση με την μεγαλύτερη υποστήριξη ανάμεσα στα υπερσυνόλά του

ABC=>D

Εξώφυλλο Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 16

### Αναπαράσταση Στοιχειοσυνόλων

Maximal vs Κλειστά στοιχειοσύνολα

Minimum support = 2

Closed but not maximal

Closed and maximal

# Closed = 9  
# Maximal = 4

Εξώφυλλο Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 17

### Μέτρηση Ενδιαφέροντος

Στην αρχική διατύπωση του προβλήματος της εξόρυξης κανόνων συσχέτισης χρησιμοποιήθηκαν ως μέτρα μόνο η **υποστήριξη** και η **εμπιστοσύνη**

Μια σειρά μέτρων βασισμένα στη στατιστική

- $P(S \wedge B) = P(S) \times P(B) \Rightarrow$  Στατιστική ανεξαρτησία
- $P(S \wedge B) > P(S) \times P(B) \Rightarrow$  Positively correlated (θετική συσχέτιση)
- $P(S \wedge B) < P(S) \times P(B) \Rightarrow$  Negatively correlated (αρνητική συσχέτιση)

A=>BCD

Εξώφυλλο Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 18

### Μέτρηση Ενδιαφέροντος

Έστω ένας κανόνας,  $X \rightarrow Y$ , η πληροφορία που χρειάζεται για τον υπολογισμό του ενδιαφέροντος του κανόνα μπορεί να υπολογιστεί από τον **contingency table**

Contingency table for  $X \rightarrow Y$

	Y	$\bar{Y}$	
X	$f_{11}$	$f_{10}$	$f_{1+}$
$\bar{X}$	$f_{01}$	$f_{00}$	$f_{0+}$
	$f_{+1}$	$f_{+0}$	T

$f_{11}$ : support of X and Y  
 $f_{10}$ : support of X and  $\bar{Y}$   
 $f_{01}$ : support of  $\bar{X}$  and Y  
 $f_{00}$ : support of  $\bar{X}$  and  $\bar{Y}$

Μέτρηση συχνότητας εμφάνισης → Χρησιμοποιείται για τον ορισμό διαφόρων μέτρων

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 19

### Μέτρηση Ενδιαφέροντος

Μέτρα που λαμβάνουν υπ' όψιν τους στη στατιστική εξάρτηση  $X \rightarrow Y$

$$Lift = \frac{P(Y | X)}{P(Y)}$$

$$Interest = \frac{P(X, Y)}{P(X)P(Y)}$$

$$PS = P(X, Y) - P(X)P(Y)$$

$$\phi - coefficient = \frac{P(X, Y) - P(X)P(Y)}{\sqrt{P(X)[1 - P(X)]P(Y)[1 - P(Y)]}}$$

Γενικά έχουν προταθεί 21 τέτοια μέτρα

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 20

### Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 21

### Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

Ο Αρρίγιο από τους παλιότερους, αλλά:

- Συχνά μεγάλο I/O επειδή κάνει πολλαπλά περάσματα στη βάση των δοσοληψιών
- Κακή απόδοση όταν οι δοσοληψίες έχουν μεγάλο πλάτος

Άλλες μέθοδοι:

- Διαφορετικές διασχίσεις του πλέγματος των στοιχειοσυνόλων
- Αναπαράσταση Συνόλου Δοσοληψιών

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 22

### Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

Αρρίγιο: Γενικά προς Συγκεκριμένα  $k-1 \rightarrow k$

Itemset Lattice - Πλέγμα Στοιχειοσυνόλων

Αν αυτό είναι το συχνό, το βρίσκουμε ελέγχουμε όλα τα υποσυνόλά του

Αν τα συχνά είναι προς το κατώτατο σημείο (bottom) τους πλέγματος, ίσως συμφέρει

Συγκεκριμένα προς Γενικά  $k \rightarrow k-1$

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 23

### Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

Διάσχιση του Πλέγματος των Στοιχειοσυνόλων: Συγκεκριμένα-προς-Γενικά vs Γενικά-προς-Συγκεκριμένα

$k \rightarrow k-1$  (συγκεκριμένο-προς-γενικό)

Ποιο χρήσιμο για τον εντοπισμό **maximal** συχνών στοιχειοσυνόλων σε πυκνές (δηλ. με μεγάλο πλάτος δοσοληψίες) όπου το συχνό στοιχειοσύνολο βρίσκεται κοντά στο κατώτατο σημείο του πλέγματος

Αν συχνό, δε χρειάζεται να ελέγξουμε κανένα από τα υποσυνόλά του

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 24

### Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

#### Διάσχιση του Πλέγματος των Στοιχειοσυνόλων: Κλάσεις Ισοδυναμίας

Χωρισμός των στοιχειοσυνόλων του πλέγματος σε ξένες μεταξύ τους ομάδες (κλάσεις ισοδυναμίας) και εξέταση των στοιχειοσυνόλων ανά κλάση

Αpriori: ορίζει τις κλάσεις με βάση το μήκος K των στοιχειοσυνόλων, πρώτα αυτά μήκους 1, μετά μήκους 2 κ.ο.κ

Prefix (Suffix): Δύο στοιχειοσύνολα ανήκουν στην ίδια κλάση αν έχουν κοινό πρόθεμα (επίθημα) μήκους K

Εξέλιξη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 25

### Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

#### Διάσχιση του Πλέγματος των Στοιχειοσυνόλων: BFS vs DFS

Αpriori

BFS: Breadth-First-Search

Χρήσιμο για την εύρεση maximal συχνών στοιχειοσυνόλων γιατί τα εντοπίζει πιο γρήγορα από το BFS

Μόλις εντοπιστεί το maximal, είναι δυνατόν να κλαδευτούν πολλά υποσύνολα του

DFS: Depth-First-Search

Εξέλιξη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 26

### Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

#### Διάσχιση του Πλέγματος των Στοιχειοσυνόλων: BFS vs DFS

Figure 6.22. Generating candidate itemsets using the depth-first approach.

Εξέλιξη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 27

### Άλλοι Τρόποι Υπολογισμού

#### Αναπαράσταση της Βάσης Δεδομένων: Οριζόντια vs Κάθετη

Αυτό χρησιμοποιεί ο Apriori

Εναλλακτικά: Για κάθε στοιχείο σε ποιες δοσοληψίες εμφανίζεται

Horizontal Data Layout

TID	Items
1	A,B,E
2	B,C,D
3	C,E
4	A,C,D
5	A,B,C,D
6	A,E
7	A,B
8	A,B,C
9	A,C,D
10	B

Vertical Data Layout

	A	B	C	D	E
1	1	2	2	1	
2		1	3	4	3
3			5	4	5
4			6	7	8
5			7	8	9
6			8	9	
7			9		
8			10		
9					
10					

Η υποστήριξη υπολογίζεται παίρνοντας τις τομές των TID-λίστλων

Εξέλιξη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 28

null

### Άλλοι Τρόποι Υπολογισμού

Η υποστήριξη υπολογίζεται παίρνοντας τις τομές των TID-λίστων

A
1
4
5
6
7
8
9

 $\wedge$ 

B
1
2
5
7
8
10

 $\rightarrow$ 

AB
1
5
7
8

- Η υποστήριξη ενός K-στοιχειοσυνόλου υπολογίζεται παίρνοντας τις τομές των TID-λίστων δύο από τα (K-1)-ύπο-στοιχειοσύνολα του.
- Πλεονέκτημα: πολύ γρήγορος υπολογισμός της υποστήριξης
- Πρόβλημα, αν οι TID-λίστες είναι μεγάλες και δε χωρούν στη μνήμη

Εξέλιξη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 29

### Ο Αλγόριθμος FP-Growth

A

AC

ABC

B

AC

ABC

Εξέλιξη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 30

(a) Prefix tree

### Αλγόριθμος FP-Growth



Χρησιμοποιεί μια συμπίεσμένη αναπαράσταση της βάσης με τη μορφή ενός FP-δέντρου

- Το δέντρο μοιάζει με prefix tree (trie)
- Ο αλγόριθμος κατασκευής διαβάζει μια δοσοληψία τη φορά, απεικονίζει τη δοσοληψία σε ένα μονοπάτι του FP-δέντρου
- Μερικά μονοπάτια μπορεί να επικαλύπτονται: όσο περισσότερα μονοπάτια επικαλύπτονται, τόσο καλύτερη συμπίεση

Μόλις το FP-δέντρο κατασκευαστεί, ο αλγόριθμος χρησιμοποιεί μια αναδρομική διαίρει-και-βασίλευε (divide-and-conquer) προσέγγιση για την εξόρυξη των συχνών στοιχειοσυνόλων

### Αλγόριθμος FP-Growth



#### Κατασκευή FP-δέντρου

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

Το FP-δέντρο είναι ένα προθεματικό δέντρο

Επειδή έχουμε σύνολα, κάπως πρέπει να τα διατάξουμε ώστε να βρίσκουμε προθέματα

Δηλαδή δε μπορεί το ένα σύνολο να είναι {A, B} και το άλλο {B, C, A} γιατί χάνουμε το κοινό πρόθεμα AB (ή BA)

Άρα τα στοιχεία σε κάθε σύνολο πρέπει να ακολουθούν κάποια διάταξη, έστω τη λεξικογραφική (θα δούμε αργότερα αν κάτι άλλο συμφέρει καλύτερα)

Αρχικά, το δέντρο κενό  null

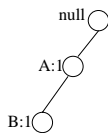
### Αλγόριθμος FP-Growth



#### Κατασκευή FP-δέντρου

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

Διάβασμα TID=1:



Κάθε κόμβος ετικέτα: ποιο στοιχείο και τη συχνότητα εμφάνισης (υποστήριξη) - πόσες δοσοληψίες φτάνουν σε αυτόν

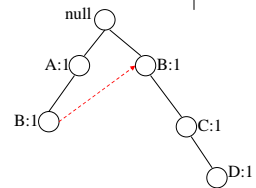
### Αλγόριθμος FP-Growth



#### Κατασκευή FP-δέντρου

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

Διάβασμα TID=1:



Διάβασμα TID=2:

Κάθε κόμβος ετικέτα: ποιο στοιχείο και τη συχνότητα εμφάνισης (υποστήριξη) - πόσες δοσοληψίες φτάνουν σε αυτόν

Επίσης, δείκτες μεταξύ των κόμβων που αναφέρονται στο ίδιο στοιχείο

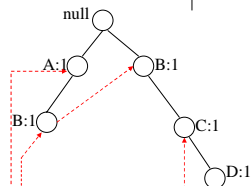
### Αλγόριθμος FP-Growth



#### Κατασκευή FP-δέντρου

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

Διάβασμα TID=1, 2:



Πίνακας Δεικτών

Item	Pointer
A	.....
B	.....
C	.....
D	.....
E	.....

Επίσης, κρατάμε δεικτες για να βοηθήσουν στον υπολογισμό των συχνών στοιχειοσυνόλων

### Αλγόριθμος FP-Growth

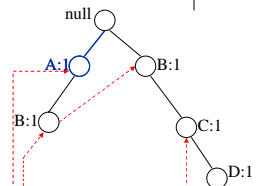


#### Κατασκευή FP-δέντρου

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

Διάβασμα TID=1, 2:

Διάβασμα TID=3:



Πίνακας Δεικτών

Item	Pointer
A	.....
B	.....
C	.....
D	.....
E	.....

### Αλγόριθμος FP-Growth

**Κατασκευή FP-δέντρου**

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

Διάβαση TID=1, 2:  
Διάβαση TID=3

Item	Pointer
A	.....
B	.....
C	.....
D	.....
E	.....

Πίνακας Δεικτών

37

### Αλγόριθμος FP-Growth

**Κατασκευή FP-δέντρου**

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

Διάβαση TID=1, 2:  
Διάβαση TID=3

Item	Pointer
A	.....
B	.....
C	.....
D	.....
E	.....

Πίνακας Δεικτών

38

### Αλγόριθμος FP-Growth

**Κατασκευή FP-δέντρου**

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

Item	Pointer
A	.....
B	.....
C	.....
D	.....
E	.....

Header table

39

### Αλγόριθμος FP-Growth

**Μέγεθος FP-δέντρου**

Κάθε δοσοληψία αντιστοιχεί σε ένα μονοπάτι από τη ρίζα

Το μέγεθος του δέντρου συνήθως μικρότερο των δεδομένων, αν υπάρχουν κοινά προθέματα

Αν όλες οι δοσοληψίες τα ίδια δεδομένα, μόνο ένα κλαδί

Αν όλες διαφορετικές, ο χώρος μεγαλύτερος (γιατί αποθηκεύεται περισσότερη πληροφορία, όπως δείκτες μεταξύ των κόμβων αλλά και συχνότητες εμφάνισης)

40

### Αλγόριθμος FP-Growth

**Κατασκευή FP-δέντρου**

Το τελικό δέντρο, εξαρτάται από τη διάταξη: άλλη διάταξη -> άλλα προθέματα

(Συνήθως) Μικρότερο δέντρο, αν όχι λεξικογραφικά, αλλά με βάση τη συχνότητα εμφάνισης -> Αρχικά, διαβάζουμε όλα τα δεδομένα μια φορά ώστε να υπολογιστεί ο μετρητής υποστήριξης κάθε στοιχείου, και διατάσσουμε τα στοιχεία με βάση αυτό

Για το παράδειγμα,  $\sigma(A)=7, \sigma(B)=8, \sigma(C)=7, \sigma(D)=5, \sigma(E)=3$

Άρα, διάταξη B, A, C, D, E

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

TID	Items
1	{B,A}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{B,A,C}
6	{B,A,C,D}
7	{B,C}
8	{B,A,C}
9	{B,A,D}
10	{B,C,E}

41

### Αλγόριθμος FP-Growth

**Χρήση FP-δέντρου για εύρεση συχνών στοιχειοσυνόλων**

TID	Items
1	{A,B}
2	{B,C,D}
3	{A,C,D,E}
4	{A,D,E}
5	{A,B,C}
6	{A,B,C,D}
7	{B,C}
8	{A,B,C}
9	{A,B,D}
10	{B,C,E}

Item	Pointer
A	.....
B	.....
C	.....
D	.....
E	.....

Header table

Πως:  
Bottom-up traversal του δέντρου  
Αυτά που τελειώνουν σε E, μετά αυτά που τελειώνουν σε D, C, B και τέλος A - suffix-based classes

42

**Αλγόριθμος FP-Growth**

Υποπρόβλημα: Βρες συχνά στοιχεία που τελειώνουν σε **E**

Item	Pointer
A	A:7
B	B:5
C	C:3
D	D:1
E	E:1

Θα δούμε στη συνέχεια πως υπολογίζεται η υποστήριξη για τα πιθανά στοιχειοσύνολα

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 43

**Αλγόριθμος FP-Growth**

Για το **D**

Item	Pointer
A	A:7
B	B:5
C	C:3
D	D:1
E	E:1

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 44

**Αλγόριθμος FP-Growth**

Για το **C**

Item	Pointer
A	A:7
B	B:5
C	C:3
D	D:1
E	E:1

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 45

**Αλγόριθμος FP-Growth**

Για το **B**

Item	Pointer
A	A:7
B	B:5
C	C:3
D	D:1
E	E:1

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 46

**Αλγόριθμος FP-Growth**

Για το **A**

Item	Pointer
A	A:7
B	B:5
C	C:3
D	D:1
E	E:1

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 47

**Αλγόριθμος FP-Growth**

**Ώαση 1**  
Όλα τα μονοπάτια που περιέχουν το E

Προθεματικά Μονοπάτια (prefix paths)

Προθεματικά μονοπάτια του E:  
{E}, {D,E}, {C,D,E}, {A,D,E}, {A,C,D,E}, {C,E}, {B,C,E}

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 48



**Αλγόριθμος FP-Growth**

**Ξάση 1**  
Όλα τα μονοπάτια που περιέχουν το E  
Προθεματικά Μονοπάτια (prefix paths)

**Header table**

Item	Pointer
A	
B	
C	
D	
E	

Προθεματικά μονοπάτια του E:  
{E}, {D,E}, {C,D,E}, {A,D,E}, {A,C,D,E}, {C,E}, {B,C,E}

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 49

**Αλγόριθμος FP-Growth**

**Ξάση 1**  
Όλα τα μονοπάτια που περιέχουν το E  
Προθεματικά Μονοπάτια (prefix paths)

Προθεματικά μονοπάτια του E:  
{E}, {D,E}, {C,D,E}, {A,D,E}, {A,C,D,E}, {C,E}, {B,C,E}

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 50

**Αλγόριθμος FP-Growth**

Έστω  $\text{minsup} = 2$   
Βρες την υποστήριξη του {E}  
Πως:  
Ακολούθησε τους συνδέσμους αθροίζοντας  $1+1+1=3 > 2$   
Οπότε {E} συχνό

{E} συχνό άρα προχωράμε για DE, CE, BE, AE

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 51

**Αλγόριθμος FP-Growth**

Μετατροπή των προθεματικών δέντρων σε FP-δέντρο υπό συνθήκες (conditional FP-tree)  
Δύο αλλαγές  
(1) Αλλαγή των μετρητών  
(2) Περικοπή

{E} συχνό άρα προχωράμε για DE, CE, BE, AE

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 52

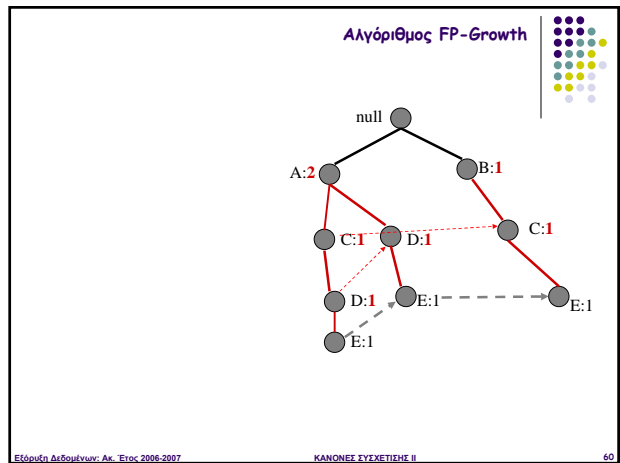
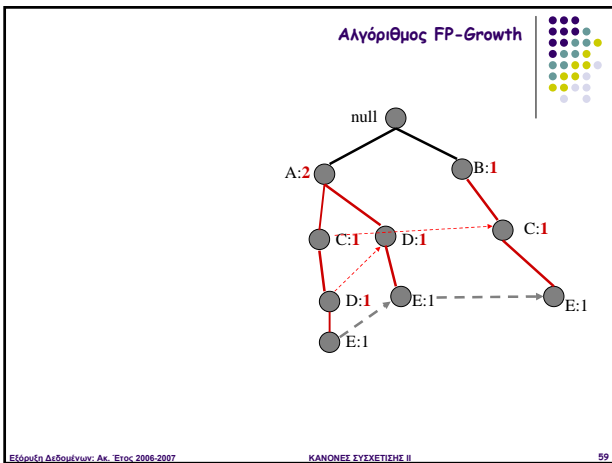
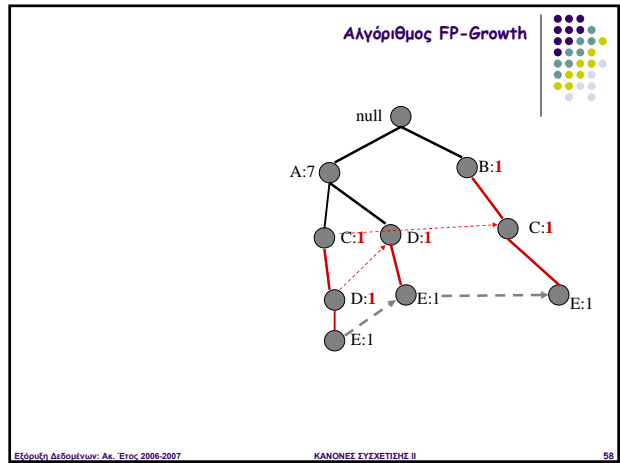
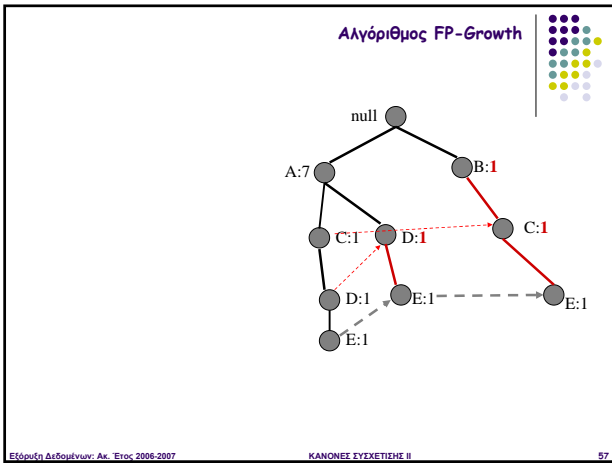
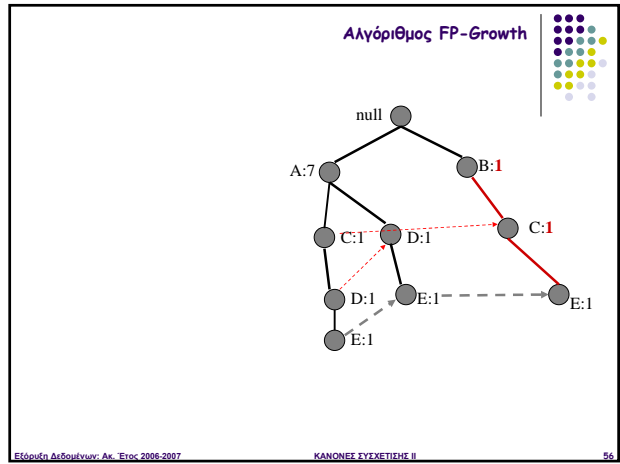
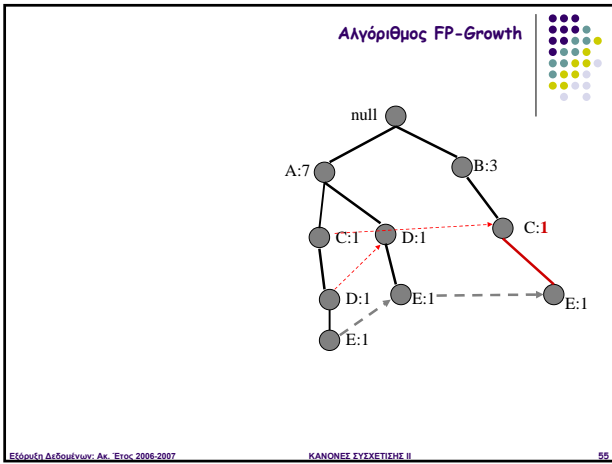
**Αλγόριθμος FP-Growth**

Αλλαγή μετρητών  
Οι μετρητές σε κάποιους κόμβους περιλαμβάνουν ποσολλήμεις που δεν έχουν το E  
Πχ στο  $\text{null} \rightarrow B \rightarrow C \rightarrow E$  μετράμε και την {B,C}

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 53

**Αλγόριθμος FP-Growth**

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 54



Αλγόριθμος FP-Growth

Περικοπή (truncate)  
Σβήσε τους κόμβους του E

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 61

Αλγόριθμος FP-Growth

Περικοπή (truncate)  
Σβήσε τους κόμβους του E

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 62

Αλγόριθμος FP-Growth

Περικοπή (truncate)  
Σβήσε τους κόμβους του E

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 63

Αλγόριθμος FP-Growth

Πιθανή περαιτέρω περικοπή  
Κάποια στοιχεία μπορεί να έχουν υποστήριξη μικρότερη της ελάχιστης  
Πχ το B → περικοπή

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 64

Αλγόριθμος FP-Growth

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 65

Αλγόριθμος FP-Growth

Εξώφυλλο Διδασκαλίας: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 66

**Αλγόριθμος FP-Growth**

Υπο-συνθήκη FP-δέντρο για το E  
 Ο αλγόριθμος επαναλαμβάνεται για το {D, E}, {C, E}, {A, E}

Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 67

**Αλγόριθμος FP-Growth**

**Ξάση 1**  
 Όλα τα μονοπάτια που περιέχουν το D (DE)  
 Προθεματικά Μονοπάτια (prefix paths)

Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 68

**Αλγόριθμος FP-Growth**

**Ξάση 1**  
 Όλα τα μονοπάτια που περιέχουν το D (DE)  
 Προθεματικά Μονοπάτια (prefix paths)

Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 69

**Αλγόριθμος FP-Growth**

Βρες την υποστήριξη του {D, E}  
 Πως:  
 Ακολουθήσε τους συνδέσμους αθροίζοντας  $1+1=2 \geq 2$   
 Οπότε {D, E} συχνό

Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 70

**Αλγόριθμος FP-Growth**

Κατασκεύασε το υπο-συνθήκη FP-δέντρο για το {D, E}

1. Αλλαγή υποστήριξης
2. Περικοπές κόμβων

Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 71

**Αλγόριθμος FP-Growth**

1. Αλλαγή υποστήριξης

Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 72

**Αλγόριθμος FP-Growth**

2. Περικοπές κόμβων

Εύρωπη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 73

**Αλγόριθμος FP-Growth**

2. Περικοπές κόμβων

Εύρωπη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 74

**Αλγόριθμος FP-Growth**

2. Περικοπές κόμβων

Εύρωπη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 75

**Αλγόριθμος FP-Growth**

Τελικό υπο-συνθήκη FP-δέντρο για το {D, E}

Υποστήριξη του A είναι  $\geq \text{minsup} \rightarrow \{A, D, E\}$  συχνό  
Αφού μόνο έναν κόμβο, επιστροφή στο επόμενο υποπρόβλημα

Εύρωπη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 76

**Αλγόριθμος FP-Growth**

Υπο-συνθήκη FP-δέντρο για το E  
Ο αλγόριθμος επαναλαμβάνεται για το {D, E}, {C, E}, {A, E}

Εύρωπη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 77

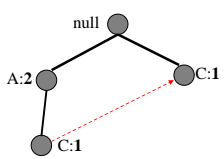
**Αλγόριθμος FP-Growth**

**Ψάξη 1**  
Όλα τα μονοπάτια που περιέχουν το C (CE)  
Προθεματικά Μονοπάτια (prefix paths)

Εύρωπη Διδασκόντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 78

**Αλγόριθμος FP-Growth**

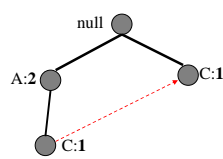
**Ψάξη 1**  
 Όλα τα μονοπάτια που περιέχουν το C (CE)  
 Προθεματικά Μονοπάτια (prefix paths)



Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 79

**Αλγόριθμος FP-Growth**

Βρες την υποστήριξη του {C, E}  
 Πως:  
 Ακολούθησε τους συνδέσμους  
 αθροίζοντας  $1+1=2 \geq 2$   
 Οπότε {C, E} συχνά

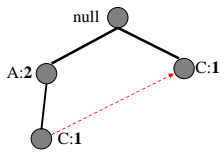


Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 80

**Αλγόριθμος FP-Growth**

Κατασκεύασε το υπο-συνθήκη FP-δέντρο για το {C, E}

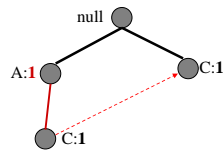
1. Αλλαγή υποστήριξης
2. Περικοπές κόμβων



Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 81

**Αλγόριθμος FP-Growth**

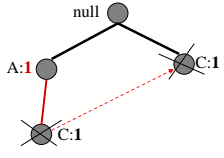
1. Αλλαγή υποστήριξης



Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 82

**Αλγόριθμος FP-Growth**

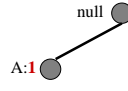
2. Περικοπή Κόμβων



Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 83

**Αλγόριθμος FP-Growth**

2. Περικοπή Κόμβων



Εξώφυλλο Διδασκάλων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 84

**Αλγόριθμος FP-Growth**

2. Περικοπή Κόμβων

Εύρωτη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 85

**Αλγόριθμος FP-Growth**

2. Περικοπή Κόμβων

null

Άρα, επιστροφή στο επόμενο υποπρόβλημα

Εύρωτη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 86

**Αλγόριθμος FP-Growth**

Υπο-συνθήκη FP-δέντρο για το E  
Ο αλγόριθμος επαναλαμβάνεται για το {B, E}, {C, E}, {A, E}

Εύρωτη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 87

**Αλγόριθμος FP-Growth**

**Ψάση 1**  
Όλα τα μονοπάτια που περιέχουν το A (AE)  
Προθεματικά Μονοπάτια (prefix paths)

Εύρωτη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 88

**Αλγόριθμος FP-Growth**

**Ψάση 1**  
Όλα τα μονοπάτια που περιέχουν το A (AE)  
Προθεματικά Μονοπάτια (prefix paths)

Εύρωτη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 89

**Αλγόριθμος FP-Growth**

Βρες την υποστήριξη του {A, E}  
Οπότε {A, E} συχνό

Δε χρειάζεται να φτιάξουμε υπο-συνθήκη FP-δέντρο για το {A, E}

Εύρωτη Διδασκντων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 90





1. Αλλαγή υποστήριξης

Αλγόριθμος FP-Growth

Εύρηνη Διδασκάλου: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΙΣΗΣ II 97

1. Αλλαγή υποστήριξης

Αλγόριθμος FP-Growth

Εύρηνη Διδασκάλου: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΙΣΗΣ II 98

2. Περικοπή Κόμβων

Αλγόριθμος FP-Growth

Εύρηνη Διδασκάλου: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΙΣΗΣ II 99

2. Περικοπή Κόμβων

Αλγόριθμος FP-Growth

Εύρηνη Διδασκάλου: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΙΣΗΣ II 100

Προθεματικά δέντρα και υποσυνθήκη δέντρα  
Για τα AD, BD και CD κοκ

Αλγόριθμος FP-Growth

Εύρηνη Διδασκάλου: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΙΣΗΣ II 101

Αλγόριθμος FP-Growth

Παράδειγμα τεχνικής διαιρεί-και-βασίλευε

Σε κάθε αναδρομικό βήμα, λύνεται και ένα υπο-πρόβλημα:

- Κατασκευάζεται το προθεματικό δέντρο
- Υπολογίζεται η νέα υποστήριξη για τους κόμβους του
- Περικόβονται οι κόμβοι με μικρή υποστήριξη

Επειδή τα υποπρόβλήματα είναι ξένα μεταξύ τους, δεν δημιουργούνται τα ίδια συχνά στοιχειοσύνολα δυο φορές

Ο υπολογισμός της υποστήριξης είναι αποδοτικός - γίνεται ταυτόχρονα με τη δημιουργία των συχνών στοιχειοσυνόλων

Εύρηνη Διδασκάλου: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΧΕΤΙΣΗΣ II 102

## Αλγόριθμος FP-Growth



Η απόδοση του FP-Growth εξαρτάται από τον παράγοντα συμπίεσης του συνόλου των δεδομένων (compression factor)

Αν τα τελικά δέντρα είναι «θαμνώδη» (bushy) τότε δε δουλεύει καλά, αυξάνεται ο αριθμός των υποπροβλημάτων (οι αναδρομικές κλήσεις)

## Επίδραση της Υποστήριξης



## Κατανομή Υποστήριξης



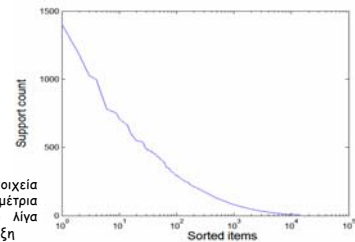
- Η απόδοση των αλγορίθμων εξαρτάται από τα δεδομένα εισόδου, πχ ο αριθμός από τον αριθμό των στοιχείων, το πλάτος των δασοληψιών, ο FP-Growth από την τομή των δασοληψιών
- Επίσης, από την τιμή της ελάχιστης υποστήριξης (minsup). Πώς θα προσδιοριστεί μια κατάλληλη τιμή για το *minsup*:
  - Αν η τιμή είναι πολύ υψηλή, μπορεί να χαθούν στοιχειοσύνολα που περιέχουν ενδιαφέροντα σπάνια στοιχεία (πχ ακριβά προϊόντα)
  - Αν η τιμή είναι πολύ χαμηλή, οι μέθοδοι γίνονται ακριβοί γιατί ο αριθμός των υποψηφίων στοιχειοσυνόλων είναι πολύ μεγάλος
  - Ο αριθμός των συχνών στοιχειοσυνόλων γίνεται πολύ μεγάλος

## Κατανομή Υποστήριξης



Επιπρόσθετα, η χρήση μόνο μιας ελάχιστης υποστήριξης μπορεί να μην αρκεί  
Για πολλά πραγματικά δεδομένα η κατανομή της υποστήριξης δεν είναι ομοιόμορφη (skewed support distribution)

Support distribution of a retail data set



Τα περισσότερα στοιχεία έχουν μικρή ή μέτρια υποστήριξη και μόνο λίγα έχουν μεγάλη υποστήριξη

## Κατανομή Υποστήριξης



Ομάδα	G1	G2	G3
Υποστήριξη	<1%	1%-90%	>90%
Αριθμός στοιχείων	1735	358	20

Παράδειγμα κανόνες μεταξύ G1 και G3 (χαβιάρι και γάλα)  
Cross-support patterns

## Πολλαπλές Τιμές Υποστήριξης



### Πολλαπλές Ελάχιστες Τιμές Υποστήριξης

$MS(i)$ : ελάχιστη υποστήριξη για το στοιχείο  $i$

- Π.χ.:  $MS(\text{Milk})=5\%$ ,  $MS(\text{Coke}) = 3\%$ ,  
 $MS(\text{Broccoli})=0.1\%$ ,  $MS(\text{Salmon})=0.5\%$
- $MS(\{\text{Milk}, \text{Broccoli}\}) = \min(MS(\text{Milk}), MS(\text{Broccoli})) = 0.1\%$

Πρόβλημα: Η υποστήριξη παύει να είναι αντιμονότονη:

- Έστω:  $\text{Support}(\text{Milk}, \text{Coke}) = 1.5\%$  and  $\text{Support}(\text{Milk}, \text{Coke}, \text{Broccoli}) = 0.5\%$
- $\{\text{Milk}, \text{Coke}\}$  είναι μη συχνό αλλά το  $\{\text{Milk}, \text{Coke}, \text{Broccoli}\}$  είναι συχνό

Λόγω του Broccoli που κατεβάζει το minsup

## Πολλαπλές Τιμές Υποστήριξης

### Multiple Minimum Support (Liu 1999)

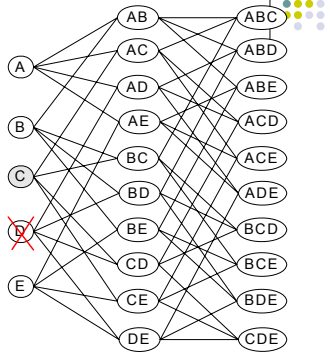
- Ταξινόμηση τα στοιχεία με βάση την ελάχιστη τιμή υποστήριξης (σε αύξουσα διάταξη)
  - π.χ.:  $MS(\text{Milk})=5\%$ ,  $MS(\text{Coke})=3\%$ ,  
 $MS(\text{Broccoli})=0.1\%$ ,  $MS(\text{Salmon})=0.5\%$
  - Διάταξη: Broccoli, Salmon, Coke, Milk
- Τροποποίηση του Αργιορί (Βήμα Φαλιδισμού):
  - $L_1$ : σύνολο συχνών στοιχειοσυνόλων
  - $F_1$ : σύνολο στοιχείων που η υποστήριξη τους είναι  $\geq MS(1)$  όπου  $MS(1)$  είναι  $\min_i(MS(i))$
  - $C_2$ : τα υποψήφια στοιχειοσύνολα μεγέθους 2 παράγονται από το  $F_1$  αντί του  $L_1$

## Πολλαπλές Τιμές Υποστήριξης

- Τροποποιήσεις στον Αργιορί (Βήμα Φαλιδισμού):
  - Στον παραδοσιακό Αργιορί,
    - Ένα υποψήφιο  $(k+1)$ -στοιχειοσύνολο δημιουργείται συγχωνεύοντας δύο συχνά  $k$ -στοιχειοσύνολα
    - Το υποψήφιο φαλιδίζεται αν περιέχει ένα (οποιοδήποτε) μη συχνό  $k$ -στοιχειοσύνολο
  - Τροποποίηση βήματος φαλιδισμού:
    - Φαλιδίσει μόνο αν το υποσύνολο περιέχει το πρώτο στοιχείο π.χ. Candidate={Broccoli, Coke, Milk} (διατεταγμένα με βάση την μικρότερη ελάχιστη υποστήριξη) {Broccoli, Coke} και {Broccoli, Milk} είναι συχνά αλλά {Coke, Milk} είναι μη συχνό
      - Candidate δε σβήνεται γιατί το {Coke, Milk} δεν περιέχει το πρώτο στοιχείο, δηλαδή, Broccoli.

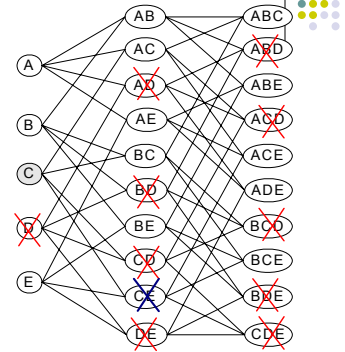
## Πολλαπλές Τιμές Υποστήριξης

Item	MS(i)	Sup(i)
A	0.10%	0.25%
B	0.20%	0.26%
C	0.30%	0.29%
D	0.50%	0.05%
E	3%	4.20%



## Πολλαπλές Τιμές Υποστήριξης

Item	MS(i)	Sup(i)
A	0.10%	0.25%
B	0.20%	0.26%
C	0.30%	0.29%
D	0.50%	0.05%
E	3%	4.20%



## Ακολουθιακά Δεδομένα

Μέχρι στιγμής, οι δοσοληψίες *σύνολα από στοιχεία*, δεν έχει σημασία η σειρά εμφάνισης των στοιχείων σε κάθε δοσοληψία - επίσης, *σύνολα από δοσοληψίες*, δεν έχει σημασία η σειρά εμφάνισης κάθε δοσοληψίας

Ωστόσο, πολλά δεδομένα στο «Καλάθι της νοικοκυράς» περιέχουν χρονική πληροφορία, π.χ. ποιες δοσοληψίες κάνει ένας πελάτης σε μια συγκεκριμένη χρονική περίοδο

Επίσης, γεγονότα που είναι αποτελέσματα επιστημονικών πειραμάτων

Δηλαδή, σχέση διάταξης είτε χρονική, είτε χωρική

### Ακολουθίες

**Sequence Database (Ακολουθιακή Βάση Δεδομένων):**

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

Ταξινόμηση με βάση τη χρονιά → ακολουθία

Γεγονότα (Events) (~στοιχεία) σχετιζόμενα με Αντικείμενα (Objects) (~όσοοψηφίες) και τότε αυτά εμφανίζονται

Εύρεση Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 115

### Ακολουθίες

**Παραδείγματα Ακολουθιακών Δεδομένων**

- Ακολουθία προσπελάσεων Web:
  - < {Homepage} {Electronics} {Digital Cameras} {Canon Digital Camera} {Shopping Cart} {Order Confirmation} {Return to Shopping} >
- Ακολουθία γεγονότων που οδήγησαν σε πυρηνικό ατύχημα στο 3-mile Island: ([http://stellar-one.com/nuclear/staff\\_reports/summary\\_SOE\\_the\\_initiating\\_event.htm](http://stellar-one.com/nuclear/staff_reports/summary_SOE_the_initiating_event.htm))
  - < {clogged resin} {outlet valve closure} {loss of feedwater} {condenser polisher outlet valve shut} {booster pumps trip} {main waterpump trips} {main turbine trips} {reactor pressure increases}>
- Ακολουθία βιβλίων δανεισμού από βιβλιοθήκη:
  - < {Fellowship of the Ring} {The Two Towers} {Return of the King}>

Εύρεση Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 116

### Ακολουθίες

**Τυπικός Ορισμός Ακολουθίας**

Μια ακολουθία (sequence) είναι μια διατεταγμένη λίστα από στοιχεία (elements) (~transactions)

$$s = \langle e_1 e_2 e_3 \dots \rangle$$

Κάθε στοιχείο αποτελείται από μια συλλογή από γεγονότα (events) (~items)

$$e_i = \{i_1, i_2, \dots, i_k\}$$

Κάθε στοιχείο αντιστοιχεί σε μια συγκεκριμένη χρονική στιγμή ή τοποθεσία

Μήκος (length) μιας ακολουθίας, |s|, είναι ο αριθμός των στοιχείων της ακολουθίας

Μια k-ακολουθία είναι μια ακολουθία που περιέχει k γεγονότα (items)

Στοιχεία (Transaction) E1 E2 E3 E4 E2 E3 E4 Γεγονός (Item)

Ακολουθία

Μήκος = 5, 8-ακολουθία

Εύρεση Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 117

### Ακολουθίες

**Παραδείγματα Ακολουθιακών Δεδομένων**

Sequence Database	Sequence	Element (Transaction)	Event (Item)
Customer	Purchase history of a given customer	A set of items bought by a customer at time t	Books, diary products, CDs, etc
Web Data	Browsing activity of a particular Web visitor	A collection of files viewed by a Web visitor after a single mouse click	Home page, index page, contact info, etc
Event data	History of events generated by a given sensor	Events triggered by a sensor at time t	Types of alarms generated by sensors
Genome sequences	DNA sequence of a particular species	An element of the DNA sequence	Bases A, T, G, C

Χρονική διάταξη

Χωρική διάταξη

Element (Transaction) E1 E2 E3 E4 Event (Item)

Sequence

Εύρεση Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 118

Object A:

### Ακολουθίες

**Τυπικός Ορισμός Υπο-ακολουθίας**

Μια ακολουθία  $\langle a_1 a_2 \dots a_n \rangle$  περιέχεται σε μια άλλη ακολουθία (είναι υπο-ακολουθία της)  $\langle b_1 b_2 \dots b_m \rangle$  ( $m \geq n$ ) αν υπάρχουν ακέραιοι  $i_1 < i_2 < \dots < i_n$  τέτοιοι ώστε  $a_1 \subseteq b_{i_1}, a_2 \subseteq b_{i_2}, \dots, a_n \subseteq b_{i_n}$

Παραδείγματα

Ακολουθία Δεδομένων	Υπο-ακολουθία	Περιέχεται;
$\langle \{2,4\} \{3,5,6\} \{8\} \rangle$	$\langle \{2\} \{3,5\} \rangle$	Ναι
$\langle \{1,2\} \{3,4\} \rangle$	$\langle \{1\} \{2\} \rangle$	Όχι
$\langle \{2,4\} \{2,4\} \{2,5\} \rangle$	$\langle \{2\} \{4\} \rangle$	Ναι

Εύρεση Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 119

### Ακολουθίες

**Εύρεση Ακολουθιακών Προτύπων**

Έστω D ένα σύνολο που περιέχει μια ή περισσότερες ακολουθίες.

Η υποστήριξη (support) μιας ακολουθίας w ορίζεται ως το ποσοστό των ακολουθιών στο D που περιέχουν το w

Σύνολο 5 ακολουθιών

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

Παραδείγματα:  $\langle \{1,2\} \rangle$   $s=60\%$

Εύρεση Δεδομένων: Ακ. Έτος 2006-2007 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 120

line

10

2  
3  
5

Object B:

4  
5  
6

Object C:

## Εύρεση Ακολουθιακών Προτύπων

Έστω D ένα σύνολο που περιέχει μια ή περισσότερες ακολουθίες.

Η υποστήριξη (support) μιας ακολουθίας w ορίζεται ως το ποσοστό των ακολουθιών στο D που περιέχουν το w

Σύνολο 5 ακολουθιών

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
D	1	2,4,5
D	2	3,4
D	3	4,5
E	1	1,3
E	2	2,4,5

Παραδείγματα:

$\langle \{3\}\{5\} \rangle$  s=80%  
 $\langle \{1,2\}\{5\} \rangle$  s=40%  
 $\langle \{5\}\{1,2\} \rangle$  s=0%

Επίσης, οποιαδήποτε ακολουθία με μήκος μεγαλύτερο του 2

## Εύρεση Ακολουθιακών Προτύπων

Ένα ακολουθιακό πρότυπο (sequential pattern) είναι μια συχνή υπο-ακολουθία (δηλαδή μια ακολουθία με υποστήριξη  $\geq$  minsup)

Minsup = 50%

Παραδείγματα συχνών υπο-ακολουθιών

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

$\langle \{1,2\} \rangle$  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%

## Εύρεση Ακολουθιακών Προτύπων

### Ορισμός Προβλήματος Εξόρυξης Ακολουθιακών Προτύπων (Sequential Pattern Mining)

Είσοδος:

Μια βάση από ακολουθίες  
Ένα ελάχιστο κατώφλι υποστήριξης, minsup

Πρόβλημα:

Βρες όλες τις υπο-ακολουθίες με υποστήριξη  $\geq$  minsup

## Εύρεση Ακολουθιακών Προτύπων

Έστω μια ακολουθία:  $\langle \{a\} \{c\} \{d\} \{e\} \{f\} \{g\} \{h\} \{i\} \rangle$

Παραδείγματα υπο-ακολουθίας:

$\langle \{a\} \{c\} \{f\} \{g\} \rangle$ ,  $\langle \{c\} \{d\} \{e\} \rangle$ ,  $\langle \{b\} \{g\} \rangle$ , etc  
Ο αριθμός τους είναι εκθετικός

Πόσες k-υποακολουθίες μπορεί να εξαχθούν από μια n-ακολουθία:

$\langle \{a\} \{b\} \{c\} \{d\} \{e\} \{f\} \{g\} \{h\} \{i\} \rangle$  n = 9

Παράδειγμα για k = 4

$\{a\}\{c\}\{f\}\{g\}$

$\{a\}$ ,  $\{f\}$ ,  $\{h\} \{i\}$

$\{c\} \{d\} \{e\} \{g\}$

$$\binom{n}{k} = \binom{9}{4} = 126$$

$$\sum_{i=1,9} \binom{9}{i} = 511$$

## Εύρεση Ακολουθιακών Προτύπων

### Brute-Force Μέθοδος

Απαρίθμηση όλων των πιθανών υπο-ακολουθιών και υπολογισμός της υποστήριξής τους

▪ Έστω n γεγονότα:  $i_1, i_2, i_3, \dots, i_n$

▪ Υποψήφιος 1-υπο-ακολουθίες:

▪  $\langle \{i_1\} \rangle, \langle \{i_2\} \rangle, \langle \{i_3\} \rangle, \dots, \langle \{i_n\} \rangle$

▪ Υποψήφιος 2-υπο-ακολουθίες:

▪  $\langle \{i_1, i_2\} \rangle, \langle \{i_1, i_3\} \rangle, \dots, \langle \{i_1\} \{i_2\} \rangle, \langle \{i_1\} \{i_3\} \rangle, \dots, \langle \{i_n\} \{i_1\} \rangle, \dots$

▪ Υποψήφιος 3-υπο-ακολουθίες:

▪  $\langle \{i_1, i_2, i_3\} \rangle, \langle \{i_1, i_2, i_4\} \rangle, \dots, \langle \{i_1, i_2\} \{i_3\} \rangle, \langle \{i_1, i_2\} \{i_4\} \rangle, \dots, \langle \{i_1\} \{i_1, i_2\} \rangle, \langle \{i_1\} \{i_1, i_3\} \rangle, \dots, \langle \{i_1\} \{i_1\} \{i_2\} \rangle, \langle \{i_1\} \{i_1\} \{i_3\} \rangle, \dots$

## Εύρεση Ακολουθιακών Προτύπων



Παρατήρηση

Ο αριθμός των υποψηφίων ακολουθιών είναι πολύ μεγαλύτερος από τον αριθμό των υποψηφίων στοιχειοσύνολων για δύο κυρίως λόγους:

- Ένα στοιχείο μπορεί να εμφανιστεί μόνο μια φορά σε ένα στοιχειοσύνολο, ενώ ένα γεγονός μπορεί να εμφανιστεί περισσότερο από μια φορά σε μια ακολουθία

Πχ το στοιχειοσύνολο  $\{i_1, i_2\} \rightarrow \langle \{i_1, i_1\} \rangle \langle \{i_1\}, \{i_1\} \rangle$ , κλπ

- Η διάταξη δεν έχει σημασία στα στοιχειοσύνολα, αλλά έχει στις ακολουθίες

Πχ το στοιχειοσύνολο  $\{i_1, i_2\} \rightarrow \langle \{i_1\}, \{i_2\} \rangle \langle \{i_2\}, \{i_1\} \rangle$ , κλπ



Αρτιοί για ακολουθίες

Μια  $k$ -ακολουθία πρέπει να περιέχει όλες τις πιθανές  $k-1$  υπο-ακολουθίες της

Οπότε, παρόμοιος αλγόριθμος



Βήμα 1:

- Κάνε το πρώτο πέρασμα στη βάση των ακολουθιών  $D$  και παράγνε όλες τις συχνές ακολουθίες ενός στοιχείου

Βήμα 2:

Επανάλαβε μέχρι να μην παράγονται νέες συχνές ακολουθίες

- **Δημιουργία Υποψηφίων - Candidate Generation:**
  - Συγχώνευση συχνών ακολουθιών που βρέθηκαν στο  $(k-1)$ η πέρασμα για δημιουργία υποψηφίων ακολουθιών με  $k$  στοιχεία
- **Ψαλίδισμα Υποψηφίων - Candidate Pruning:**
  - Ψαλίδισε τις  $k$ -ακολουθίες που περιέχουν μη συχνές Prune candidate  $(k-1)$ -υπο-ακολουθίες
- **Υπολογισμός Υποστήριξης - Support Counting:**
  - Κάνε ένα νέο πέρασμα στη βάση  $D$  για τον υπολογισμό της υποστήριξης των νέων υποψηφίων
- **Υπολογισμός Υποψηφίων - Candidate Elimination:**
  - Διώξε τις υποψήφιες  $k$ -ακολουθίες που η πραγματική τους υποστήριξη είναι μικρότερη του  $\text{minsup}$



Δημιουργία υποψηφίων

- Βάση ( $k=2$ ):  
Συγχώνευση δύο συχνών 1-ακολουθιών  $\langle \{i_1\} \rangle$  and  $\langle \{i_2\} \rangle$  θα παράξει δύο υποψήφιες 2-ακολουθίες:  $\langle \{i_1\}, \{i_2\} \rangle$  and  $\langle \{i_1\}, \{i_2\} \rangle$
- Γενική περίπτωση ( $k>2$ ):  
**Συνθήκη για συγχώνευση:** Μια συχνή  $(k-1)$ -ακολουθία  $w_1$  συγχωνεύεται με μια άλλη συχνή  $(k-1)$ -ακολουθία  $w_2$  για να παραχθεί μια υποψήφια  $k$ -ακολουθία αν η υπο-ακολουθία που παίρνουμε αν *αφήσουμε το πρώτο γεγονός της  $w_1$*  είναι το ίδιο με την υπο-ακολουθία που παίρνουμε αν *αφήσουμε το τελευταίο γεγονός της  $w_2$*   
Το **αποτέλεσμα** μετά τη συγχώνευση είναι η ακολουθία  $w_1$  επεκταμένη με το τελευταίο γεγονός της  $w_2$ .
  - Αν τα τελευταία δύο γεγονότα της  $w_2$  ανήκουν στο ίδιο στοιχείο τότε το τελευταίο γεγονός της  $w_2$  γίνεται μέρος του τελευταίου στοιχείου της  $w_1$
  - Αλλιώς, το τελευταίο γεγονός της  $w_2$  γίνεται ένα διαφορετικό στοιχείο appended στο τέλος της  $w_1$

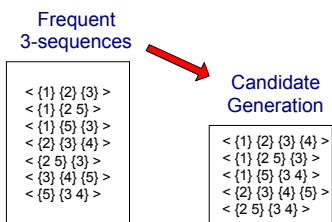


Παραδείγματα Δημιουργίας Υποψηφίων

- Συγχώνευση των ακολουθιών  $w_1 = \langle \{ \} \{2\} \{3\} \{4\} \rangle$  και  $w_2 = \langle \{2\} \{3\} \{4\} \{ \} \rangle$  μας δίνει την υποψήφια ακολουθία  $\langle \{1\} \{2\} \{3\} \{4\} \{5\} \rangle$  γιατί τα 2 τελευταία γεγονότα της  $w_2$  (4 και 5) ανήκουν στο ίδιο στοιχείο
- Συγχώνευση των ακολουθιών  $w_1 = \langle \{ \} \{2\} \{3\} \{4\} \rangle$  και  $w_2 = \langle \{2\} \{3\} \{4\} \{ \} \rangle$  μας δίνει την υποψήφια ακολουθία  $\langle \{1\} \{2\} \{3\} \{4\} \{5\} \rangle$  γιατί τα 2 τελευταία γεγονότα της  $w_2$  (4 και 5) δεν ανήκουν στο ίδιο στοιχείο
- Δε χρειάζεται να συγχωνεύσουμε τις ακολουθίες  $w_1 = \langle \{ \} \{2\} \{6\} \{4\} \rangle$  και  $w_2 = \langle \{1\} \{2\} \{4\} \{ \} \rangle$  για να πάρουμε το υποψήφιο  $\langle \{1\} \{2\} \{6\} \{4\} \{5\} \rangle$  συγχώνευση  $w_1$  με  $\langle \{2\} \{6\} \{4\} \{5\} \rangle$



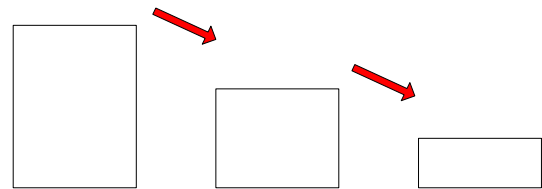
Παραδείγματα Δημιουργίας Υποψηφίων





### Ψαλίδισμα υποψηφίων

Μια υποψήφια  $k$ -ακολουθία σβήνεται αν έχει τουλάχιστον μια μη συχνή ( $k-1$ )-υπο-ακολουθία



Τέλος