

## 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 (transaction)
- στοιχείο (item)
- Παρουσία 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 (συχνό) στοιχειοσύνολο, δημιουργήσε κανόνες με μεγάλη υποστήριξη, όπου κάθε κανόνες είναι μια δυαδική διαμερίση (δηλ. χωρισμός στα δύο) του συχνού στοιχειοσυνόλου

### Εύρεση Συχνών Στοιχειοσυνόλων

**Itemset Lattice - Πλέγμα Στοιχειοσυνόλων**

Όλα τα δυνατά στοιχειοσύνολα έχουμε 5 στοιχεία

Για d στοιχεία,  $2^d$  πιθανά στοιχειοσύνολα

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 7

### Εύρεση Συχνών Στοιχειοσυνόλων: Στρατηγική αργίορι

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

Η ισοδύναμη αν ένα στοιχειοσύνολο είναι μη συχνό, όλα τα υπερασύνολα του είναι μη συχνά

βρέθηκε μη συχνό

ψαλιδισμένα υπερασύνολα

Support-based pruning  
Ψαλίδισμα με βάση την υποστήριξη

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 8

### Στρατηγική αργίορι

**Γενικός Αλγόριθμος για την Εύρεση Συχνών Στοιχειοσυνόλων**

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 δοσοληψιών
4. Σβήσε τα υποψήφια στοιχειοσύνολα που δεν είναι συχνά
5.  $k = k + 1$

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 9

### Στρατηγική αργίορι: Δημιουργία Στοιχειοσυνόλων

Για την παραγωγή υποψηφίων 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) στοιχειοσυνόλων αν τα πρώτα k-2 στοιχεία τους είναι τα ίδια

Για να αποφύγουμε τη δημιουργία του ίδιου στοιχειοσυνόλου, κρατάμε κάθε στοιχειοσύνολο (λεξιλογιαφικά) ταξινομημένο

**Ψαλίδισμα**

- 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

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 10

### Στρατηγική αργίορι: Υπολογισμός Υποστήριξης

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

Σε κάθε βήμα k+1

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

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 11

### Παραγωγή Κανόνων

Δοθέντος ενός συχνού στοιχειοσυνόλου L, βρες όλα τα μη κενά υποσύνολα  $f \subseteq L$  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

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 12

### Παραγωγή Κανόνων για τον Αλγόριθμο αργιορί

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

Ψαλίδισμα με βάση την εμπιστοσύνη

Εστω κόμβος με μικρή εμπιστοσύνη

Ψαλίδισμα μένει κανόνες

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

Υπολογίζουμε την εμπιστοσύνη

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

Σημείωση: Για τον υπολογισμό της εμπιστοσύνης δεν χρειάζεται να διαπεράσουμε τη βάση

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 13

### Αναπαράσταση Στοιχειοσυνόλων

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

Τι να κρατήσουμε;

Ψάχνουμε για *αντιπροσωπευτικά* συχνά στοιχειοσύνολα (δηλαδή, να μπορούμε να πάρουμε από αυτά ακριβώς όλα τα συχνά και ιδεατά να μπορούμε να υπολογίσουμε και την υποστήριξη όλων των συχνών):

- Μαximal συχνά
- Κλειστά συχνά

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 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 υποσύνολά τους

ΟΜΩΣ: Δεν προσφέρουν καμιά πληροφορία για την υποστήριξη των υποσυνόλων τους

Συχνά

Μη συχνά

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 15

### Αναπαράσταση Στοιχειοσυνόλων

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

Ένα στοιχειοσύνολο είναι *κλειστό συχνό στοιχειοσύνολο* αν είναι κλειστό και συχνό (δηλαδή, η υποστήριξη του είναι μεγαλύτερη ή ίση με 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

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 16

### Αναπαράσταση Στοιχειοσυνόλων

Maximal vs Closed Itemsets

TID	στοιχεία
1	ABC
2	ABCD
3	BCE
4	ACDE
5	DE

Δεν εμφανίζονται σε καμιά δοσολημία

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 17

### Αναπαράσταση Στοιχειοσυνόλων

Maximal vs Closed Itemsets

Ελάχιστη υποστήριξη = 2

Κλειστά αλλά όχι maximal

Κλειστά και maximal

# Closed = 9  
# Maximal = 4

Για να υπολογίσουμε όλα τα συχνά στοιχειοσύνολα και την υποστήριξη τους, ξεκινάμε από τα μεγαλύτερα κλειστά και προχωράμε

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 18

## Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

## Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

Ο Αργιότερος από τους παλιότερους, αλλά:

Συχνά μεγάλο I/O επειδή κάνει πολλαπλά περάσματα στη βάση των δοσοληψιών

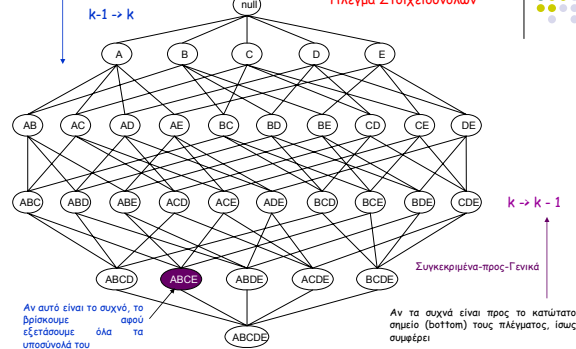
Κακή απόδοση όταν οι δοσοληψίες έχουν μεγάλο πλάτος

Άλλες μέθοδοι:

- Διαφορετικές διασχίσεις του πλέγματος των στοιχειοσυνόλων
- Αναπαράσταση Συνόλου Δοσοληψιών

## Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

Αργιότερο: Γενικά-προς-Συγκεκριμένα



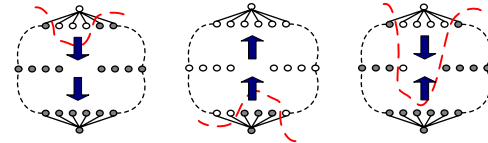
## Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

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

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

Πιο 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 υποσύνολά του



Γενικό-προς-Συγκεκριμένο Συγκεκριμένο-προς-Γενικό Διπλής Κατεύθυνσης

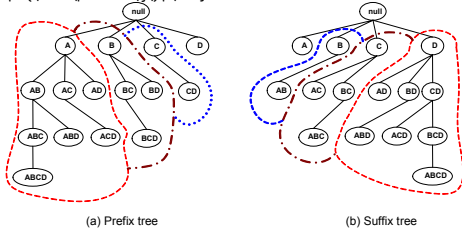
## Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

Διάσχιση του Πλέγματος των Στοιχειοσυνόλων:  
Κλάσεις Ισοδυναμίας

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

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

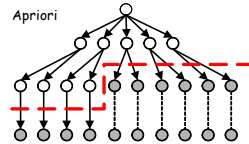
Prefix (Suffix): Δύο στοιχειοσύνολα ανήκουν στην ίδια κλάση αν έχουν κοινό πρόθεμα (ή επίθημα-κατάληξη) μήκους  $k$



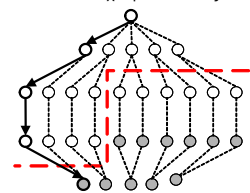
## Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

Διάσχιση του Πλέγματος των Στοιχειοσυνόλων:  
BFS vs DFS

Αργιότερο



DFS: Depth-First-Search  
Διάσχιση κατά Βάθος



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

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

### Άλλοι Μέθοδοι Υπολογισμού Συχνών Στοιχειοσυνόλων

Διάσχιση του Πλέγματος των Στοιχειοσυνόλων: BFS vs DFS

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

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 25

### Άλλοι Τρόποι Υπολογισμού

Αναπαράσταση της Βάσης Δεδομένων: Οριζόντια vs Κάθετη

Αυτό χρησιμοποιεί ο αργιότι

Εναλλακτικά: Για κάθε στοιχείο σε ποιες δσοοληψίες εμφανίζεται

Κάθετη Διάβρωση Δεδομένων

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

Οριζόντια Διάβρωση Δεδομένων

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

Η υποστήριξη υπολογίζεται παίρνοντας τις τομές των TID-λιστών

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 26

### Άλλοι Τρόποι Υπολογισμού

Η υποστήριξη υπολογίζεται παίρνοντας τις τομές των TID-λιστών

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-λιστών δύο από τα (k-1)-υπο-στοιχειοσύνολα του.
- Πλεονέκτημα: πολύ γρήγορος υπολογισμός της υποστήριξης
- Πρόβλημα, αν οι TID-λίστες είναι μεγάλες και δε χωρούν στη μνήμη

Θα δούμε τον FP-Growth που χρησιμοποιεί μια prefix-based αναπαράσταση των δσοοληψιών

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 27

### Ο Αλγόριθμος FP-Growth

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 28

### Αλγόριθμος FP-Growth

#### Με λίγα λόγια:

Ο αλγόριθμος χρησιμοποιεί μια συμπίεσμένη αναπαράσταση της βάσης με τη μορφή ενός FP-δέντρου

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

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

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 29

### 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 βρούμε 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

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΧΕΤΗΣΗΣ II 30

### 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 έχει μια **ΕΤΙΚΕΤΑ**: ποιο στοιχείο και τη συχνότητα εμφάνισης (υποστήριξη) - πόσες δοσοληψίες φτάνουν σε αυτόν

Εύρωτη Διδασκνν: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 31

### 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

Εύρωτη Διδασκνν: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 32

### 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:

Επίσης, κρατάμε **πίνακα δεικτών** για να βοηθήσουν στον υπολογισμό των συχνών στοιχειασυνόλων

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

Εύρωτη Διδασκνν: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 33

### 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:

Επίσης, κρατάμε **πίνακα δεικτών** για να βοηθήσουν στον υπολογισμό των συχνών στοιχειασυνόλων

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

Εύρωτη Διδασκνν: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 34

### 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:

Επίσης, κρατάμε **πίνακα δεικτών** για να βοηθήσουν στον υπολογισμό των συχνών στοιχειασυνόλων

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

Εύρωτη Διδασκνν: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 35

### 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:

Επίσης, κρατάμε **πίνακα δεικτών** για να βοηθήσουν στον υπολογισμό των συχνών στοιχειασυνόλων

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

Εύρωτη Διδασκνν: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 36

### 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}

Πίνακας Δεικτών

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

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 37

### 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)

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 38

### Αλγόριθμος FP-Growth

#### Κατασκευή FP-δέντρου

Το τελικό δέντρο, εξαρτάται από τη **διάταξη**: άλλη διάταξη → άλλα προθέματα

(Συνήθως) μικρότερο δέντρο, αν όχι λεξιλογιακά, 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

▪ **Επίσης, αγνοούμε όσα στοιχεία είναι μη συχνά**

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	{B,A,C}
9	{B,A,D}
10	{B,C,E}

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 39

### Αλγόριθμος FP-Growth

#### Αλγόριθμος εύρεσης συχνών στοιχειοσυνόλων

- Είσοδος: FP-δέντρο
- Έξοδος: Συχνά στοιχειοσύνολα και η υποστήριξη τους
- Μέθοδος
  - Διαίρει-και-Βασίλευε
    - ο Χωρίζουμε τα στοιχειοσύνολα σε αυτά που τελειώνουν σε E, D, C, B, A
    - ο Μετά αυτά που τελειώνουν σε E σε αυτά σε DE, CE, BE, AE κοκ

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 40

### Αλγόριθμος FP-Growth

#### Αλγόριθμος εύρεσης συχνών στοιχειοσυνόλων

Όλα τα στοιχειοσύνολα

Όλα τα δυνατά στοιχειοσύνολα!

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 41

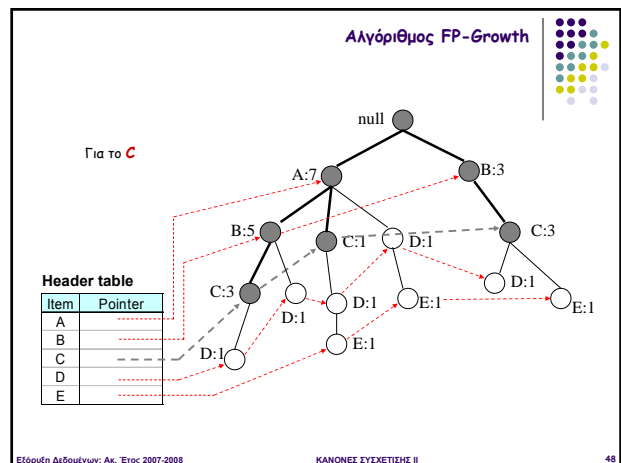
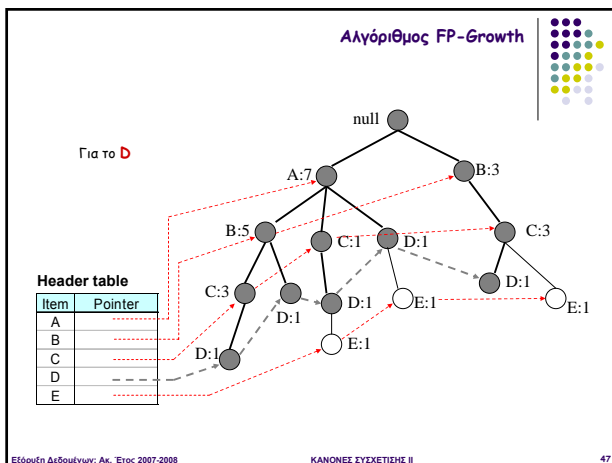
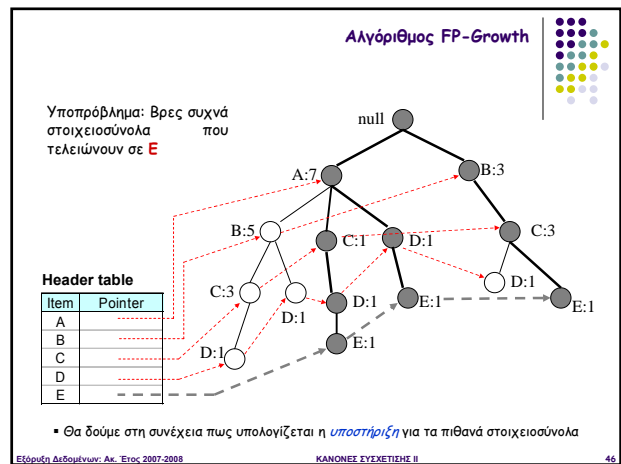
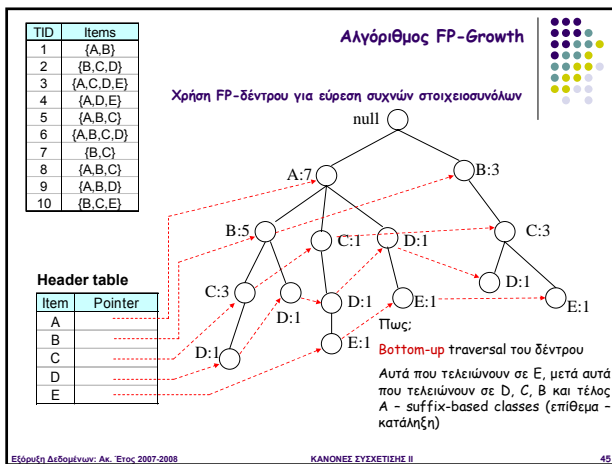
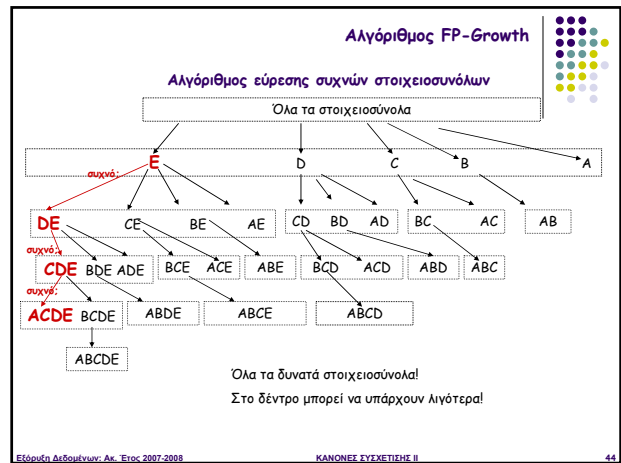
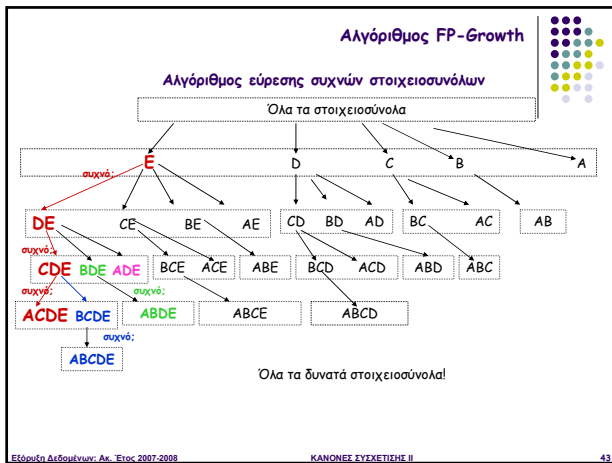
### Αλγόριθμος FP-Growth

#### Αλγόριθμος εύρεσης συχνών στοιχειοσυνόλων

Όλα τα στοιχειοσύνολα

Όλα τα δυνατά στοιχειοσύνολα!

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 42





**Αλγόριθμος FP-Growth**

Για το **B**

Item	Pointer
A	---
B	---
C	---
D	---
E	---

Εξώφυλλο Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 49

**Αλγόριθμος FP-Growth**

Για το **A**

Item	Pointer
A	---
B	---
C	---
D	---
E	---

Εξώφυλλο Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 50

**Αλγόριθμος FP-Growth**

**Συνοπτικά**

Σε κάθε βήμα, για το suffix X

- Φάση 1
  - κατασκευάζουμε το **προθεματικό δέντρο** για το X και υπολογίζουμε την υποστήριξη χρησιμοποιώντας τον πίνακα
- Φάση 2
  - Αν είναι συχνό, κατασκευάζουμε το **υπο-συνθήκη δέντρο** για το X, σε βήματα
    - επανα-υπολογισμός υποστήριξης
    - περικοπή κόμβων με μικρή υποστήριξη
    - περικοπή φύλλων

Εξώφυλλο Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 51

**Αλγόριθμος FP-Growth**

**Φάση 1 - κατασκευή προθεματικού δέντρου**

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 Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 52

**Αλγόριθμος FP-Growth**

**Φάση 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 Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 53

**Αλγόριθμος FP-Growth**

Έστω  $\text{minsup} = 2$

Βρες την υποστήριξη του {E}

Πως:  
Ακολουθήσε τους συνδέσμους αθροίζοντας  $1+1=3 > 2$   
Οπότε {E} συχνό

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

Εξώφυλλο Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 54

**Αλγόριθμος FP-Growth**

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

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

Εύρωδη Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 55

**Αλγόριθμος FP-Growth**

Αλλαγή μετρητών  
 Οι μετρητές σε κάποιους κόμβους περιλαμβάνουν ποσοτήτες που δεν έχουν το E  
 Πχ στο null->B->C->E μετράμε και την {B, C}

Εύρωδη Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 56

**Αλγόριθμος FP-Growth**

Εύρωδη Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 57

**Αλγόριθμος FP-Growth**

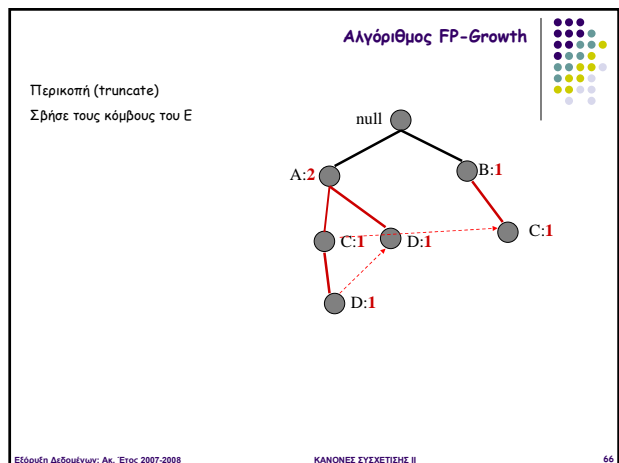
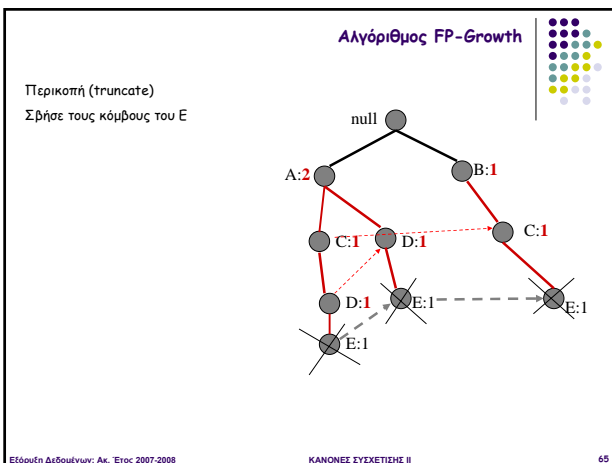
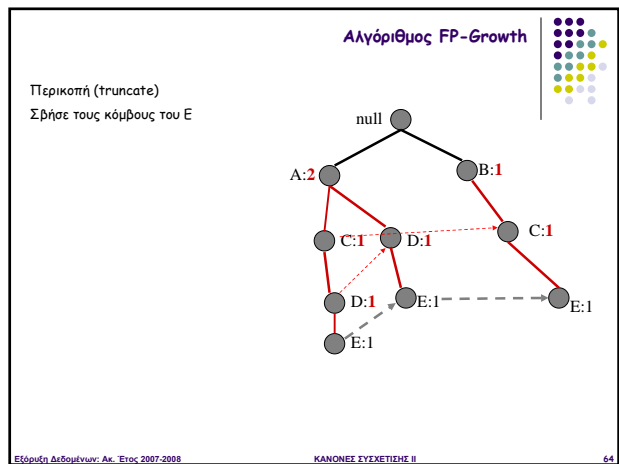
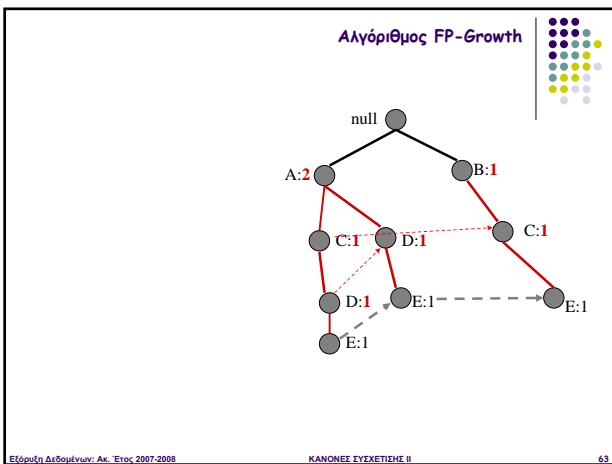
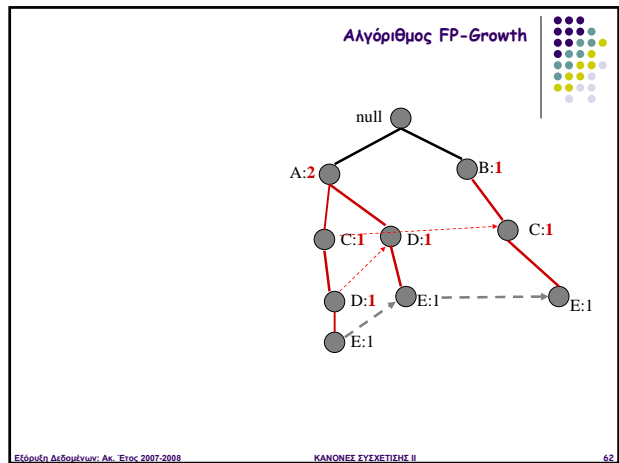
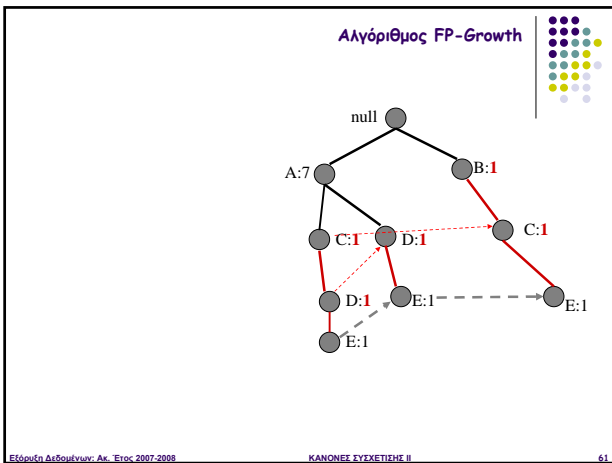
Εύρωδη Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 58

**Αλγόριθμος FP-Growth**

Εύρωδη Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 59

**Αλγόριθμος FP-Growth**

Εύρωδη Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 60



**Αλγόριθμος FP-Growth**

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

Αυτό σημαίνει ότι το B εμφανίζεται μαζί με το E λιγότερο από  $\text{minsup}$  φορές

Εξώφυλλο Διδασκίντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 67

**Αλγόριθμος FP-Growth**

Εξώφυλλο Διδασκίντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 68

**Αλγόριθμος FP-Growth**

Εξώφυλλο Διδασκίντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 69

**Αλγόριθμος FP-Growth**

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

Εξώφυλλο Διδασκίντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 70

**Αλγόριθμος FP-Growth**

**Ψάξη 1**  
 Όλα τα μονοπάτια που περιέχουν το D (DE)  
 Πρόθεματικά Μονοπάτια (prefix paths)

Εξώφυλλο Διδασκίντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 71

**Αλγόριθμος FP-Growth**

**Ψάξη 1**  
 Όλα τα μονοπάτια που περιέχουν το D (DE)  
 Πρόθεματικά Μονοπάτια (prefix paths)

Εξώφυλλο Διδασκίντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 72

**Αλγόριθμος FP-Growth**

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

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 73

**Αλγόριθμος FP-Growth**

**Ψάξη 2**  
 Κατασκευάσε το υπο-συνθήκη FP-  
 δέντρο για το {D, E}

1. Αλλαγή υποστήριξης
2. Περικοπές κόμβων

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 74

**Αλγόριθμος FP-Growth**

1. Αλλαγή υποστήριξης

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 75

**Αλγόριθμος FP-Growth**

2. Περικοπές κόμβων

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 76

**Αλγόριθμος FP-Growth**

2. Περικοπές κόμβων

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 77

**Αλγόριθμος FP-Growth**

2. Περικοπές κόμβων

Εξώφυλλο Διδασκάντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 78

**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

Εξώφυλλο Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 79

**Αλγόριθμος FP-Growth**

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

Εξώφυλλο Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 80

**Αλγόριθμος FP-Growth**

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

Εξώφυλλο Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 81

**Αλγόριθμος FP-Growth**

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

Εξώφυλλο Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 82

**Αλγόριθμος FP-Growth**

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

Εξώφυλλο Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 83

**Αλγόριθμος FP-Growth**

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

1. Αλλαγή υποστήριξης
2. Περικοπές κόμβων

Εξώφυλλο Διδασκόντων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 84

**Αλγόριθμος FP-Growth**

1. Αλλαγή υποστήριξης

Εξώφυλλο Διδασκάλων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 85

**Αλγόριθμος FP-Growth**

2. Περικοπή Κόμβων

Εξώφυλλο Διδασκάλων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 86

**Αλγόριθμος FP-Growth**

2. Περικοπή Κόμβων

Εξώφυλλο Διδασκάλων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 87

**Αλγόριθμος FP-Growth**

2. Περικοπή Κόμβων

Εξώφυλλο Διδασκάλων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 88

**Αλγόριθμος FP-Growth**

2. Περικοπή Κόμβων

null ●

Άρα, επιστροφή στο επόμενο υποπρόβλημα

Εξώφυλλο Διδασκάλων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 89

**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 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 90

**Αλγόριθμος FP-Growth**

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

Εξώφυλλο Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 91

**Αλγόριθμος FP-Growth**

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

Εξώφυλλο Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 92

**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}

Εξώφυλλο Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 93

**Αλγόριθμος FP-Growth**

Άρα για το E  
 Έχουμε τα εξής συχνά στοιχειοσύνολα  
 {E} {D, E} {A, D, E} {C, E} {A, E}

Συνεχίζουμε για το D

Εξώφυλλο Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 94

**Αλγόριθμος FP-Growth**

Για το D

**Header table**

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

Εξώφυλλο Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 95

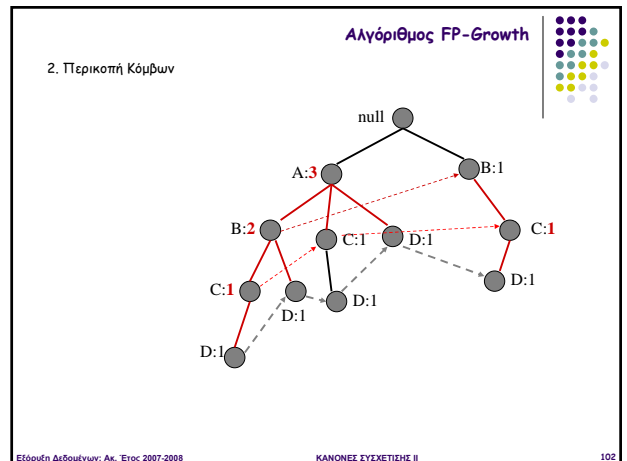
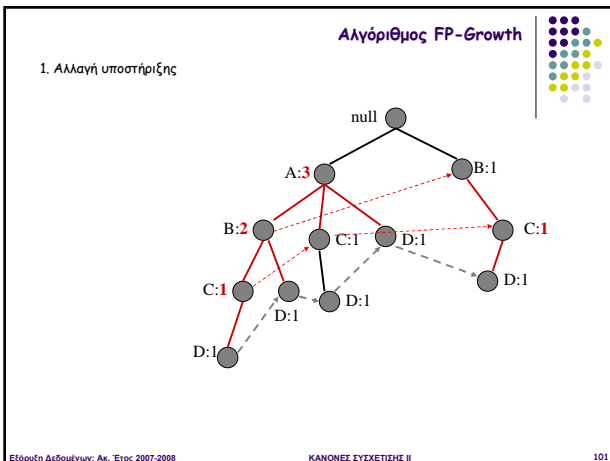
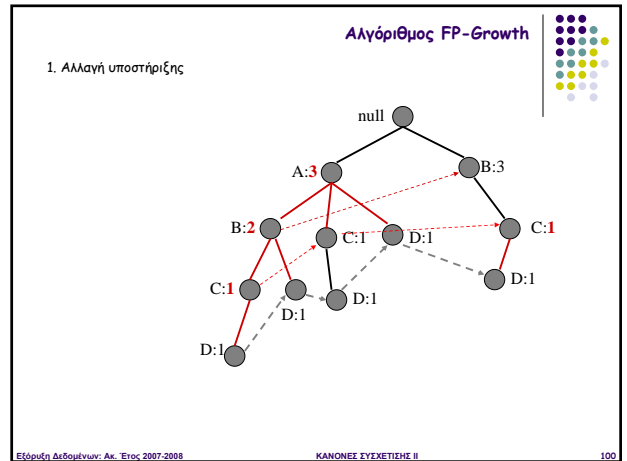
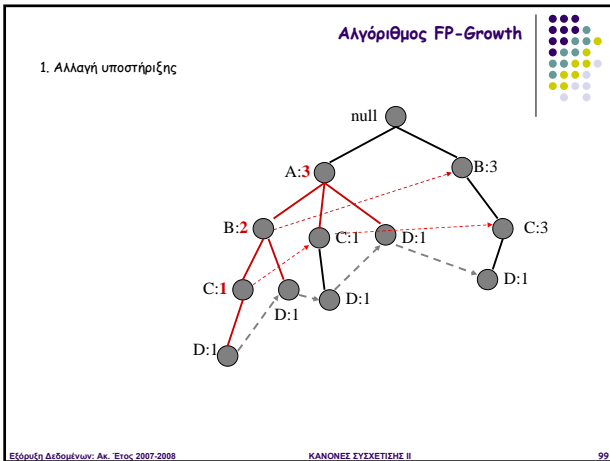
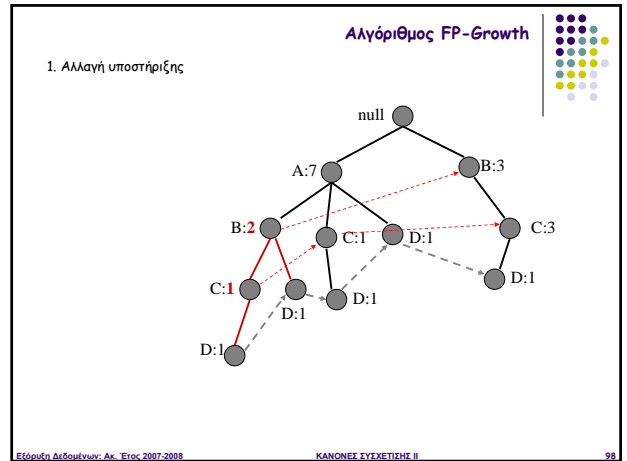
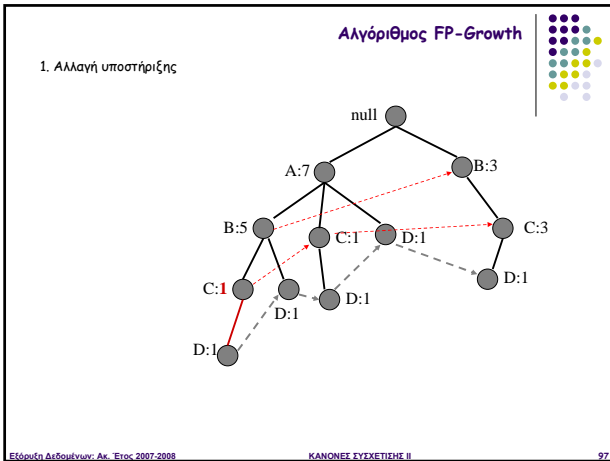
**Αλγόριθμος FP-Growth**

**Ψάξη 1**  
 Όλα τα προθεματικά μονοπάτια που περιέχουν το D  
 Υποστήριξη 5x2 -> άρα συχνό

Μετατροπή του προθεματικού δέντρου σε FP-δέντρο υπό συνθήκη

Εξώφυλλο Διδασκάλιων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 96





2. Περικοπή Κόμβων

Αλγόριθμος FP-Growth

Εύρωδη Διδασκάλου: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 103

Προθεματικά δέντρα και υποσυνθήκη δέντρα

Για τα AB, BD και CD κοκ

Αλγόριθμος FP-Growth

Εύρωδη Διδασκάλου: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 104

Αλγόριθμος FP-Growth

Παρατηρήσεις

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 υποπροβλήματα 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

Εύρωδη Διδασκάλου: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 105

Αλγόριθμος FP-Growth

Παρατηρήσεις

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)

Εύρωδη Διδασκάλου: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 106

Αποτίμηση Κανόνων Συσχέτισης

Εύρωδη Διδασκάλου: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 107

Αποτίμηση Κανόνων Συσχέτισης

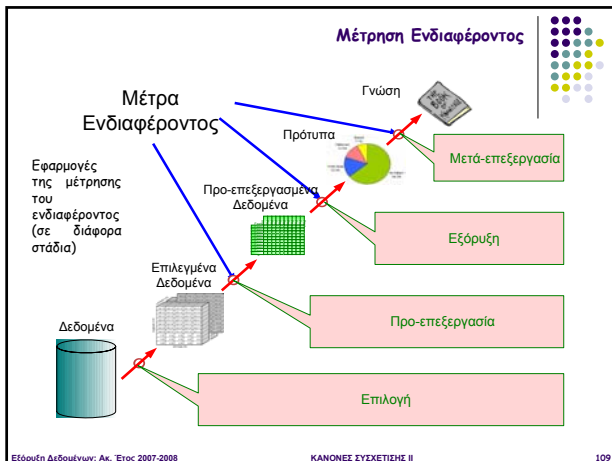
Παράγουν πάρα πολλούς κανόνες που συχνά είναι μη ενδιαφέροντες ή πλεονάζοντες (περιττοι)

Πλεονάζοντες αν  $\{A, B, C\} \rightarrow \{D\}$  και  $\{A, B\} \rightarrow \{D\}$  έχουν την ίδια υποστήριξη & εμπιστοσύνη

Μέτρα ενδιαφέροντος (interestingness) χρησιμοποιούνται για να ελαττώσουν (prune) ή να ιεραρχήσουν (rank) τα παραγόμενα πρότυπα

Χρησιμοποιούνται σε διάφορα στάδια της διαδικασίας ανάκτησης γνώσης

Εύρωδη Διδασκάλου: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 108



### Αποτίμηση Κανόνων Συσχέτισης

Γενικά: **αντικειμενικά** (objective) και **υποκειμενικά** (subjective) μέτρα ενδιαφέροντος

Ας δούμε πρώτα μερικά αντικειμενικά κριτήρια:

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

Γενικά συνήθως βασίζονται σε μετρήσεις της συχνότητας εμφάνισης που δίνονται μέσω ενός πίνακα "**contingency**" (συνάφειας)

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 110

### Μέτρηση Ενδιαφέροντος

Υπολογισμός του Μέτρου Ενδιαφέροντος (αντικειμενικά μέτρα)

**Contingency table (πίνακας συνάφειας)** Μέτρηση συχνότητας εμφάνισης

	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  $\bar{Y}$

$f_{11}$  πόσο συχνά εμφανίζεται το X και το Y (support count)  
 $f_{+1}$  μετρητής υποστήριξης (support count) του Y

Χρησιμοποιείται για τον ορισμό διαφόρων μέτρων

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

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 111

### Μέτρηση Ενδιαφέροντος

**Μειονεκτήματα της Εμπιστοσύνης**

Μεγάλες τιμές υποστήριξης μπορεί να «διώξουν» ενδιαφέροντες κανόνες. Τι γίνεται με την εμπιστοσύνη;

	Coffee	$\bar{C}$	
Tea	15	5	20
$\bar{T}$	75	5	80
	90	10	100

Ενδιαφερόμαστε για τη σχέση μεταξύ αυτών που πίνουν καφέ και αυτών που πίνουν τσάι  
Κανόνας Συσχέτισης: Tea  $\rightarrow$  Coffee

Εμπιστοσύνη =  $P(\text{Coffee}|\text{Tea}) = 0.75$

Ενώ ο κανόνας έχει υψηλή εμπιστοσύνη, ο κανόνας είναι παραπλανητικός  
 $P(\text{Coffee}|\bar{T}) = 0.9375$

$P(\text{Coffee}) = 0.9$   
*Αγνοεί την υποστήριξη του RHS*

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 112

### Μέτρηση Ενδιαφέροντος

Εξαιτίας τέτοιων προβλημάτων της υποστήριξης/εμπιστοσύνης,

Έχουν προταθεί **πολλά** αντικειμενικά μέτρα για τη μέτρηση του ενδιαφέροντος των κανόνων, που στηρίζονται κυρίως στην έννοια της στατιστικής ανεξαρτησίας

Ας δούμε ένα παράδειγμα

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 113

### Μέτρα βασισμένα στη Στατιστική

#### Στατιστική Ανεξαρτησία

Πληθυσμός 1000 σπουδαστών

- 600 σπουδαστές ξέρουν κολύμπι (S)
- 700 σπουδαστές ξέρουν ποδήλατο (B)
- 420 σπουδαστές ξέρουν κολύμπι και ποδήλατο (S,B)

- $P(S \cap B) = 420/1000 = 0.42$
- $P(S) \times P(B) = 0.6 \times 0.7 = 0.42$

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

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008 ΚΑΝΟΝΕΣ ΣΥΣΧΕΤΙΣΗΣ II 114

### Μέτρα βασισμένα στη Στατιστική

Μέτρα που λαμβάνουν υπ' όψιν τους τη στατιστική εξάρτηση

Για τη συσχέτιση:  $X \rightarrow Y$

$$Lift = \frac{P(Y|X)}{P(Y)} = \frac{f_{11}}{f_{+1}}$$

$$Interest = \frac{P(X,Y)}{P(X)P(Y)} = \frac{Nf_{11}}{f_{1+}f_{+1}}$$

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

$$\phi\text{-coefficient} = \frac{P(X,Y) - P(X)P(Y)}{\sqrt{P(X)[1-P(X)]P(Y)[1-P(Y)]}} = \frac{f_{11}f_{00} - f_{01}f_{10}}{\sqrt{f_{1+}f_{+1}f_{0+}f_{+0}}}$$

### Μέτρα βασισμένα στη Στατιστική

Παράδειγμα: Lift/Interest

	Coffee	Coffee	
Tea	15	5	20
Tea	75	5	80
	90	10	100

Κανόνας συσχέτιση: Tea  $\rightarrow$  Coffee

Εμπιστοσύνη=  $P(\text{Coffee}|\text{Tea}) = 0.75$

αλλά  $P(\text{Coffee}) = 0.9$

$\Rightarrow$  Interest =  $0.15/(0.9*0.2) = 0.8333$  ( $< 1$ , άρα αρνητικά συσχετιζόμενα)

### Μέτρα βασισμένα στη Στατιστική

Μειοεκτιμήματα του Lift & Interest

	Y	$\bar{Y}$	
X	10	0	10
$\bar{X}$	0	90	90
	10	90	100

	Y	$\bar{Y}$	
X	90	0	90
$\bar{X}$	0	10	10
	90	10	100

$$I = \frac{0.1}{(0.1)(0.1)} = 10$$

$$I = \frac{0.9}{(0.9)(0.9)} = 1.11$$

Μεγαλύτερο αν και σπάνια εμφανίζονται μαζί!

$$c = 10/100 = 0.1$$

$$s = 1$$

$$c = 90/100 = 0.9$$

$$s = 1$$

### Μέτρα βασισμένα στη Στατιστική

$\phi$ -Coefficient

$$\phi\text{-coefficient} = \frac{P(X,Y) - P(X)P(Y)}{\sqrt{P(X)[1-P(X)]P(Y)[1-P(Y)]}} = \frac{f_{11}f_{00} - f_{01}f_{10}}{\sqrt{f_{1+}f_{+1}f_{0+}f_{+0}}}$$

Κανονικοποιημένη τιμή μεταξύ του -1 και 1

Διαδική εκδοχή του Pearson's coefficient

0: στατιστική ανεξαρτησία

-1: τέλεια αρνητική συσχέτιση

1: τέλεια θετική συσχέτιση

### Μέτρα βασισμένα στη Στατιστική

$\phi$ -Coefficient

	Y	$\bar{Y}$	
X	60	10	70
$\bar{X}$	10	20	30
	70	30	100

	Y	$\bar{Y}$	
X	20	10	30
$\bar{X}$	10	60	70
	30	70	100

$$\phi = \frac{0.6 - 0.7 \times 0.7}{\sqrt{0.7 \times 0.3 \times 0.7 \times 0.3}}$$

$$= 0.5238$$

$$\phi = \frac{0.2 - 0.3 \times 0.3}{\sqrt{0.7 \times 0.3 \times 0.7 \times 0.3}}$$

$$= 0.5238$$

$\phi$  Coefficient ίδιος και για τους δύο πίνακες

### Μέτρα βασισμένα στη Στατιστική

$\phi$ -Coefficient

$$\phi\text{-coefficient} = \frac{P(X,Y) - P(X)P(Y)}{\sqrt{P(X)[1-P(X)]P(Y)[1-P(Y)]}} = \frac{f_{11}f_{00} - f_{01}f_{10}}{\sqrt{f_{1+}f_{+1}f_{0+}f_{+0}}}$$

▪ Είναι κατάλληλο για μη συμμετρικές (η απουσία και η παρουσία μετρούν το ίδιο)

▪ Λόγω κανονικοποίησης, αγνοεί το μέγεθος του δείγματος

### Μέτρα βασισμένα στη Στατιστική

#### IS-measure

$$IS(X, Y) = \frac{s(X, Y)}{\sqrt{s(X)s(Y)}} = \frac{f_{11}}{\sqrt{f_{1+}f_{+1}}} = \sqrt{I(X, Y)s(x, Y)}$$

- είναι το συνημίτονο αν θεωρηθούν διαδικές μεταβλητές
- γεωμετρικός μέσος της εμπιστοσύνης του  $X \rightarrow Y$  και  $Y \rightarrow X$

Στη βιβλιογραφία έχουν προταθεί πολλά μέτρα ανάλογα με την εφαρμογή

Με ποια κριτήρια θα επιλέξουμε ένα καλό μέτρο;

Πως έναν Αpriori-style support based pruning επηρεάζει αυτά τα μέτρα;

#	Measure	Formula
1	φ-coefficient	$\frac{P(A, B) - P(A)P(B)}{\sqrt{(P(A)P(B) - P(A)^2)(P(B) - P(B)^2)}}$
2	Goodman-Kruskal's (λ)	$\frac{\sum_{i,j} \max(P(A_i, B_j) - \sum_{k \neq i} P(A_i, B_k) - \max_j P(A_j) - \max_k P(B_k))}{3 - \max_j P(A_j) - \max_k P(B_k)}$
3	Odds ratio (ω)	$\frac{P(A, B)P(\bar{A}, \bar{B})}{P(A, \bar{B})P(\bar{A}, B)}$
4	Yule's Q	$\frac{P(A, B)P(\bar{A}, \bar{B}) - P(A, \bar{B})P(\bar{A}, B)}{P(A, B)P(\bar{A}, \bar{B}) + P(A, \bar{B})P(\bar{A}, B)}$
5	Yule's Y	$\frac{\sqrt{P(A, B)P(\bar{A}, \bar{B}) - P(A, \bar{B})P(\bar{A}, B)}}{\sqrt{P(A, B)P(\bar{A}, \bar{B}) + P(A, \bar{B})P(\bar{A}, B)}}$
6	Kappa (κ)	$\frac{P(A, B) - P(A)P(B)}{1 - P(A)P(B)}$
7	Mutual Information (M)	$-\sum_i P(A_i) \log_2 P(A_i) - \sum_j P(B_j) \log_2 P(B_j) + \sum_{i,j} P(A_i, B_j) \log_2 P(A_i, B_j)$
8	J-Measure (J)	$\max(P(A, B) \log_2 \frac{P(A, B)}{P(A)P(B)} + P(\bar{A}\bar{B}) \log_2 \frac{P(\bar{A}\bar{B})}{P(\bar{A})P(\bar{B})})$
9	Gini index (G)	$\max(P(A)[P(B A)^2 + P(\bar{B} A)^2] + P(\bar{A})[P(B \bar{A})^2 + P(\bar{B} \bar{A})^2] - P(A)^2 - P(\bar{A})^2)$
10	Support (s)	$P(A, B)$
11	Confidence (c)	$\max(P(B A), P(A B))$
12	Laplace (L)	$\max(\frac{P(A, B) + 1}{P(A) + 1}, \frac{P(A, B) + 1}{P(B) + 1})$
13	Conviction (V)	$\max(\frac{P(A, \bar{B})}{P(A)P(\bar{B})}, \frac{P(\bar{A}, B)}{P(\bar{A})P(B)})$
14	Interest (I)	$\frac{P(A, B)}{P(A)P(B)}$
15	cosine (IS)	$\frac{P(A, B)}{\sqrt{P(A)P(B)}}$
16	Piatetsky-Shapiro's (PS)	$P(A, B) - P(A)P(B)$
17	Certainty factor (F)	$\max(\frac{P(A, B) - P(A)P(B)}{1 - P(A)P(B)}, \frac{P(\bar{A}, \bar{B}) - P(\bar{A})P(\bar{B})}{1 - P(\bar{A})P(\bar{B})})$
18	Added Value (AV)	$\max(P(B A) - P(B), P(A B) - P(A))$
19	Collective strength (S)	$\frac{P(A, B) + P(\bar{A}, \bar{B})}{P(A, B) + P(\bar{A}, \bar{B}) + P(A, \bar{B}) + P(\bar{A}, B)}$
20	Jaccard (ζ)	$\frac{P(A, B)}{P(A, B) + P(\bar{A}, \bar{B})}$
21	Klosgem (K)	$\sqrt{P(\bar{A}, \bar{B}) \max(P(B A) - P(B), P(A B) - P(A))}$

### Σύγκριση Μέτρων

### Αποτίμηση Κανόνων Συσχέτισης

10 παραδείγματα contingency πινάκων:

Example	f <sub>11</sub>	f <sub>10</sub>	f <sub>01</sub>	f <sub>00</sub>
E1	8123	83	424	1370
E2	8330	2	622	1046
E3	9481	94	127	298
E4	3954	3080	5	2961
E5	2886	1363	1320	4431
E6	1500	2000	500	6000
E7	4000	2000	1000	3000
E8	4000	2000	2000	2000
E9	1720	7121	5	1154
E10	61	2483	4	7452

Ιεράρχηση των πινάκων με βάση τα διάφορα μέτρα (1 ο πιο ενδιαφέρον, 10 ο λιγότερο ενδιαφέρον):

#	φ	λ	ω	Q	Y	κ	M	J	G	s	c	L	V	I	IS	PS	F	AV	S	ζ	K
E1	1	1	3	3	3	1	2	2	1	3	5	5	4	6	2	2	4	6	1	2	5
E2	2	2	1	1	1	2	1	3	2	2	1	1	1	5	3	5	1	5	2	3	9
E3	3	3	4	4	4	3	3	8	7	1	4	4	6	10	1	8	6	10	3	1	10
E4	4	7	2	2	2	5	4	1	3	6	2	2	4	4	1	2	3	4	5	1	1
E5	5	4	8	8	8	4	7	5	4	7	9	9	9	3	6	3	9	4	5	6	3
E6	6	6	7	7	7	6	4	6	9	8	8	8	7	2	8	6	7	2	7	8	2
E7	7	5	9	9	9	6	8	6	5	4	7	7	8	5	5	4	8	5	6	4	7
E8	8	9	10	10	10	8	10	10	8	4	10	10	10	9	7	7	10	9	8	7	9
E9	9	9	5	5	5	9	7	9	8	3	3	3	3	7	9	9	3	7	9	9	8
E10	10	8	6	6	6	10	5	9	10	6	6	6	5	1	10	10	5	1	10	10	7

### Ιδιότητες ενός Καλού Μέτρου

#### Piatetsky-Shapiro:

3 γενικές ιδιότητες που πρέπει να ικανοποιεί ένα καλό μέτρο M:

- $M(A, B) = 0$  αν τα A και B είναι στατιστικά ανεξάρτητα
- $M(A, B)$  αυξάνει μονότονα με το  $P(A, B)$  όταν τα  $P(A)$  και  $P(B)$  παραμένουν αμετάβλητα
- $M(A, B)$  μειώνεται μονότονα με το  $P(A)$  [ή το  $P(B)$ ] όταν τα  $P(A, B)$  και  $P(B)$  [ή  $P(A)$ ] παραμένουν αμετάβλητα

### Ιδιότητες Μέτρων Αποτίμησης

#### Αλλαγή Διάταξης Μεταβλητών (variable permutation)

	B	$\bar{B}$	
A	p	q	
$\bar{A}$	r	s	

 $\Rightarrow$ 

	A	$\bar{A}$	
B	p	r	
$\bar{B}$	q	s	

Ισχύει  $M(A, B) = M(B, A)$

Γενικά συμμετρικά μέτρα για στοιχειοσύνολα και μη συμμετρικά για κανόνες

Συμμετρικά (symmetric) μέτρα:

- support (υποστήριξη), lift, collective strength, cosine, Jaccard, κλπ

Μη συμμετρικά (asymmetric) μέτρα:

- confidence (εμπιστοσύνη), conviction, Laplace, J-measure, κλπ

### Ιδιότητες Μέτρων Αποτίμησης

#### Κλιμάκωση Γραμμής/Στήλης (Row/Column Scaling)

Παράδειγμα Βαθμός-Φύλο (Mosteller, 1968):

		K <sub>3</sub>	K <sub>4</sub>	
		Male	Female	
K <sub>1</sub>	High	2	3	5
K <sub>2</sub>	Low	1	4	5
		3	7	10

 $\Rightarrow$ 

		Male	Female	
K <sub>1</sub>	High	4	30	34
K <sub>2</sub>	Low	2	40	42
		6	70	76

Mosteller:

Η συσχέτιση πρέπει να είναι ανεξάρτητη από το σχετικό αριθμό αγοριών-κοριτσιών στο δείγμα

Invariant under the row/column scaling operation αν  $M(T) = M(T)$  όπου T ο πίνακας contingency με μετρούς συχνότητας  $[f_{11}, f_{10}, f_{01}, f_{00}]$  και T' ο πίνακας contingency με μετρούς συχνότητας  $[k_1 k_3 f'_{11}, k_2 k_3 f'_{10}, k_1 k_4 f'_{01}, k_2 k_4 f'_{00}]$  όπου  $k_1, k_2, k_3, k_4$  θετικές σταθερές

### Ιδιότητες Μέτρων Αποτίμησης

#### Αντιστροφή (Inversion Operation)

Δοσοληψία 1 →

1	0
0	0
0	0
0	0
0	0
0	0
0	0
0	0
0	0
0	0

Δοσοληψία N →

0	1
1	1
1	1
1	0
1	1
1	1
1	1
1	1
1	1
1	1

0	0
1	0
1	0
1	0
1	0
1	0
1	0
1	0
1	0
1	0

(a)                      (b)                      (c)

**Invariant under the inversion operation** αν η τιμή της παραμένει η ίδια αν ανταλλάξουμε τις τιμές  $f_{11}$  και  $f_{00}$  και τις τιμές  $f_{10}$  και  $f_{01}$

Χρήσιμο για συμμετρικές μεταβλητές

Εξοφλή Διδασκων: Ακ. Έτος 2007-2008                      ΚΑΝΟΝΕΣ ΞΥΧΕΤΙΣΗΣ II                      127

### Ιδιότητες Μέτρων Αποτίμησης

#### Null Addition (προσθήκη μη σχετιζόμενων στοιχείων)

	B	$\bar{B}$
A	p	q
$\bar{A}$	r	s

→

	B	$\bar{B}$
A	p	q
$\bar{A}$	r	s+k

Δεν επηρεάζονται από την αύξηση του  $f_{00}$  όταν οι άλλες τιμές παραμένουν αμετάβλητες

**Invariant measures:**

- support, cosine, Jaccard, κλπ

**Non-invariant measures:**

- correlation, Gini, mutual information, odds ratio, κλπ

Εξοφλή Διδασκων: Ακ. Έτος 2007-2008                      ΚΑΝΟΝΕΣ ΞΥΧΕΤΙΣΗΣ II                      128

### Αποτίμηση Κανόνων Συσχέτισης

#### Παράδοξο του Simpson

**Buy HDTV**

Buy Exercise Machine	
Yes	No
99	81
54	66
153	147

$c((HTVS=Yes) \rightarrow (EM=Yes))=99/180=55\%$   
 $c((HTVS=No) \rightarrow (EM=Yes))=54/120=45\%$

$c((HTVS=Yes) \rightarrow (EM=Yes))=98/170=57.7\%$   
 $c((HTVS=No) \rightarrow (EM=Yes))=50/86=58.1\%$

**Students**

Buy Exercise Machine	
Yes	No
1	9
4	30
5	39

$c((HTVS=Yes) \rightarrow (EM=Yes))=1/10=10\%$   
 $c((HTVS=No) \rightarrow (EM=Yes))=4/34=11.8\%$

**Working adults**

Buy Exercise Machine	
Yes	No
98	72
50	36
148	108

$c((HTVS=Yes) \rightarrow (EM=Yes))=99/180=55\%$   
 $c((HTVS=No) \rightarrow (EM=Yes))=54/120=45\%$

Εξοφλή Διδασκων: Ακ. Έτος 2007-2008                      ΚΑΝΟΝΕΣ ΞΥΧΕΤΙΣΗΣ II                      129

### Αποτίμηση Κανόνων Συσχέτισης

#### Παράδοξο του Simpson

**Buy HDTV**

Buy Exercise Machine	
Yes	No
99 <b>a+p</b>	81 <b>b+q</b>
54 <b>c+r</b>	120 <b>d+s</b>
153	147

$a/b < c/d$   
 $p/q < r/s$  δεν συνεπάγεται ότι  
 $(a+p)/(b+q) < (c+r)/(d+s)$

**Students**

Buy Exercise Machine	
Yes	No
1 <b>a</b>	9 <b>b</b>
4 <b>c</b>	30 <b>d</b>
5	39

$c((HTVS=Yes) \rightarrow (EM=Yes))=1/10=10\%$   
 $c((HTVS=No) \rightarrow (EM=Yes))=4/34=11.8\%$

**Working adults**

Buy Exercise Machine	
Yes	No
98 <b>p</b>	72 <b>q</b>
50 <b>r</b>	36 <b>s</b>
148	108

$c((HTVS=Yes) \rightarrow (EM=Yes))=99/180=55\%$   
 $c((HTVS=No) \rightarrow (EM=Yes))=54/120=45\%$

Είναι σημαντικό πως θα γίνει διαχωρισμός (stratification) των δεδομένων

Εξοφλή Διδασκων: Ακ. Έτος 2007-2008                      ΚΑΝΟΝΕΣ ΞΥΧΕΤΙΣΗΣ II                      130

### Υποκειμενικά Μέτρα Ενδιαφέροντος

- Αντικειμενικά Μέτρα:
  - Ιεραρχούν τα αποτελέσματα με βάση στατιστικά στοιχεία που υπολογίζονται από τα δεδομένα
  - πχ., 21 μετρήσεις συσχέτισης (support, confidence, Laplace, Gini, mutual information, Jaccard, etc).
- Υποκειμενικά Μέτρα:
  - Ιεράρχηση των προτύπων με βάση την ερημνεία του χρήστη
  - Ένα πρότυπο είναι υποκειμενικά ενδιαφέρον αν είναι σε αντίθεση με αυτό που αναμένει ο χρήστης (Silberschatz & Tuzhilin)
  - Ένα πρότυπο είναι υποκειμενικά ενδιαφέρον αν μπορεί να χρησιμοποιηθεί (Silberschatz & Tuzhilin)

Εξοφλή Διδασκων: Ακ. Έτος 2007-2008                      ΚΑΝΟΝΕΣ ΞΥΧΕΤΙΣΗΣ II                      131

### Υποκειμενικά Μέτρα Ενδιαφέροντος

#### Interestingness (ενδιαφέρον) via Unexpectedness (μη αναμονή)

Domain Knowledge

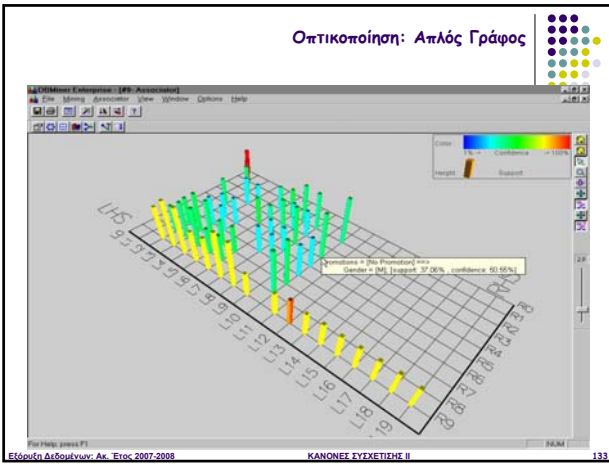
Evidence

- + Pattern expected to be frequent
- Pattern expected to be infrequent
- ⊕ Pattern found to be frequent
- ⊖ Pattern found to be infrequent
- ⊕⊖ Expected Patterns
- ⊖⊕ Unexpected Patterns

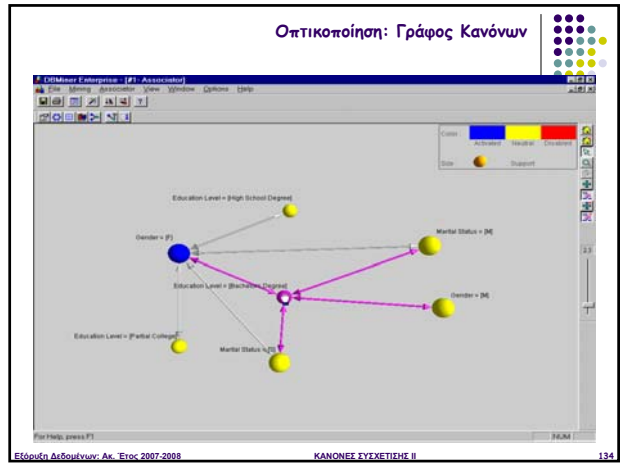
- Χρειάζεται να μοντελοποιήσουμε τι αναμένει ο χρήστης (domain knowledge)
- Χρειάζεται να συνδυάσουμε το τι αναμένεται από τους χρήστες με το τι δίνουν τα δεδομένα (δηλαδή τα πρότυπα που παίρνουμε - evidence)

Εξοφλή Διδασκων: Ακ. Έτος 2007-2008                      ΚΑΝΟΝΕΣ ΞΥΧΕΤΙΣΗΣ II                      132

### Οπτικοποίηση: Απλός Γράφος



### Οπτικοποίηση: Γράφος Κανόνων



### Οπτικοποίηση: (SGI/MineSet 3.0)

