



Αποθήκες Δεδομένων και Εξόρυξη Δεδομένων

Ενότητα 12: Κανόνες Συσχέτισης – Μέρος Β΄

Αναστάσιος Γούναρης, Επίκουρος Καθηγητής
Τμήμα Πληροφορικής ΑΠΘ



Άδειες Χρήσης

- Το παρόν εκπαιδευτικό υλικό υπόκειται σε άδειες χρήσης Creative Commons.
- Για εκπαιδευτικό υλικό, όπως εικόνες, που υπόκειται σε άλλου τύπου άδειας χρήσης, η άδεια χρήσης αναφέρεται ρητώς.



Χρηματοδότηση

- Το παρόν εκπαιδευτικό υλικό έχει αναπτυχθεί στα πλαίσια του εκπαιδευτικού έργου του διδάσκοντα.
- Το έργο «Ανοικτά Ακαδημαϊκά Μαθήματα στο Αριστοτέλειο Πανεπιστήμιο Θεσσαλονίκης» έχει χρηματοδοτήσει μόνο την αναδιαμόρφωση του εκπαιδευτικού υλικού.
- Το έργο υλοποιείται στο πλαίσιο του Επιχειρησιακού Προγράμματος «Εκπαίδευση και Δια Βίου Μάθηση» και συγχρηματοδοτείται από την Ευρωπαϊκή Ένωση (Ευρωπαϊκό Κοινωνικό Ταμείο) και από εθνικούς πόρους.





Κανόνες Συσχέτισης – Μέρος Β΄

Αλγόριθμος FP-Growth



Ευρωπαϊκή Ένωση
Ευρωπαϊκό Κοινωνικό Ταμείο



ΕΠΙΧΕΙΡΗΣΙΑΚΟ ΠΡΟΓΡΑΜΜΑ
ΕΚΠΑΙΔΕΥΣΗ ΚΑΙ ΔΙΑ ΒΙΟΥ ΜΑΘΗΣΗ
επένδυση στην κοινωνία της γνώσης

ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ ΚΑΙ ΘΡΗΣΚΕΥΜΑΤΩΝ
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



ΕΥΡΩΠΑΪΚΟ ΚΟΙΝΩΝΙΚΟ ΤΑΜΕΙΟ

Περιεχόμενα ενότητας

1. Αλγόριθμος FP-Growth.
2. Μειονεκτήματα υποστήριξης-εμπιστοσύνης.
3. Κλειστά – Maximal στοιχειοσύνολα.



Σκοποί ενότητας

- Ανάλυση των κανόνων συσχέτισης.
- Περιγραφή του αλγορίθμου FP-Growth.
- Κατασκευή του FP-δένδρου.
- Παρουσίαση των μειονεκτημάτων σχετικά με την υποστήριξη και εμπιστοσύνη.



Είναι γρήγορος ο Apriori?

Bottlenecks στην απόδοση

- Ο βασικός αλγόριθμος Apriori:
 - Χρησιμοποιεί συχνά $(k - 1)$ -στοιχειοσύνολα για την παραγωγή υποψηφίων συχνών k -στοιχειοσυνόλων.
 - Χρήση τεχνικών σαρώματος ΒΔ και ταύτισης προτύπων για τη μέτρηση της υποστήριξης των υποψηφίων συνόλων.
- Το bottleneck του Apriori: Δημιουργία υποψηφίων.
 - Πολύ μεγάλα υποψήφια σύνολα.
 - Πολλαπλές σαρώσεις της ΒΔ.



Ο Αλγόριθμος FP-Growth

- Χρησιμοποιεί μια συμπιεσμένη αναπαράσταση της βάσης
- με τη μορφή ενός **FP-δένδρου** (*FP: frequent pattern*).
- Το δένδρο μοιάζει με προθεματικό δένδρο - prefix tree (trie).
- Ο αλγόριθμος κατασκευής διαβάζει μια συναλλαγή τη φορά, και απεικονίζει τη συναλλαγή σε ένα μονοπάτι του FP-δένδρου.
- Μερικά μονοπάτια μπορεί να επικαλύπτονται: όσο περισσότερα μονοπάτια επικαλύπτονται, τόσο καλύτερη συμπίεση.
- Τα συχνά στοιχειosύνολα βρίσκονται με μια αναδρομική διαίρει-και-βασίλευε προσέγγιση.



Κατασκευή FP-δένδρου (1/7)

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-δένδρο είναι ένα προθεματικό δένδρο.
- Άρα τα στοιχεία σε κάθε σύνολο πρέπει να ακολουθούν κάποια **διάταξη**, έστω τη λεξικογραφική.
- Θα δούμε αργότερα ότι κάτι άλλο συμφέρει περισσότερο.

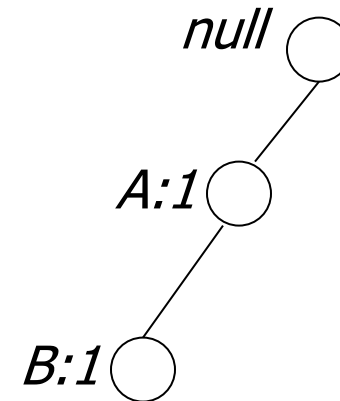
Αρχικά, το δένδρο είναι κενό.



Κατασκευή FP-δένδρου (2/7)

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}

Διάγραμμα TID=1:



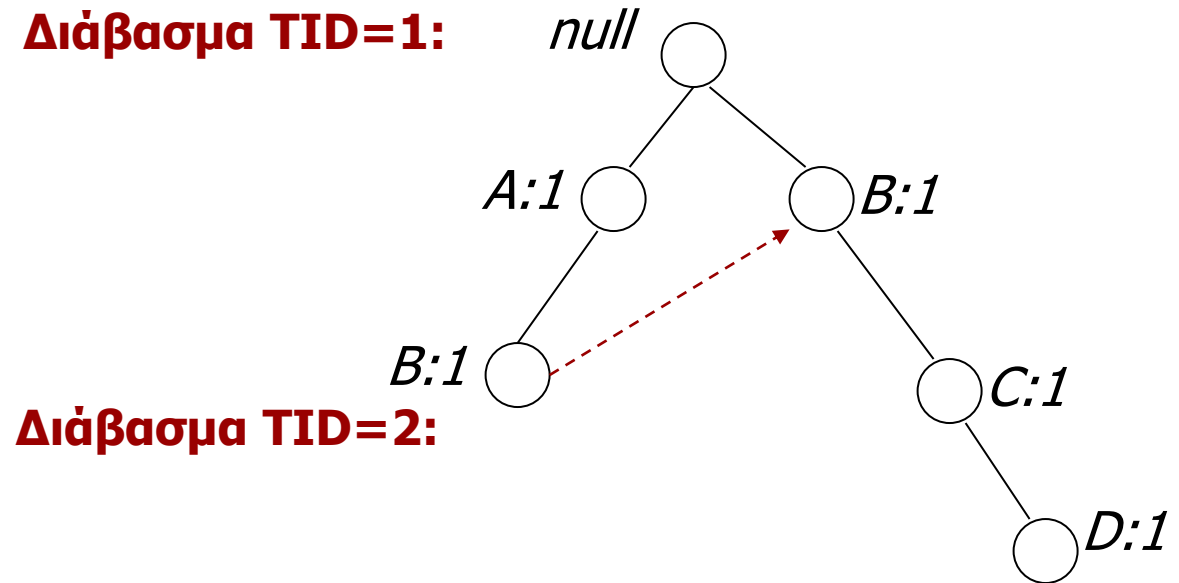
- Κάθε κόμβος έχει μια ετικέτα που δείχνει πόσες συναλλαγές φτάνουν σε αυτόν, δηλαδή πόσα μονοπάτια καταλήγουν σε αυτόν τον κόμβο.



Κατασκευή FP-δένδρου (3/7)

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:



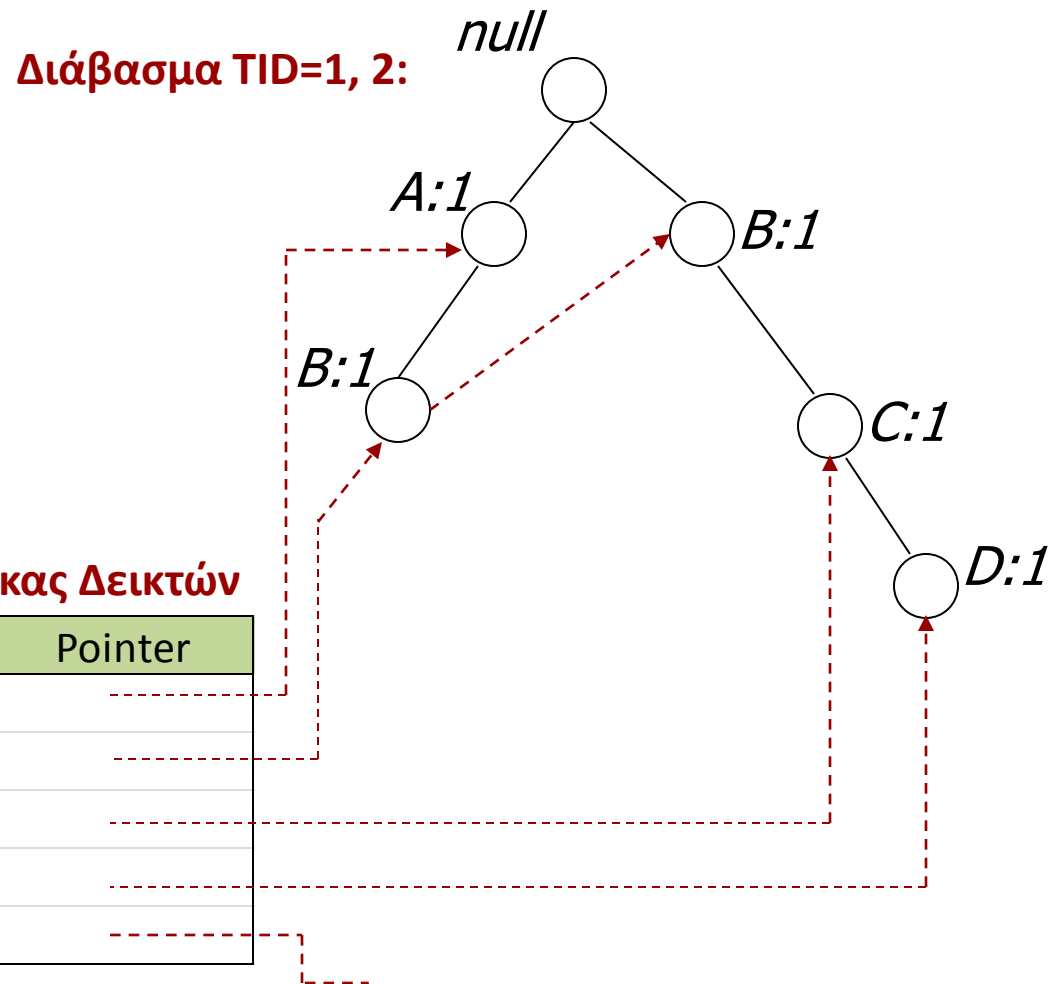
Διάβασμα TID=2:

- Κάθε κόμβος έχει μια ετικέτα που δείχνει πόσες συναλλαγές φτάνουν σε αυτόν.
- Επίσης, υπάρχουν δείκτες μεταξύ των κόμβων που αναφέρονται στο ίδιο στοιχείο.



Κατασκευή FP-δένδρου (4/7)

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}

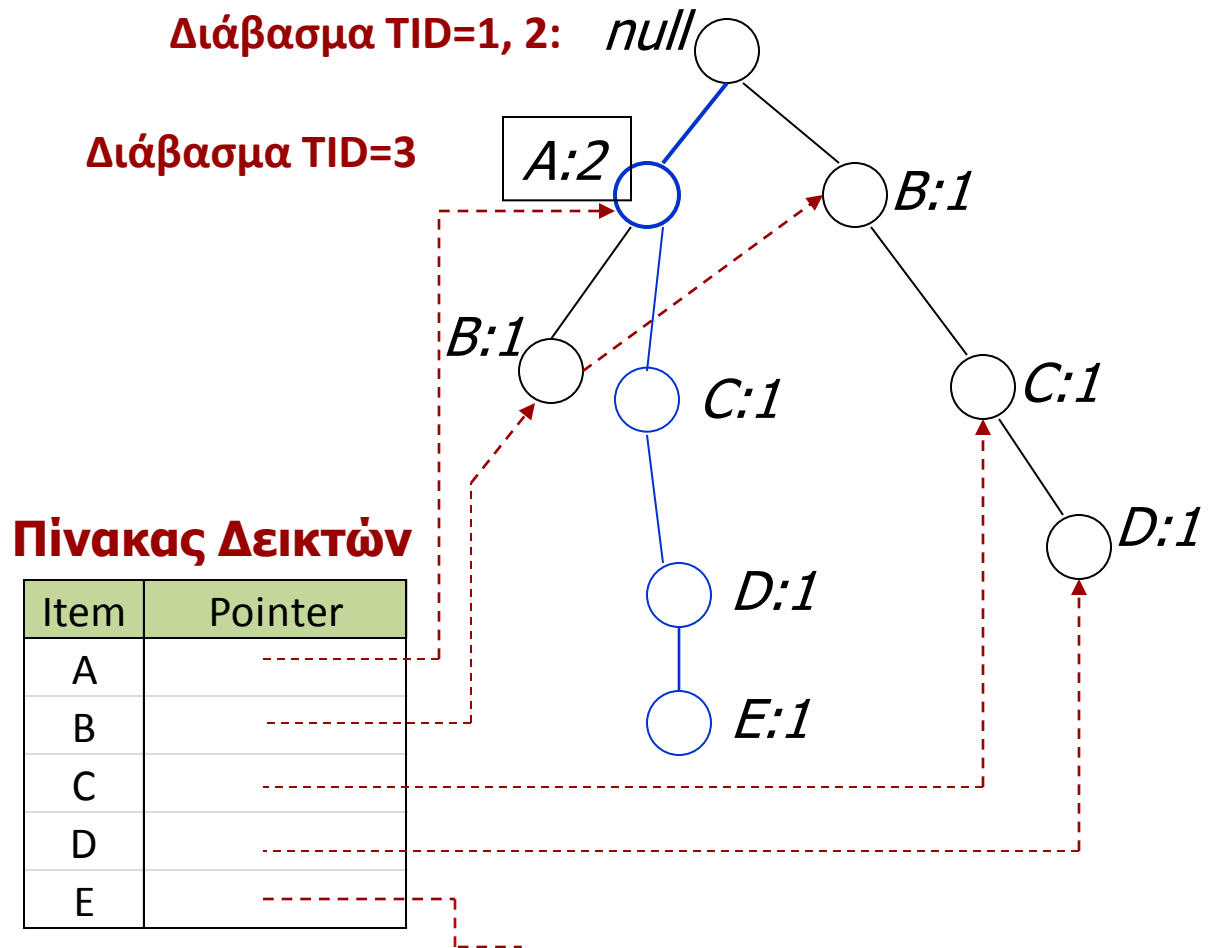


- Κρατάμε πίνακα δεικτών για να βοηθήσουν στον υπολογισμό των συχνών στοιχειοσυνόλων.



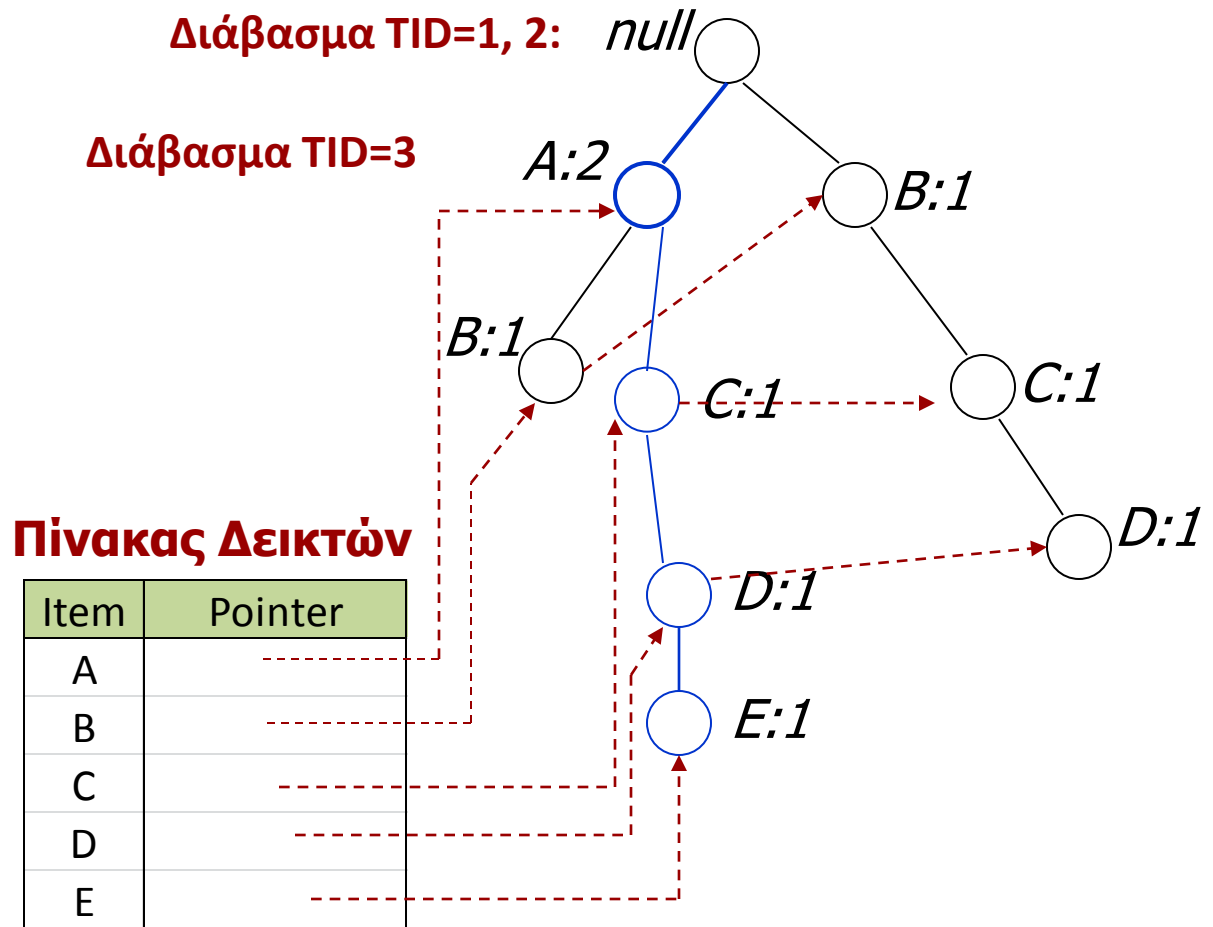
Κατασκευή FP-δένδρου (5/7)

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}



Κατασκευή FP-δένδρου (6/7)

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}

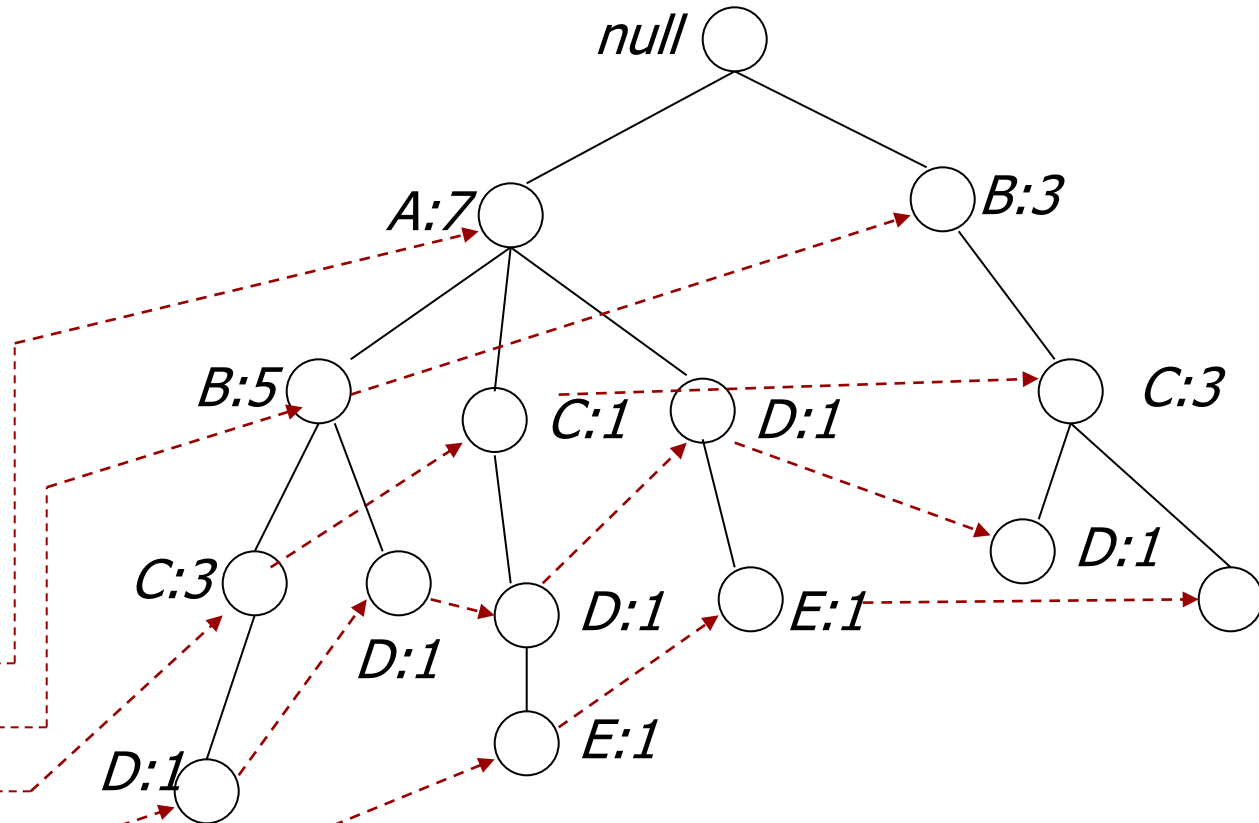


Κατασκευή FP-δένδρου (7/7)

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	



- Αφού έχουν διαβαστεί όλες οι συναλλαγές...



Μέγεθος FP-δένδρου


- Κάθε συναλλαγή αντιστοιχεί σε ένα μονοπάτι από τη ρίζα.
- Το μέγεθος του δένδρου είναι συνήθως μικρότερο των δεδομένων, αν υπάρχουν κοινά προθέματα.
 - Αν όλες οι συναλλαγές περιέχουν τα ίδια δεδομένα, τότε υπάρχει μόνο ένα κλαδί.
 - Αν όλες είναι διαφορετικές, ο χώρος είναι μεγαλύτερος...
 - ...γιατί αποθηκεύεται περισσότερη πληροφορία, όπως δείκτες μεταξύ των κόμβων αλλά και συχνότητες εμφάνισης.



Επιλογή προθέματος

- Το τελικό δένδρο, εξαρτάται από τη **διάταξη**:
 - άλλη διάταξη → άλλα προθέματα.
- (Συνήθως) μικρότερο δένδρο, αν δεν διατάσουμε τα αντικείμενα λεξικογραφικά, αλλά σύμφωνα με τη συχνότητα εμφάνισης.
- Αρχικά, διαβάζουμε όλα τα δεδομένα μια φορά ώστε να υπολογιστεί ο μετρητής υποστήριξης κάθε στοιχείου, και διατάσουμε τα στοιχεία με βάση αυτό (αγνοούμε όσα στοιχεία είναι μη συχνά).

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}



Εύρεση συχνών στοιχειοσυνόλων

- Είσοδος: FP-δένδρο.
- Έξοδος: Συχνά στοιχειοσύνολα και η υποστήριξη τους.
- Μέθοδος Διαίρει-και-Βασίλευε:
 - Χωρίζουμε τα στοιχειοσύνολα σε αυτά που τελειώνουν σε E, D, C, B, A.
 - Μετά αυτά που τελειώνουν σε E σε αυτά σε DE, CE, BE, AE κ.ο.κ.
 - Αν η διάταξη είναι βάσει της συχνότητας εμφάνισης, τότε χωρίζουμε τα στοιχειοσύνολα σε αυτά που τελειώνουν στο πιο σπάνιο στοιχείο, μετά στο δεύτερο πιο σπάνιο κ.ο.κ.

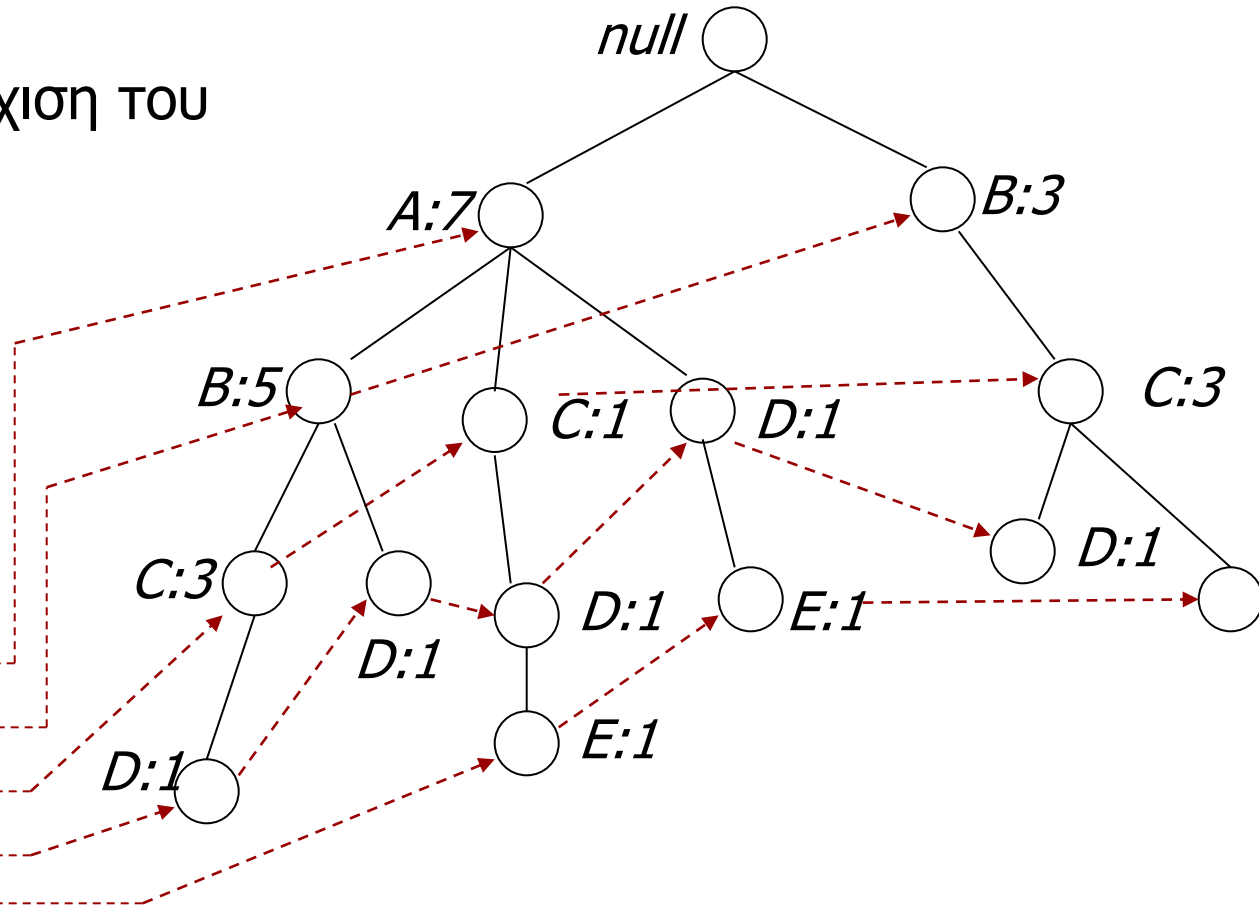


Εύρεση συχνών στοιχειοσυνόλων με χρήση του FP-δένδρου

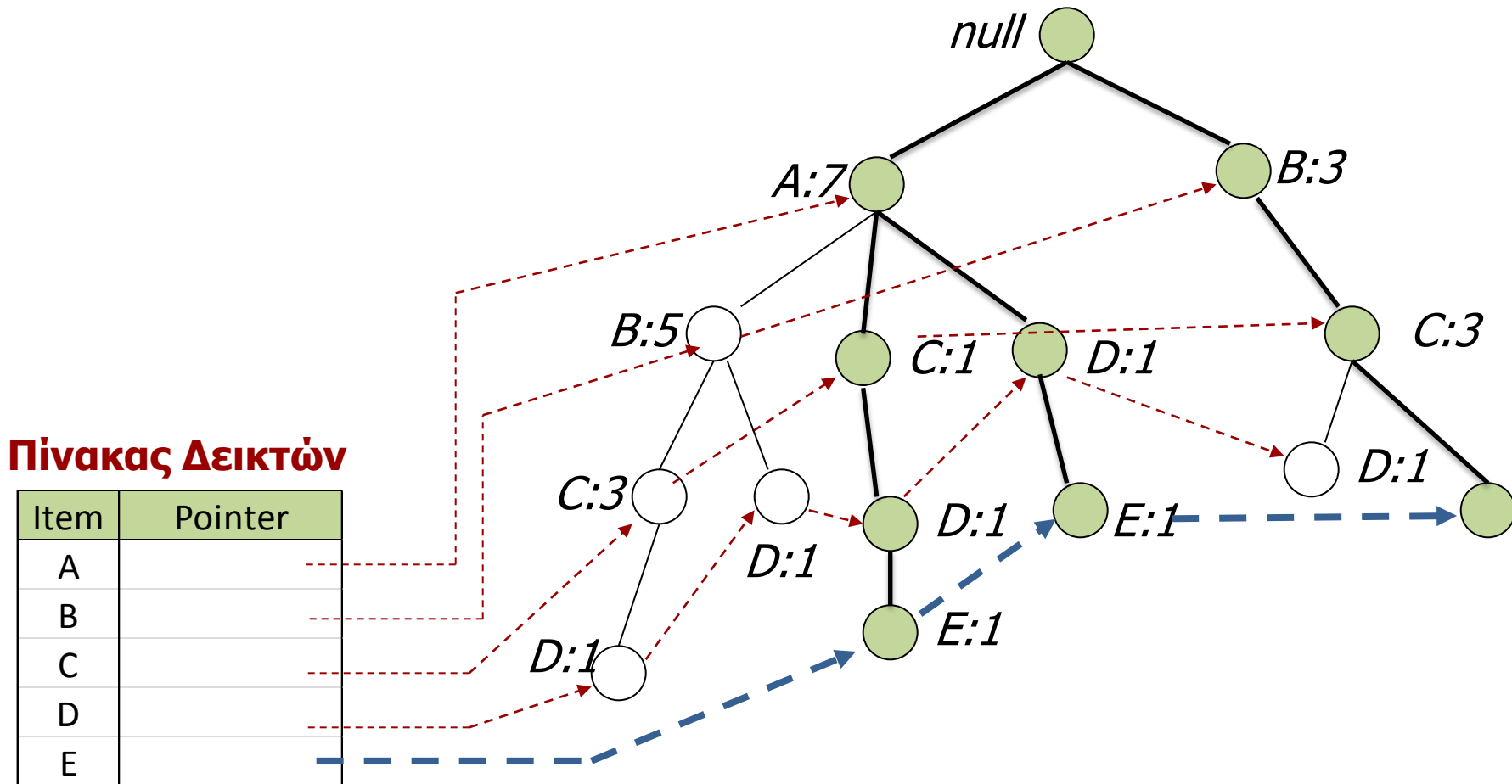
- Bottom-up διάσχιση του δένδρου.

Πίνακας Δεικτών

Item	Pointer
A	
B	
C	
D	
E	



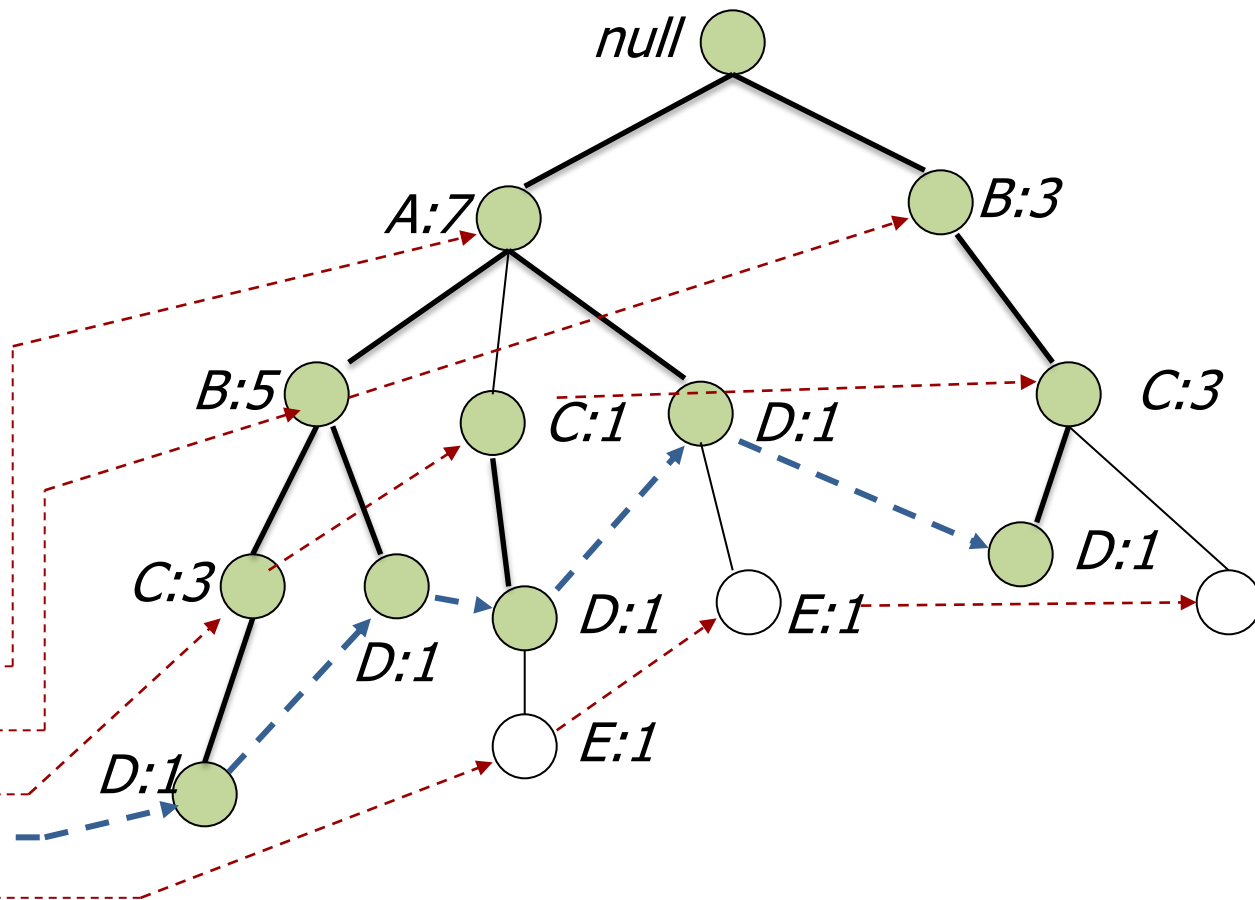
Συχνά στοιχειοσύνολα που τελειώνουν σε E



Για το D

Πίνακας Δεικτών

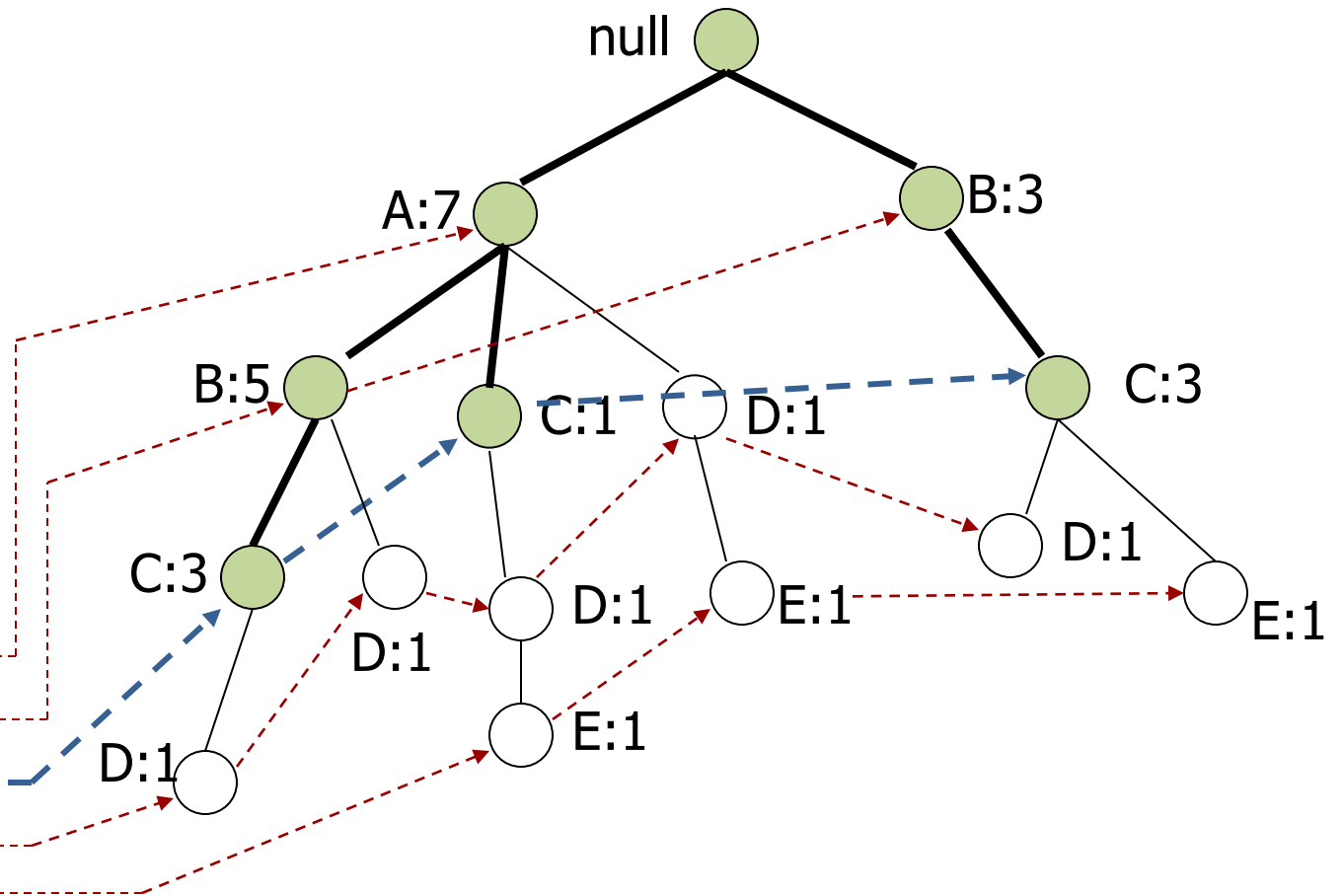
Item	Pointer
A	
B	
C	
D	
E	



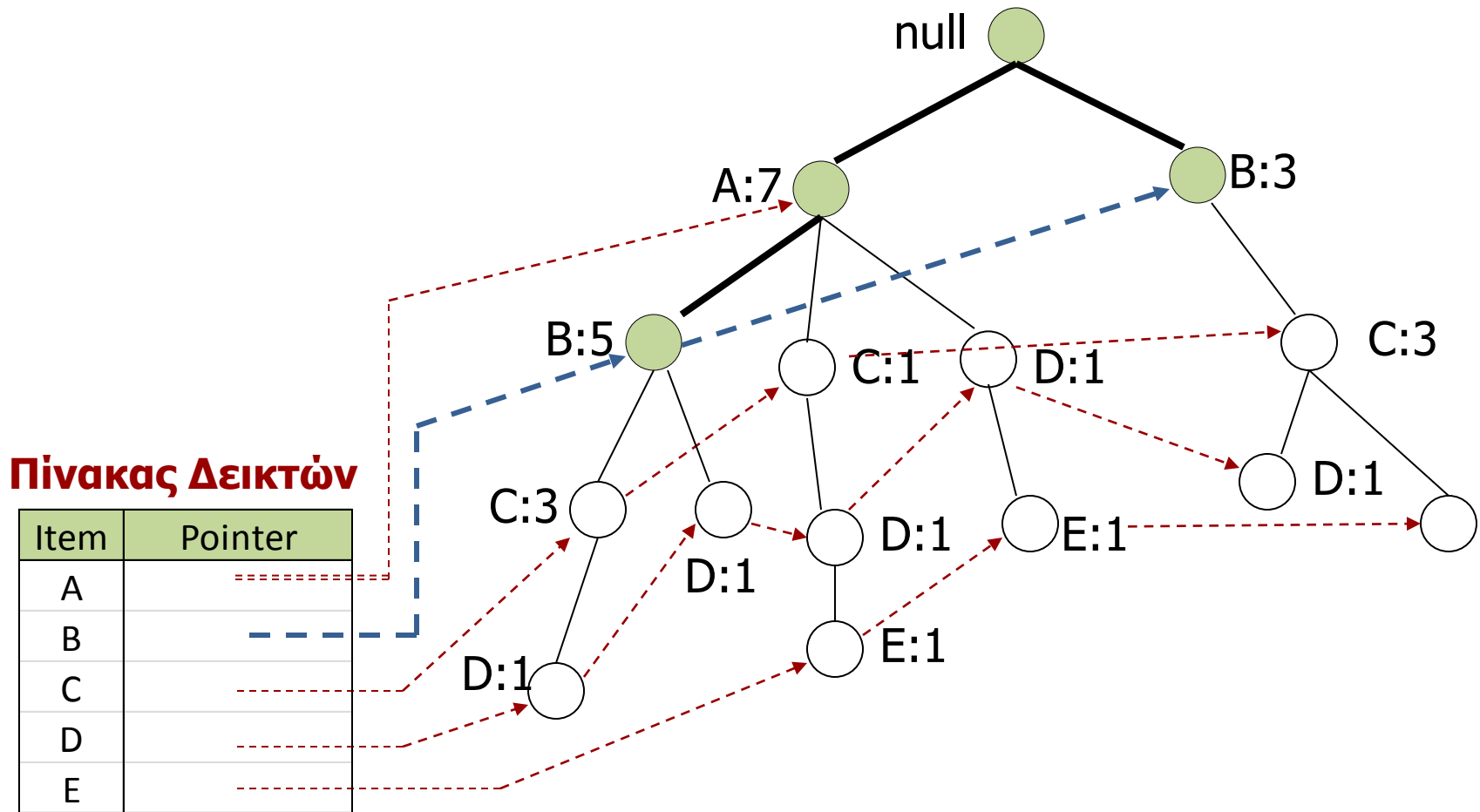
Για το C

Πίνακας Δεικτών

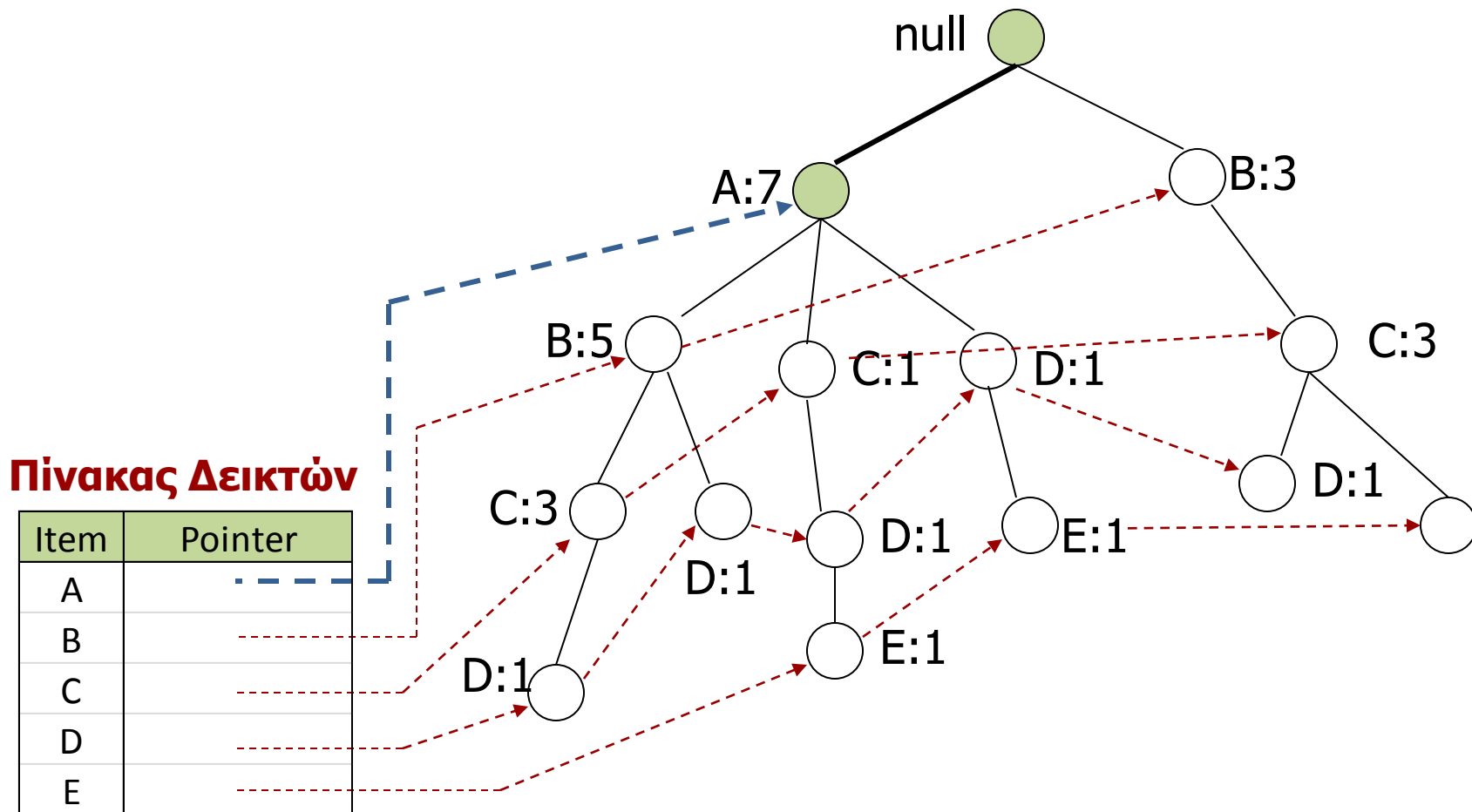
Item	Pointer
A	
B	
C	
D	
E	



Για το Β



Για το A



Συνοπτικά ο αλγόριθμος

- Σε κάθε βήμα, για το επίθεμα (suffix) X :
- **Φάση 1**
 - Κατασκευάζουμε το **προθεματικό δένδρο** για το X και υπολογίζουμε την υποστήριξη χρησιμοποιώντας τον πίνακα.
- **Φάση 2**
 - Αν είναι συχνό, κατασκευάζουμε το **υπο-συνθήκη δένδρο** για το X , σε βήματα:
 - Επανα-υπολογισμός υποστήριξης.
 - Περικοπή κόμβων με μικρή υποστήριξη.
 - Περικοπή φύλλων.

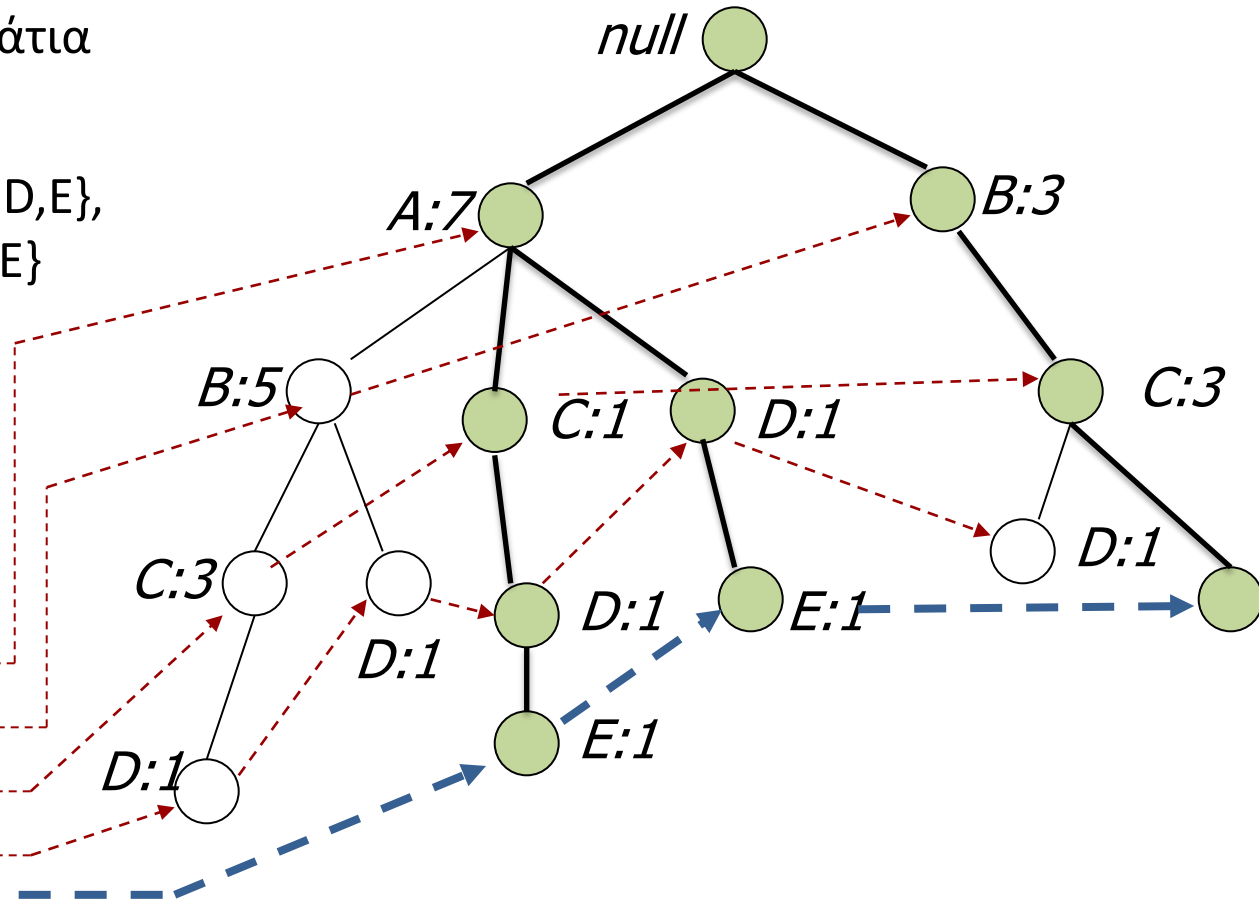


Φάση 1 (1/2)

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

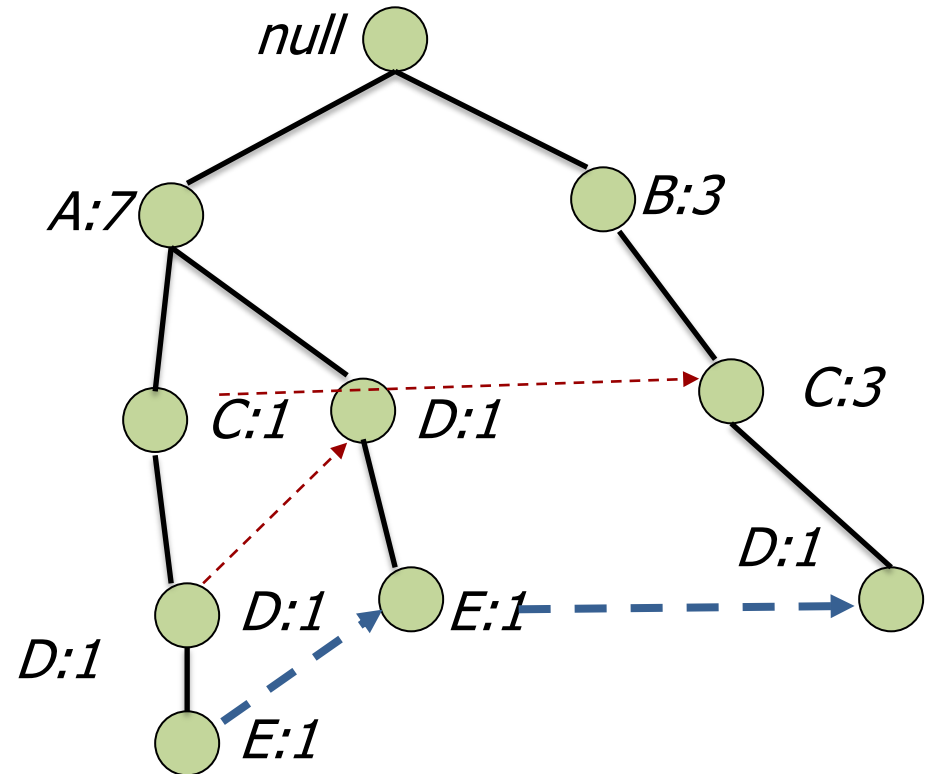
Πίνακας Δεικτών

Item	Pointer
A	
B	
C	
D	
E	



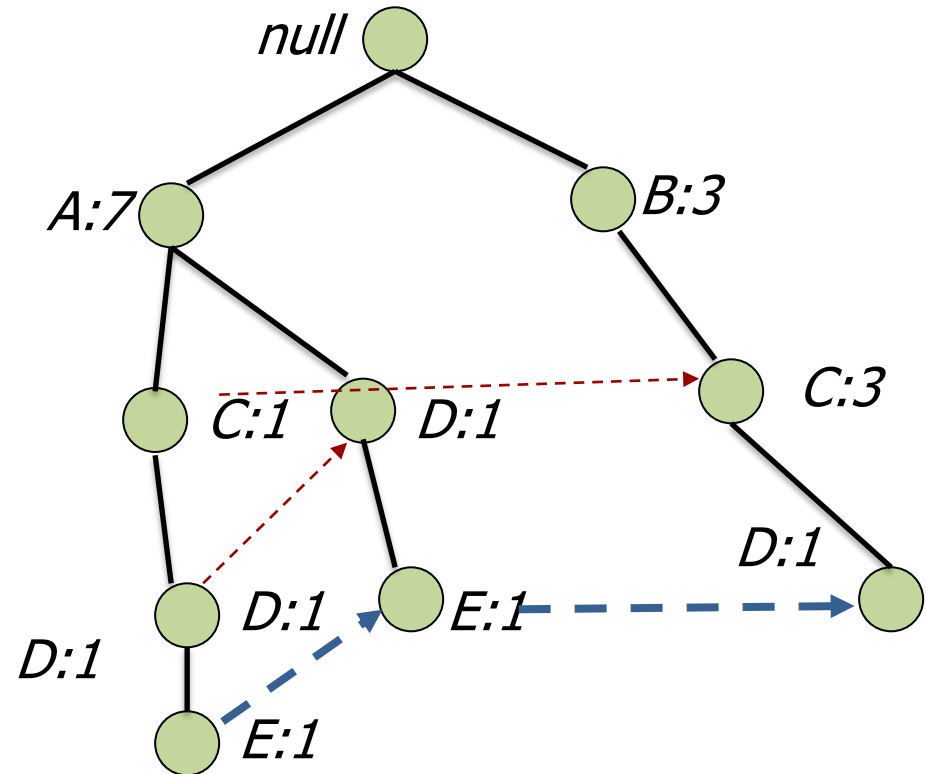
Φάση 1 (2/2)

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



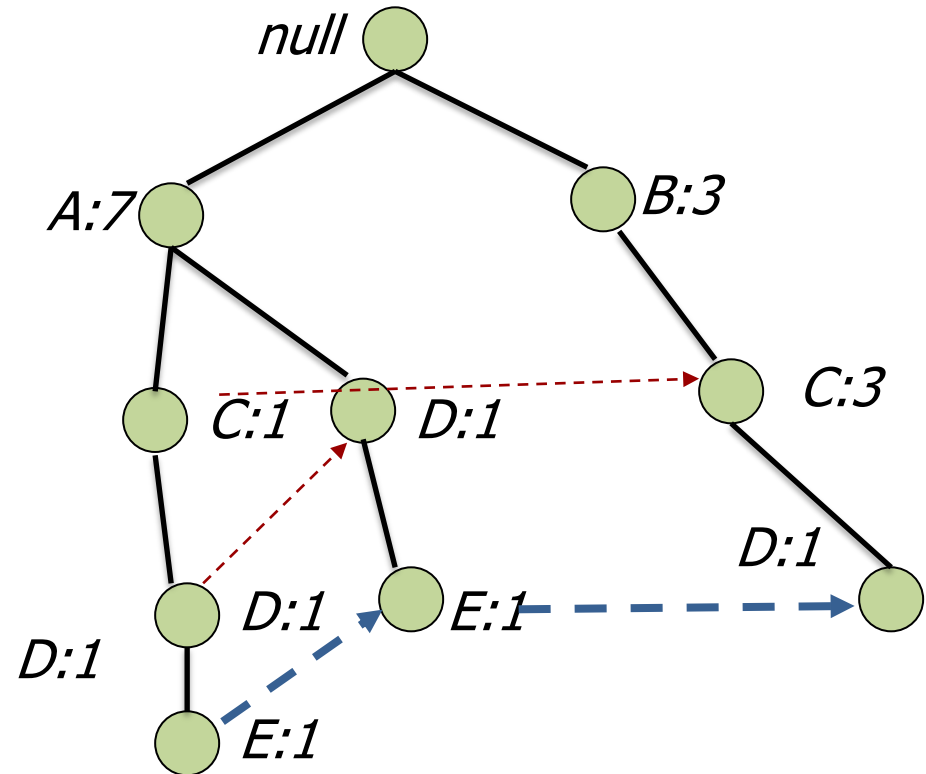
Μέτρηση Υποστήριξης

- Έστω $\text{minsup} = 2$
- Ακολουθούμε τους συνδέσμους αθροίζοντας $1+1+1=3 > 2$
- Οπότε $\{E\}$ συχνό
- ... άρα προχωράμε για DE, CE, BE, AE

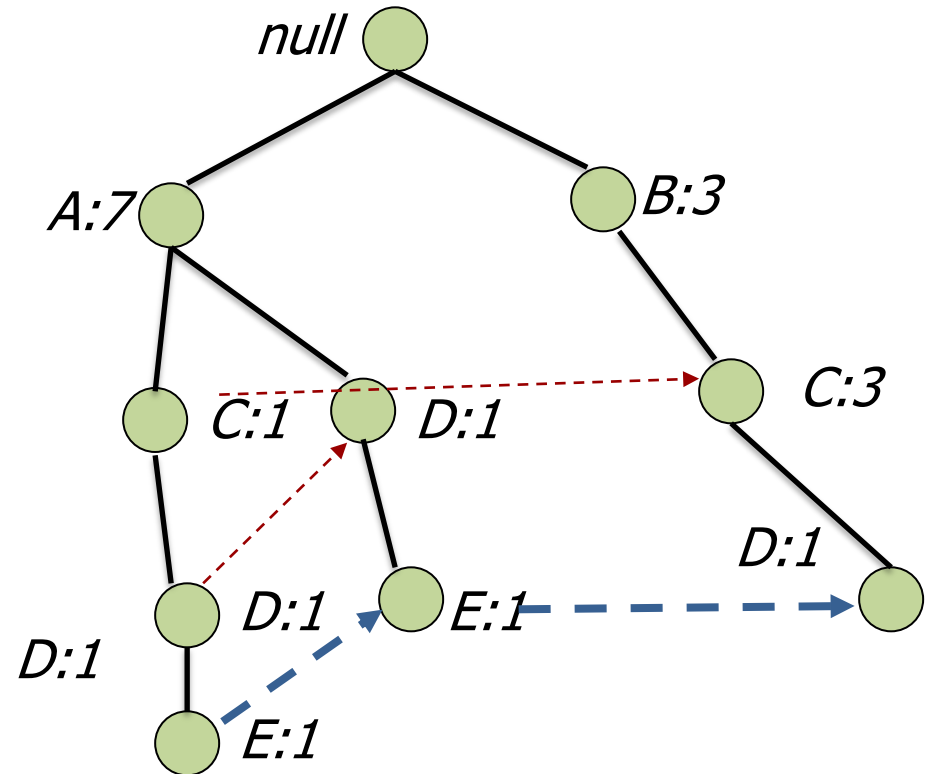


Φάση 2

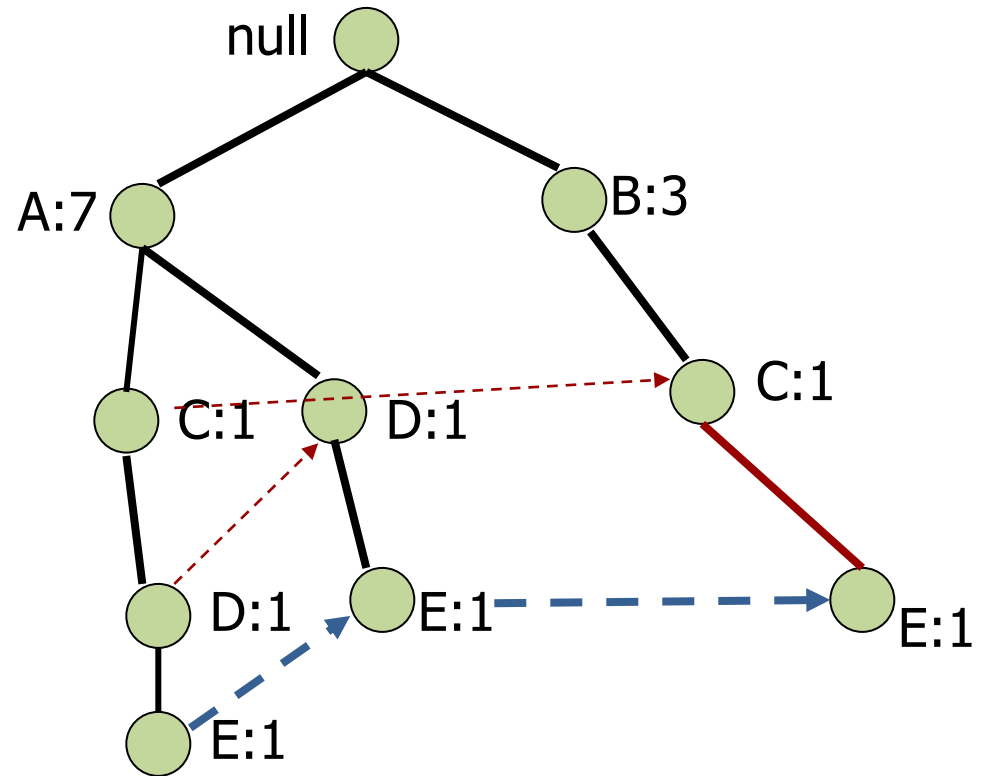
- {E} συχνό άρα προχωράμε για DE, CE, BE, AE.
- Μετατροπή των προθεματικών δένδρων σε FP-δένδρο υπό συνθήκες (conditional FP-tree).
- Δύο αλλαγές:
 - Αλλαγή των μετρητών.
 - Περικοπή.



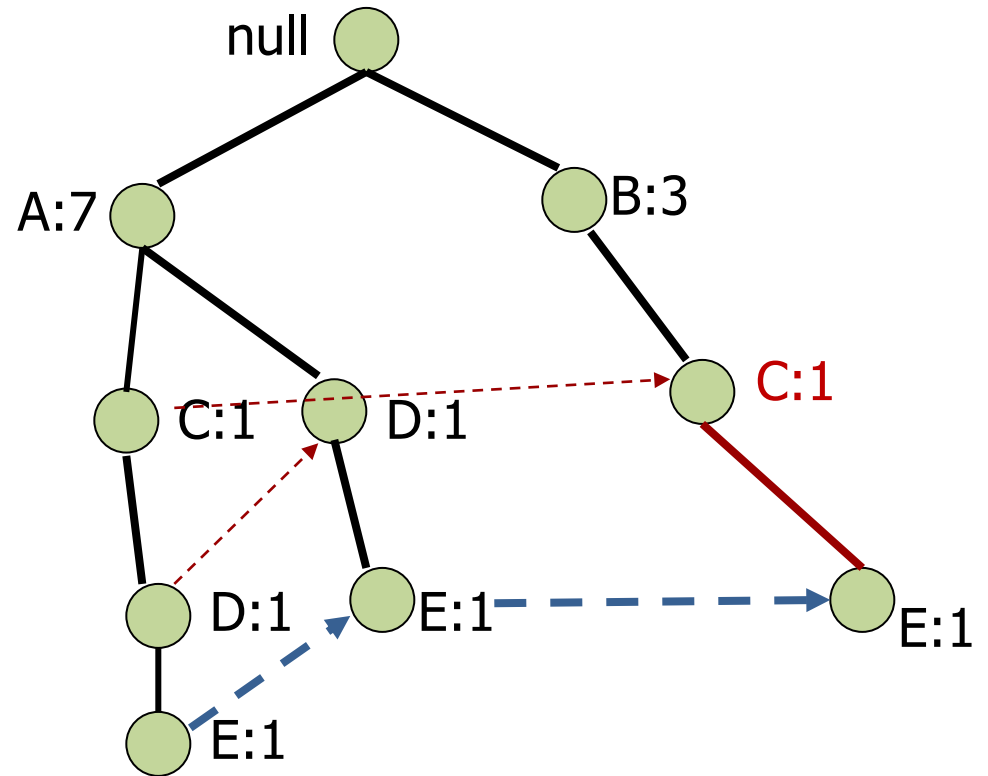
Αλλαγή μετρητών (1/7)



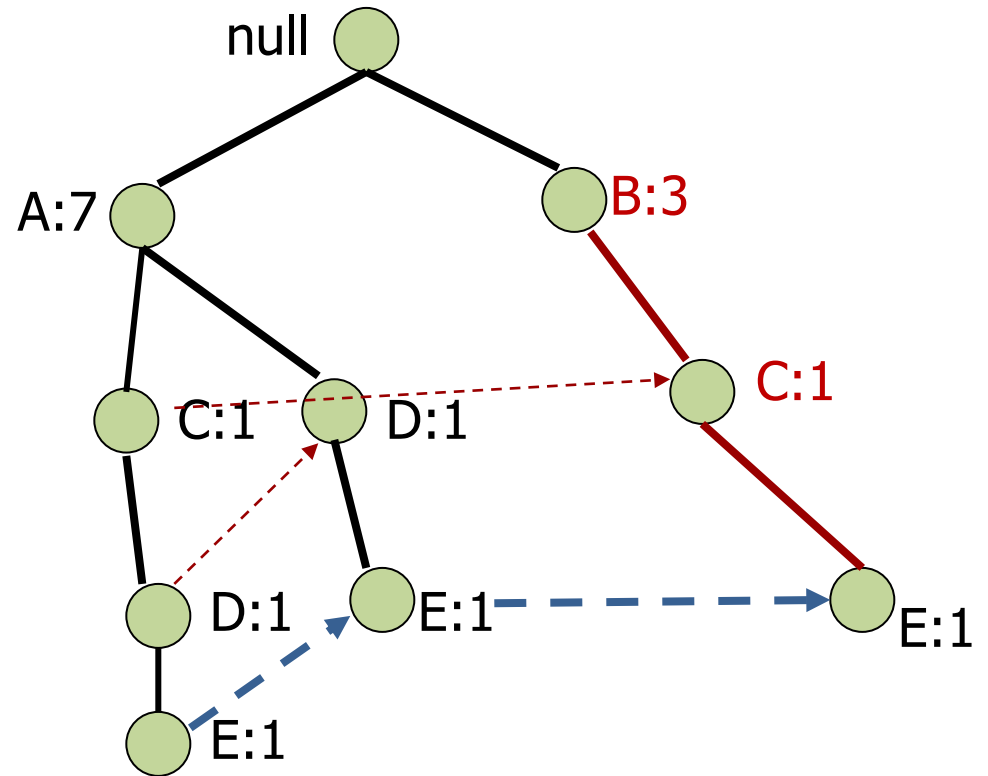
Αλλαγή μετρητών (2/7)



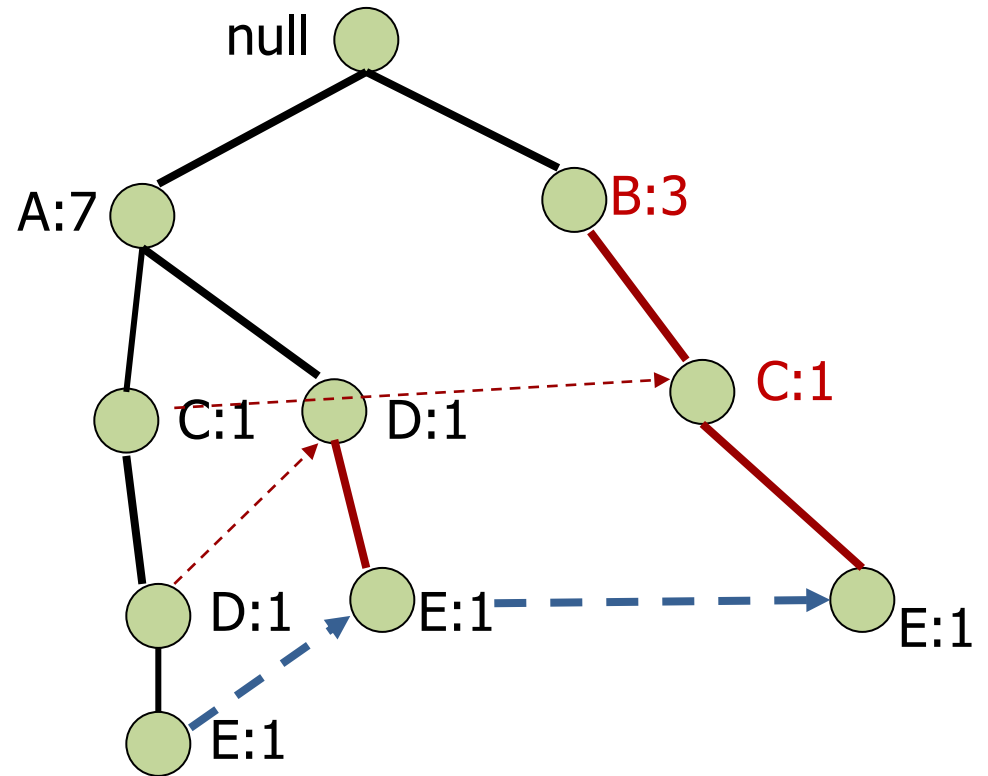
Αλλαγή μετρητών (3/7)



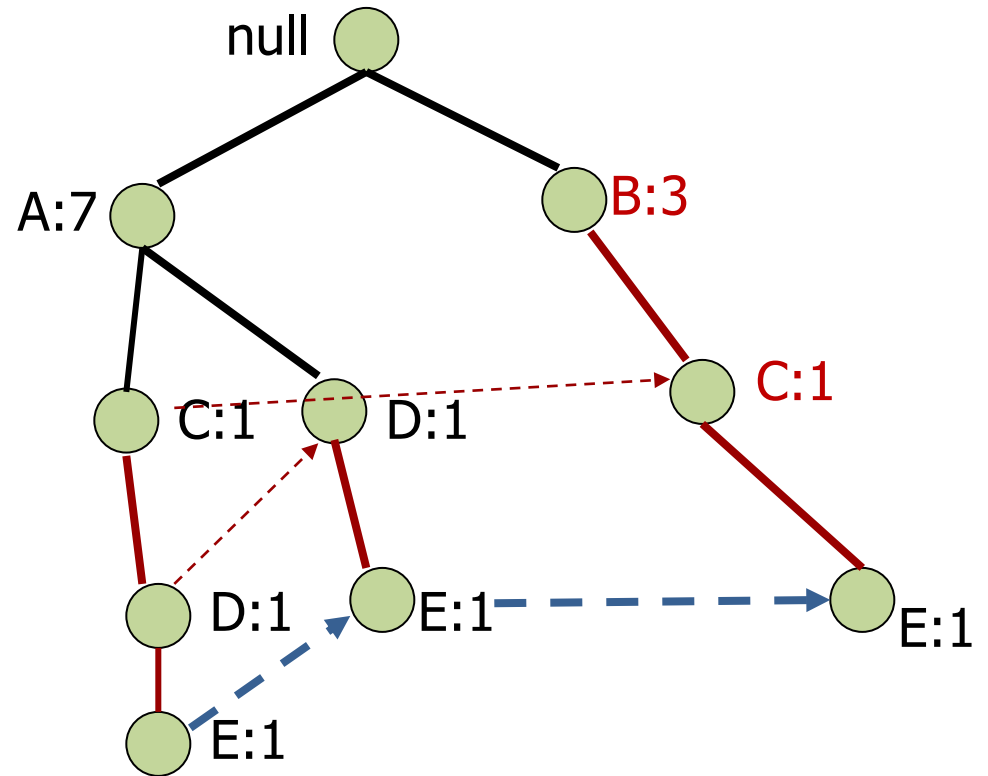
Αλλαγή μετρητών (4/7)



Αλλαγή μετρητών (5/7)

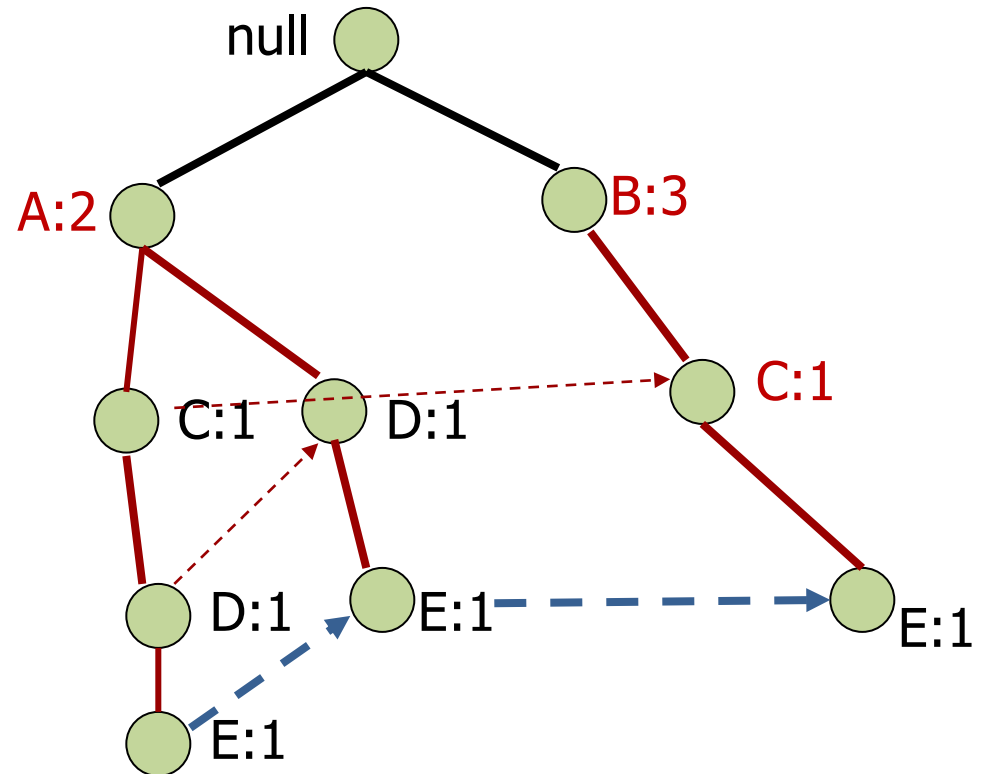


Αλλαγή μετρητών (6/7)



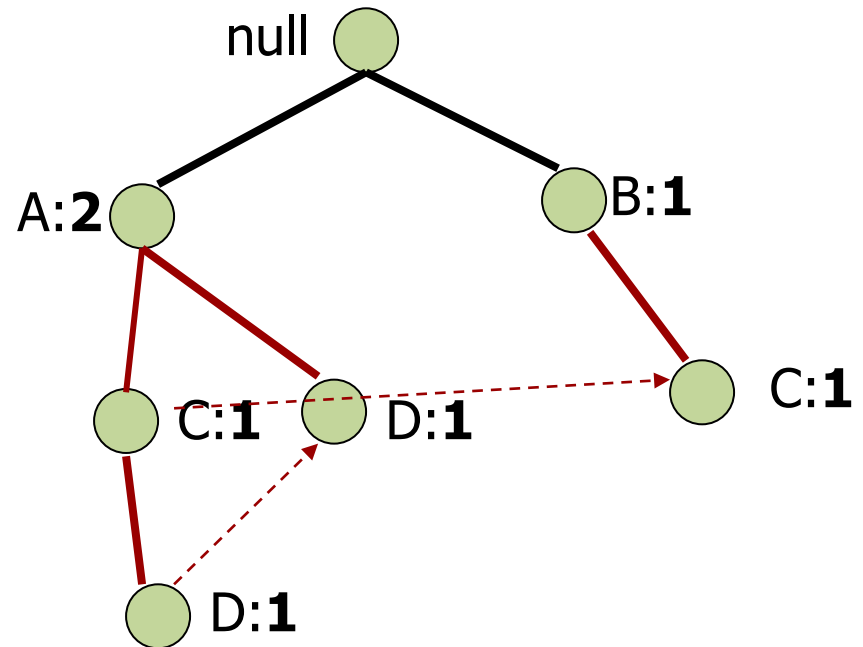
Αλλαγή μετρητών (7/7)

- Περικοπή (truncate):
Σβήσε τους κόμβους του E.

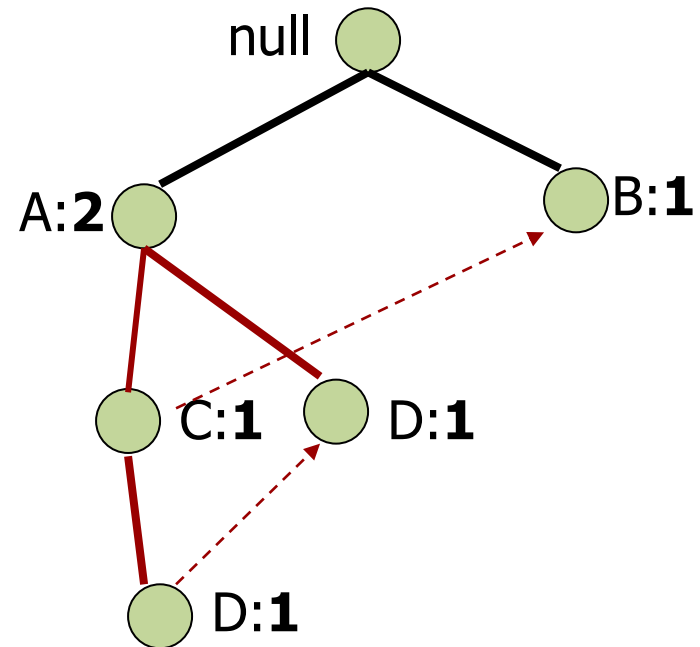


Περικοπή (1/2)

- Κάποια στοιχεία μπορεί να έχουν υποστήριξη μικρότερη της ελάχιστης (π.χ., B).
- Αυτό σημαίνει ότι το B εμφανίζεται μαζί με το E λιγότερο από \minsup φορές.
- Άρα $B \rightarrow$ περικοπή.

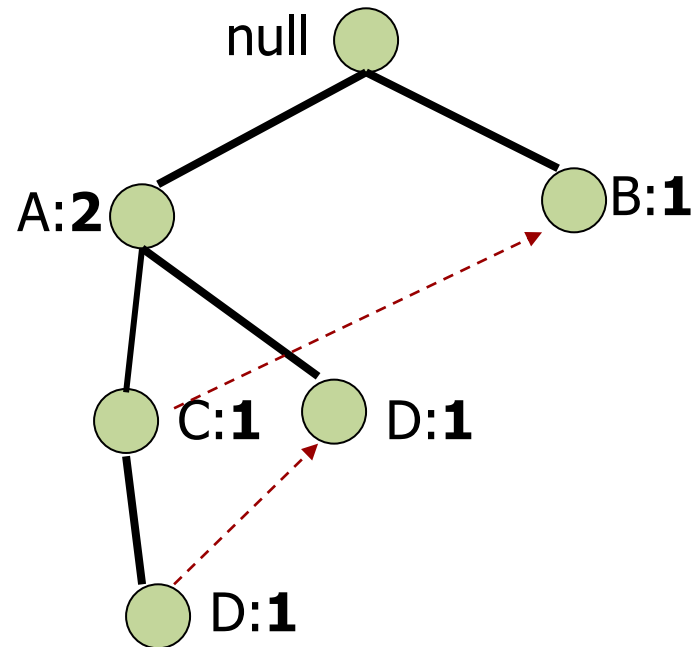


Περικοπή (2/2)



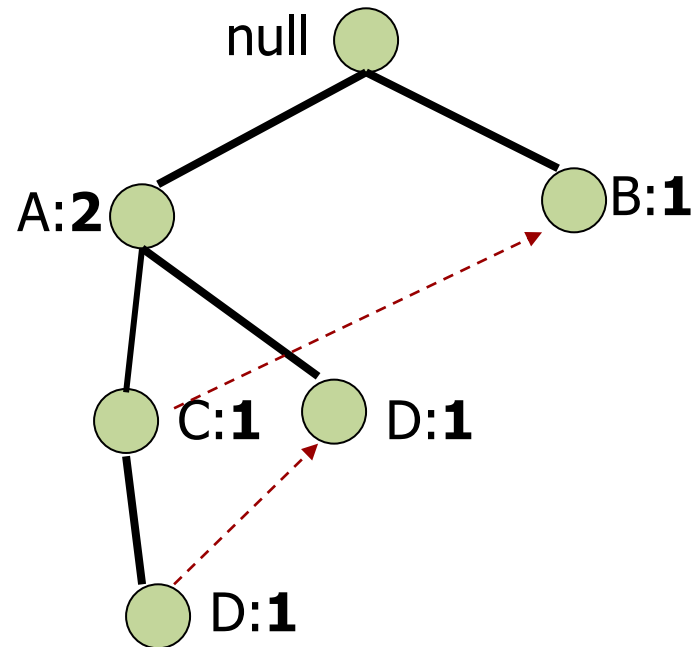
Αναδρομή

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



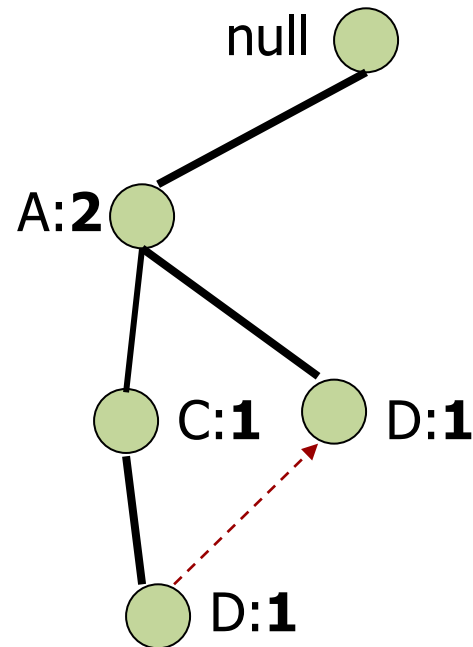
Φάση 1 (1/2)

- Βρίσκουμε όλα τα μονοπάτια που περιέχουν το D (DE).



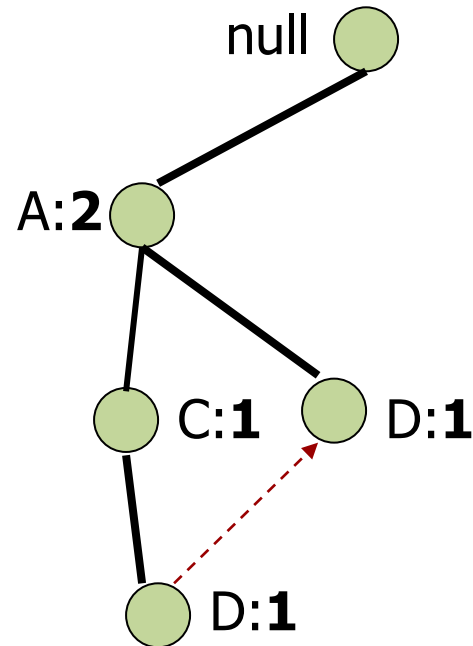
Φάση 1 (2/2)

- Βρίσκουμε όλα τα μονοπάτια που περιέχουν το D (DE).



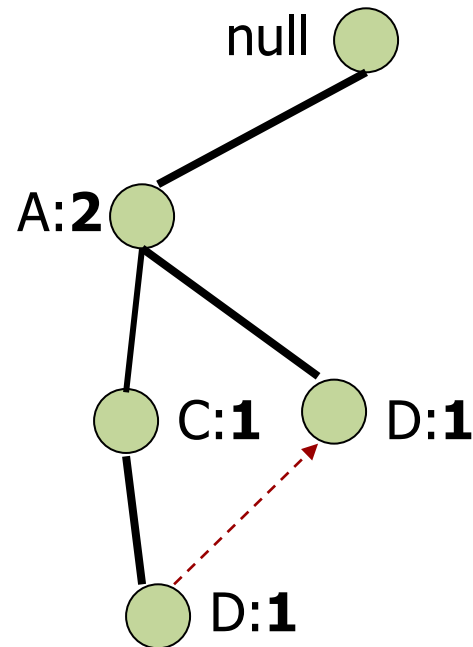
Υποστήριξη DE

- Ακολουθούμε τους συνδέσμους αθροίζοντας: $1+1=2 \geq 2$
- Οπότε $\{D, E\}$ συχνό.



Φάση 2: Υπο-συνθήκη δένδρο (1/6)

- Κατασκεύασε το υπο-συνθήκη FP-δένδρο για το {D, E}.
1. Αλλαγή υποστήριξης.
 2. Περικοπές κόμβων.



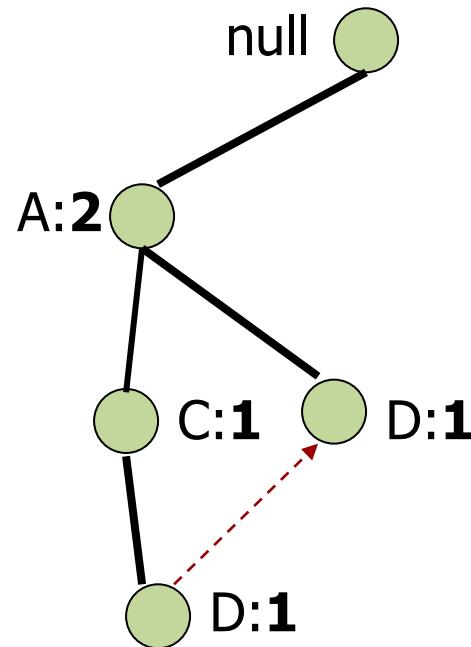
Φάση 2: Υπο-συνθήκη δένδρο (2/6)

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

1. Αλλαγή υποστήριξης:

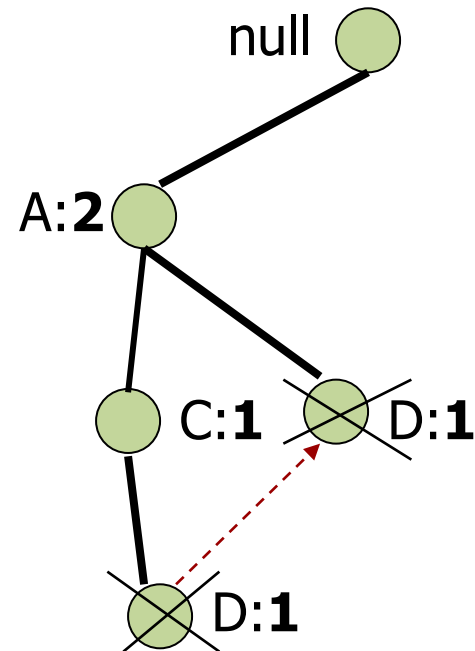
Δεν υπάρχει καμία.

2. Περικοπές κόμβων.



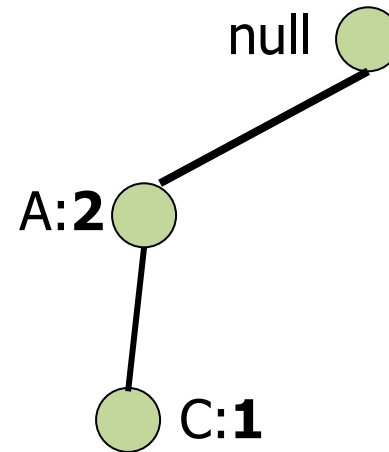
Φάση 2: Υπο-συνθήκη δένδρο (3/6)

2. Περικοπές κόμβων.



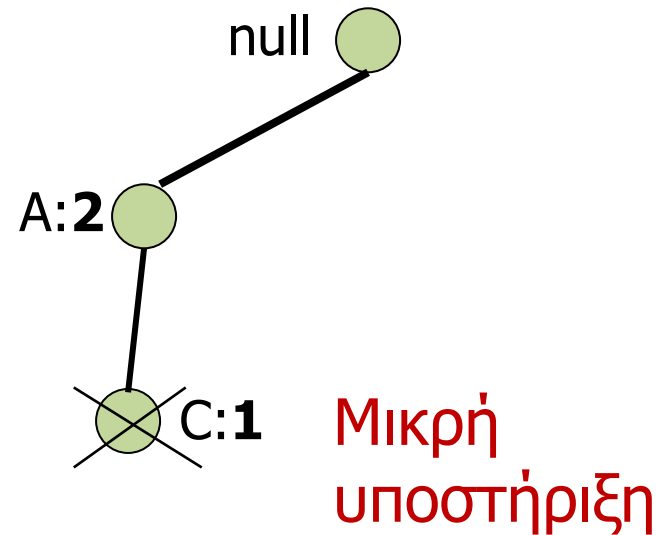
Φάση 2: Υπο-συνθήκη δένδρο (4/6)

2. Περικοπές κόμβων.

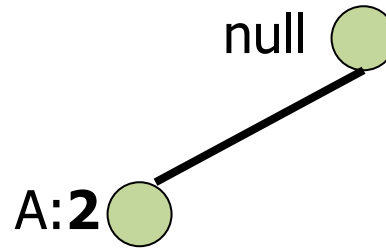


Φάση 2: Υπο-συνθήκη δένδρο (5/6)

2. Περικοπές κόμβων.



Φάση 2: Υπο-συνθήκη δένδρο (6/6)

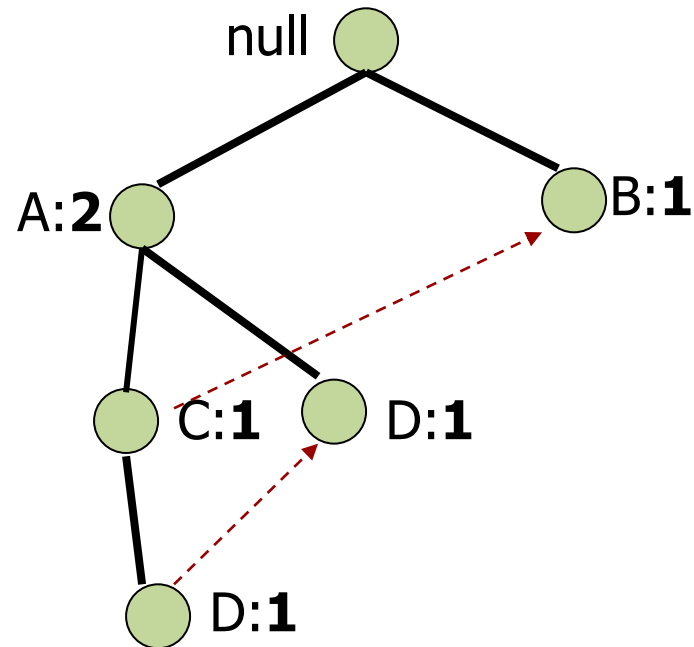


- Τελικό υπο-συνθήκη FP-δένδρο για το $\{D, E\}$
- Υποστήριξη του A είναι $\geq \text{minsup}$ $\rightarrow \{A, D, E\}$ συχνό.
- Αφού μόνο ένας κόμβος απέμεινε, επιστροφή στο επόμενο υποπρόβλημα.



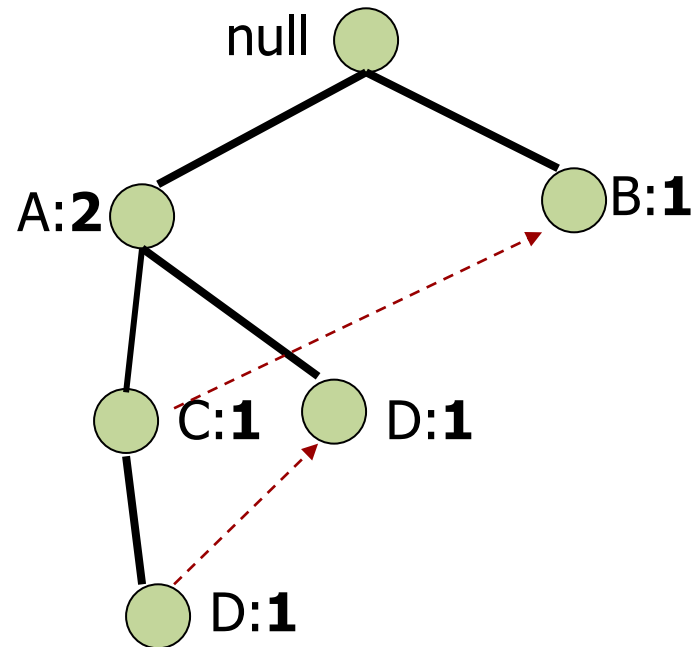
Αναδρομή

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



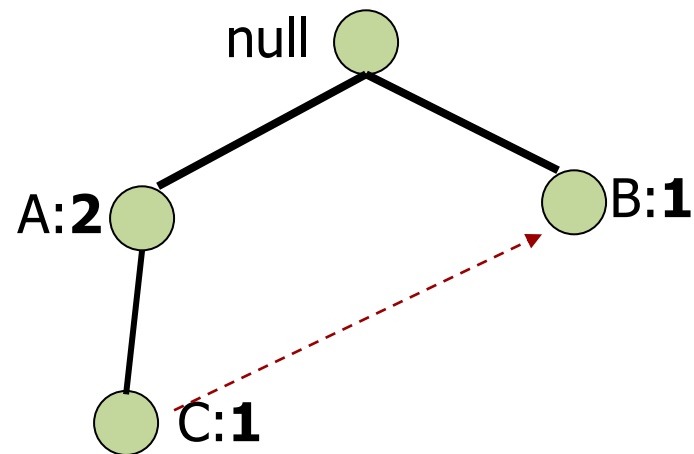
Φάση 1 (1/2)

- Όλα τα μονοπάτια που περιέχουν το C (CE).



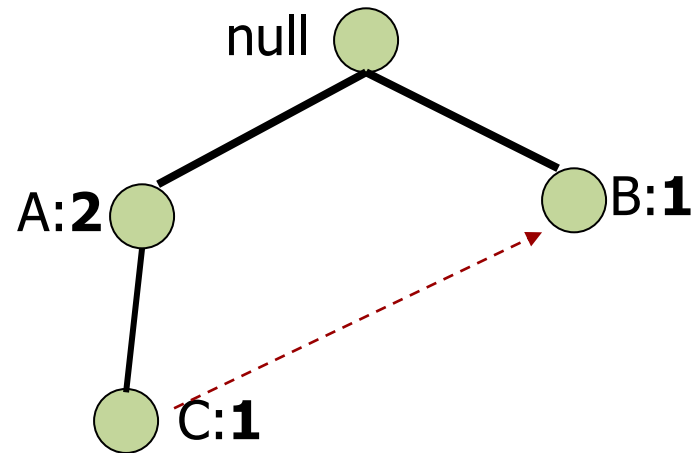
Φάση 1 (2/2)

- Όλα τα μονοπάτια που περιέχουν το C (CE).



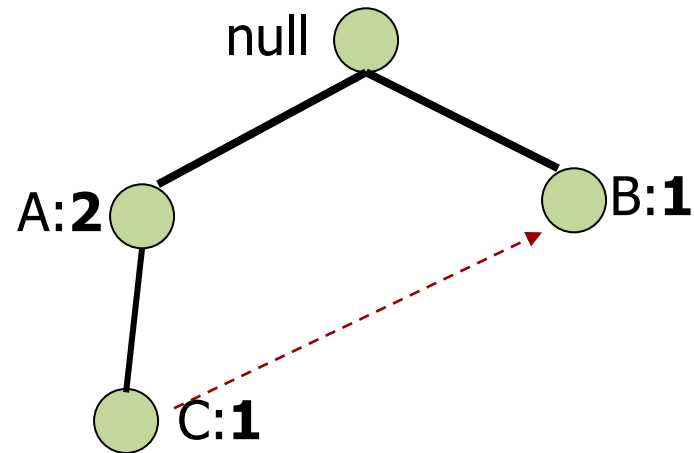
Υποστήριξη CE

- {C, E} συχνό.



Φάση 2 (1/4)

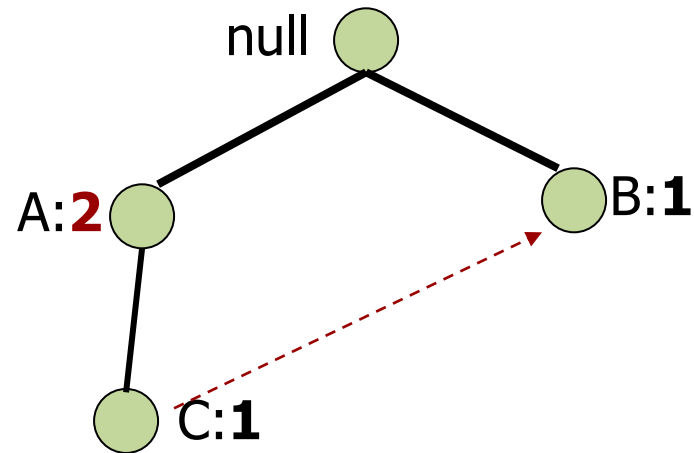
- Κατασκεύασε το υπο-συνθήκη FP-δένδρο για το {C, E}.
1. Αλλαγή υποστήριξης.
 2. Περικοπές κόμβων.



Φάση 2 (2/4)

- Κατασκεύασε το υποσυνθήκη FP-δένδρο για το {C, E}.

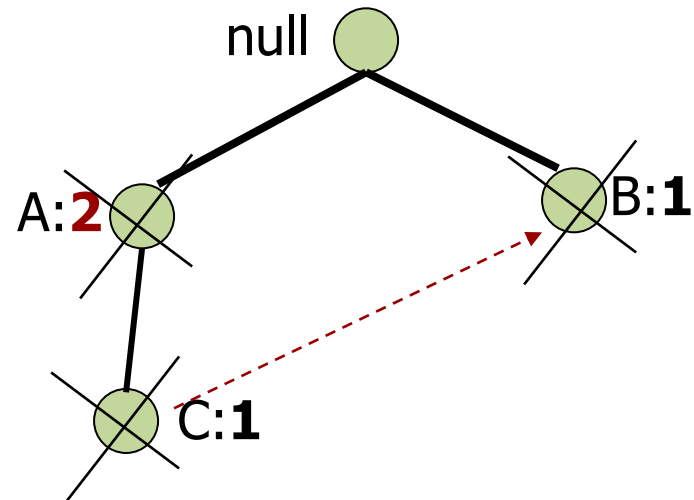
- Αλλαγή υποστήριξης
- Περικοπές κόμβων



Φάση 2 (3/4)

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

1. Αλλαγή υποστήριξης
2. Περικοπές κόμβων



Φάση 2 (4/4)

2. Περικοπές κόμβων

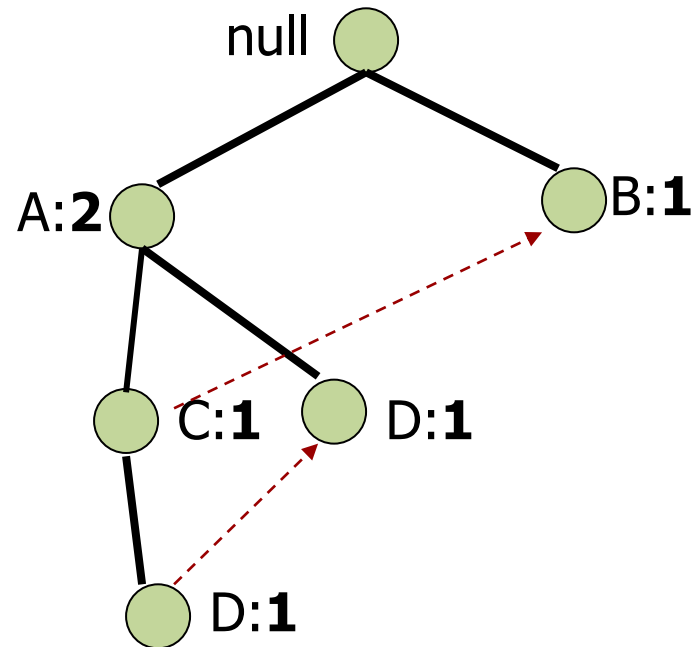
null 

- Άρα, επιστροφή στο επόμενο υποπρόβλημα.



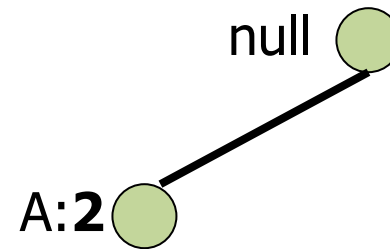
Αναδρομή

- Υπο-συνθήκη FP-δένδρο για το E.
- Ο αλγόριθμος επαναλαμβάνεται για το $\{D, E\}$, $\{C, E\}$, $\{A, E\}$.



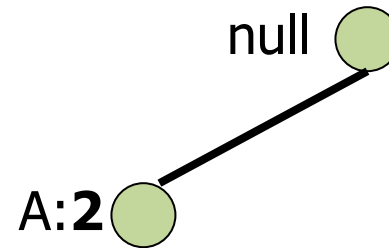
Φάση 1

- Όλα τα μονοπάτια που περιέχουν το A (ΑΕ).



Υποστήριξη ΑΕ

- $\{A, E\}$ συχνό.
- Δε χρειάζεται να φτιάξουμε υπο-συνθήκη FP-δένδρο για το $\{A, E\}$.

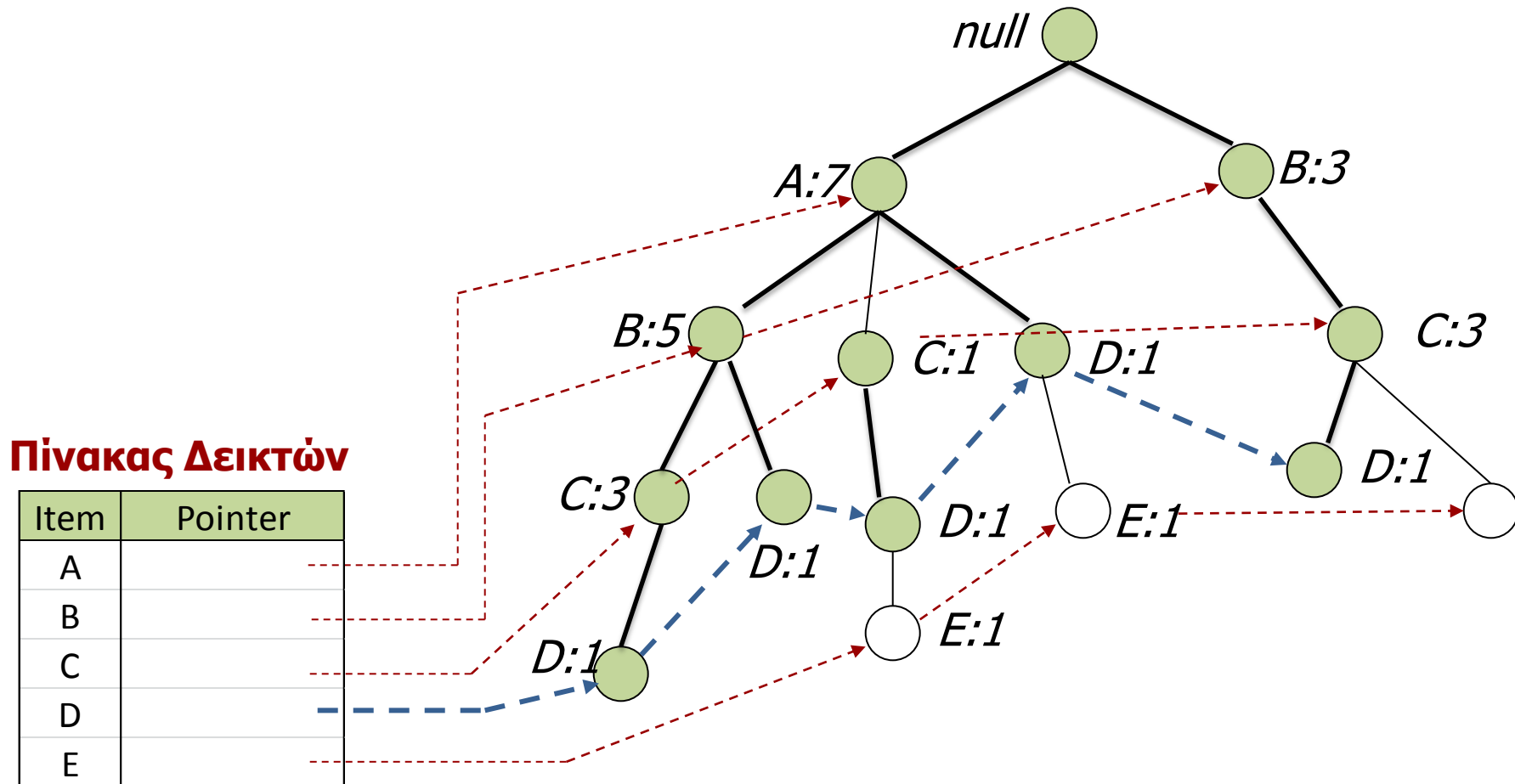


Συνολικά για το E

- Έχουμε τα εξής συχνά στοιχειοσύνολα:
 - $\{E\}$ $\{D, E\}$ $\{A, D, E\}$ $\{C, E\}$ $\{A, E\}$.
 - Συνεχίζουμε για το D.

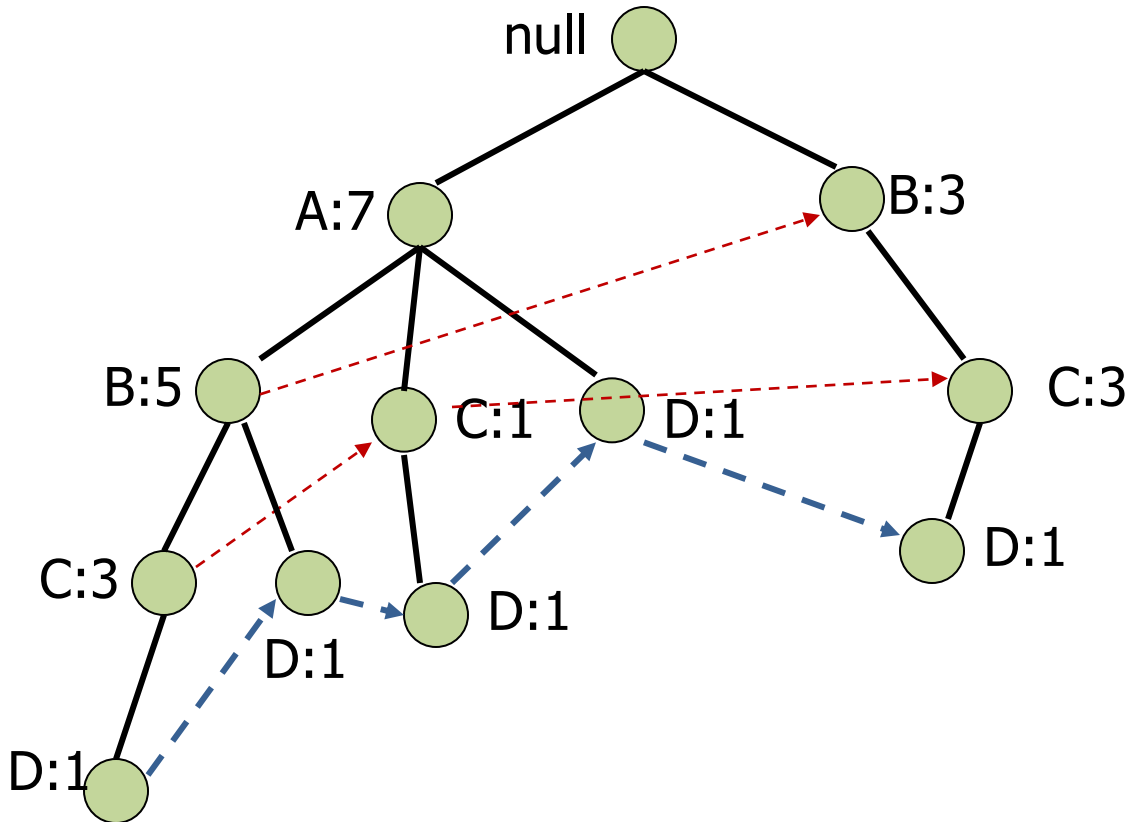


Συχνά στοιχειοσύνολα που λήγουν σε D

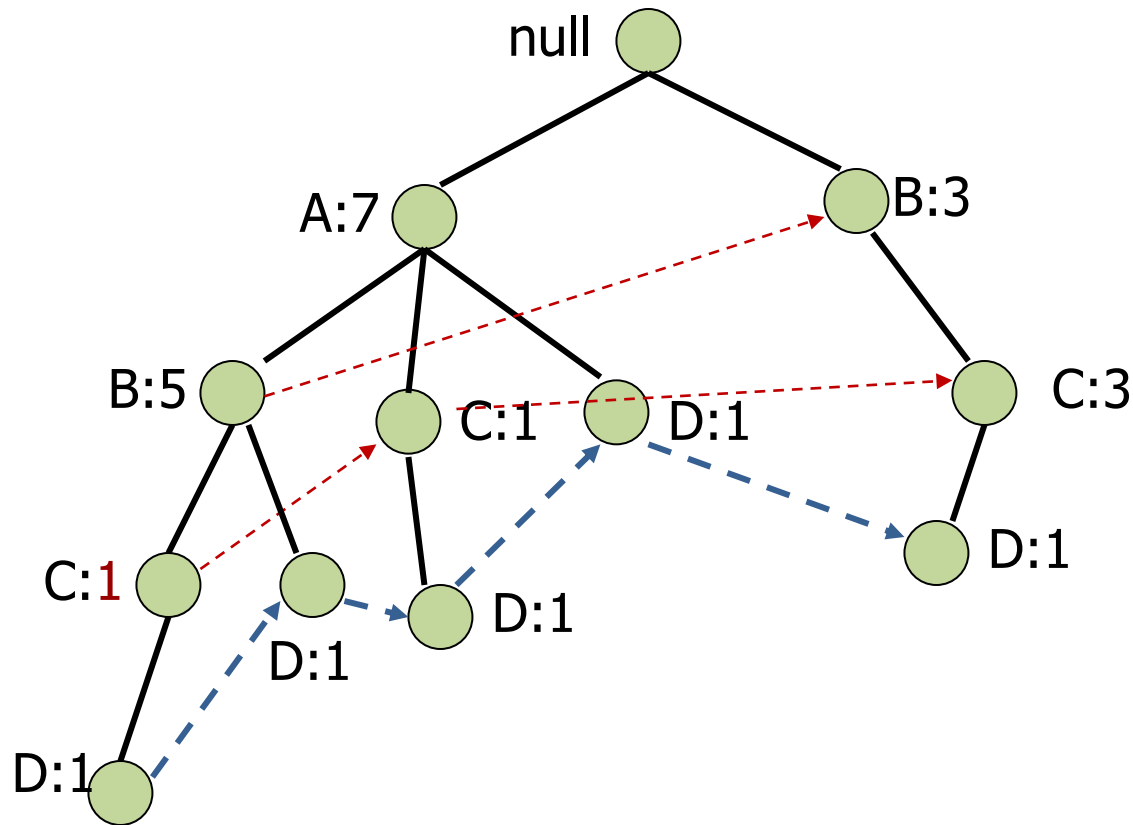


Φάση 1

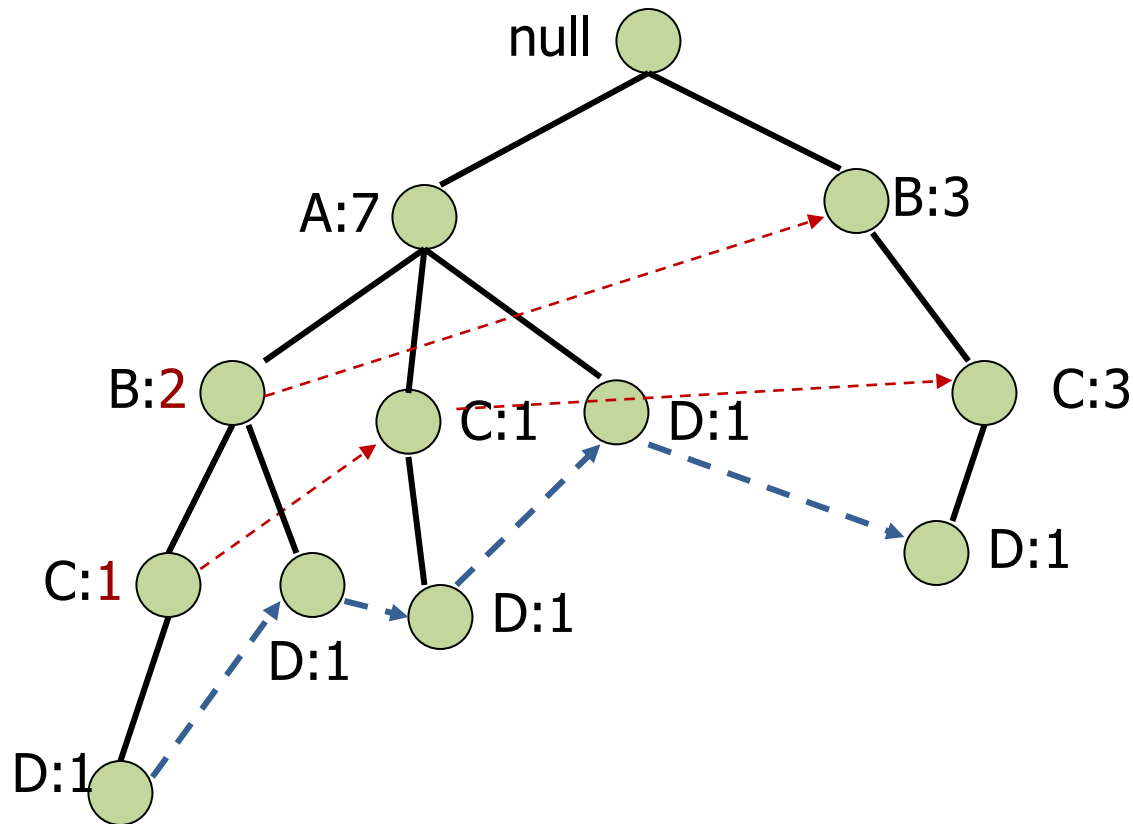
- Κρατάμε όλα τα προθεματικά μονοπάτια που περιέχουν το D.
- Υποστήριξη $5 > 2$: άρα συχνό το D.
- Μετατροπή του προθεματικού δένδρου σε FP-δένδρο υπό συνθήκη.



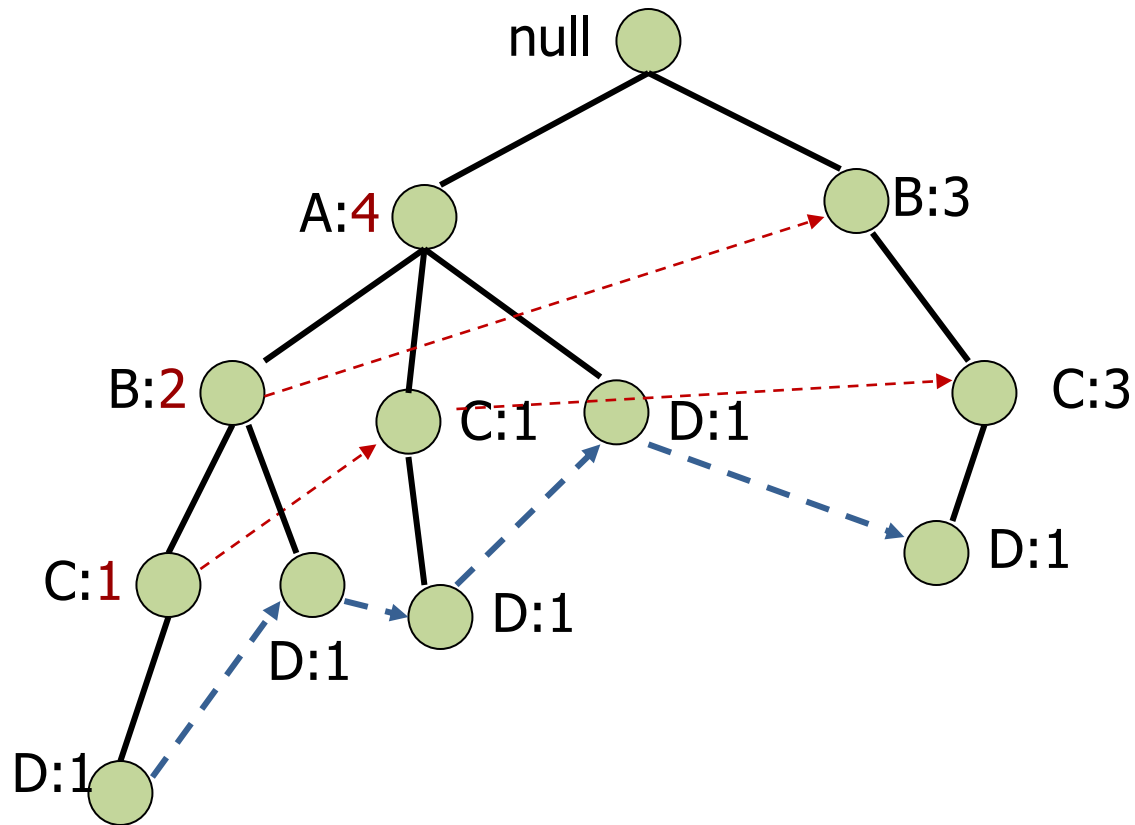
Αλλαγή υποστήριξης (1/5)



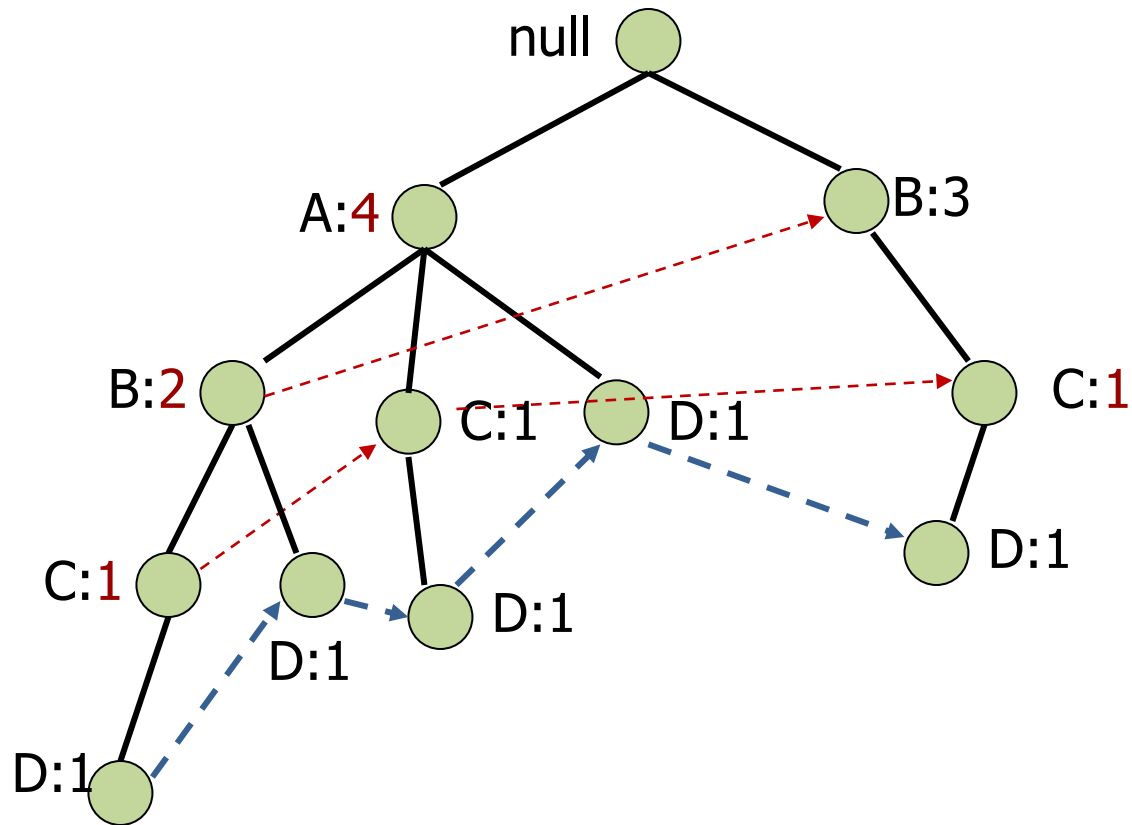
Αλλαγή υποστήριξης (2/5)



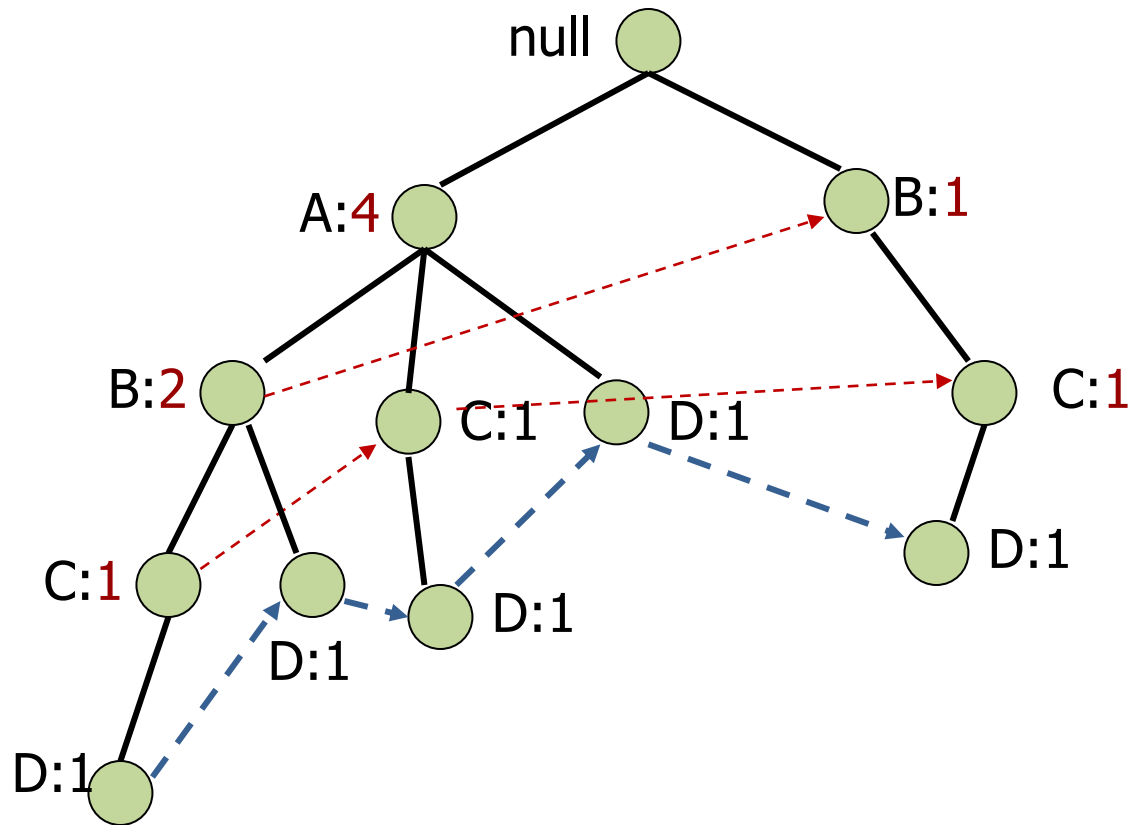
Αλλαγή υποστήριξης (3/5)



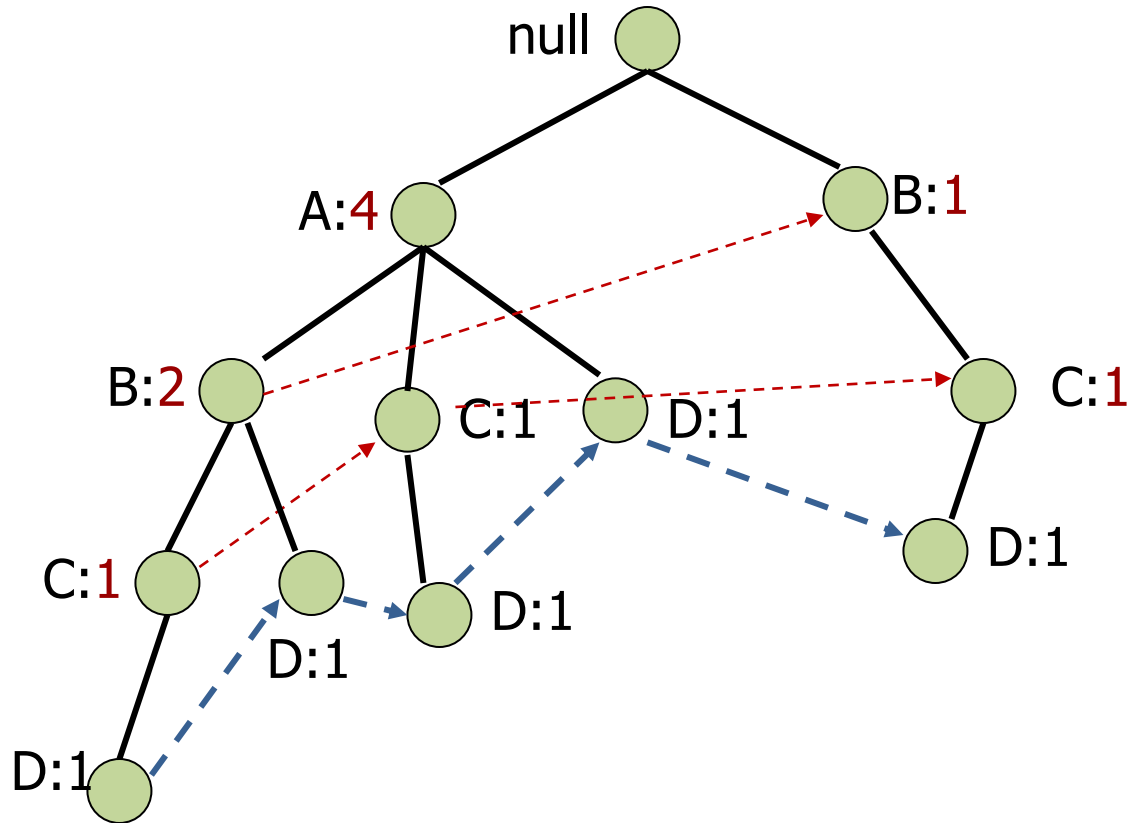
Αλλαγή υποστήριξης (4/5)



Αλλαγή υποστήριξης (5/5)



Επόμενο βήμα: Περικοπή κόμβων



Παρατηρήσεις

- Εφαρμογή τεχνικής διαίρει-και-βασίλευε.
- Σε κάθε αναδρομικό βήμα, λύνεται και ένα υπο-πρόβλημα:
 - Κατασκευάζεται το προθεματικό δένδρο.
 - Υπολογίζεται η νέα υποστήριξη για τους κόμβους του.
 - Περικόπτονται 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.
- Η απόδοση του FP-Growth εξαρτάται από τον παράγοντα συμπίεσης του συνόλου των δεδομένων (compaction factor).
 - Βοηθάει η ταξινόμηση αντικειμένων κατά φθίνουσα σειρά υποστήριξης.



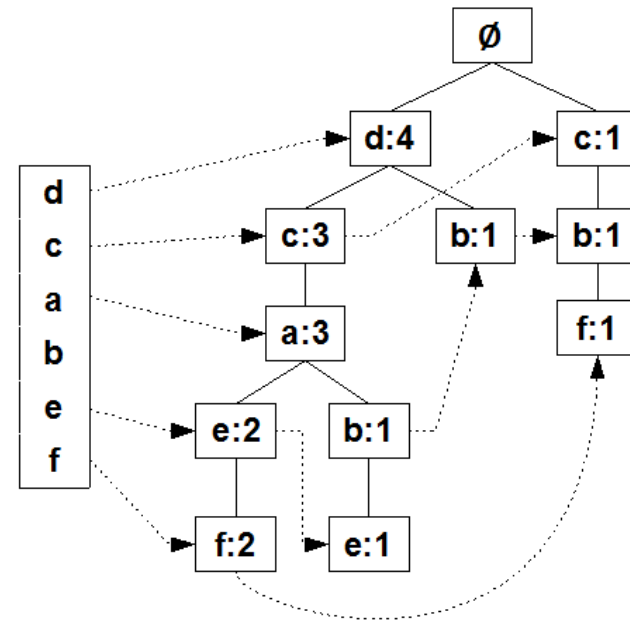
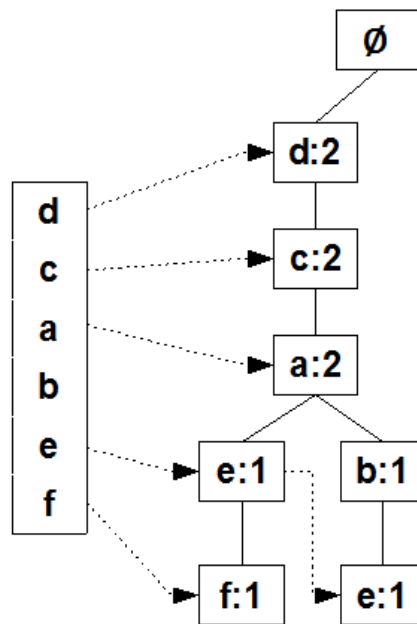
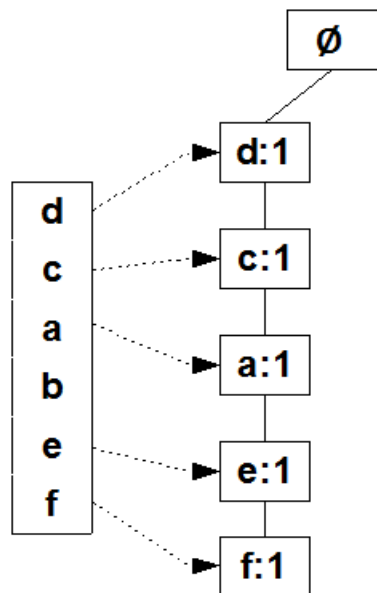
Άλλο ένα παράδειγμα (1/6)

Κωδ.	Εγγραφή	Υποστ. Αντικειμένων	Αναταξινόμηση
1	$\{a, c, d, e, f\}$	$d:4, c:4$	$\{d, c, a, e, f\}$
2	$\{a, b, c, d, e\}$	$a:3, b:3, e:3, f:3$	$\{d, c, a, b, e\}$
3	$\{b, d, g\}$	$g:1, h:1$	$\{d, b\}$
4	$\{b, c, f\}$		$\{c, b, f\}$
5	$\{a, c, d, e, f, h\}$		$\{d, c, a, e, f\}$

Minsup = 2



Άλλο ένα παράδειγμα (2/6)

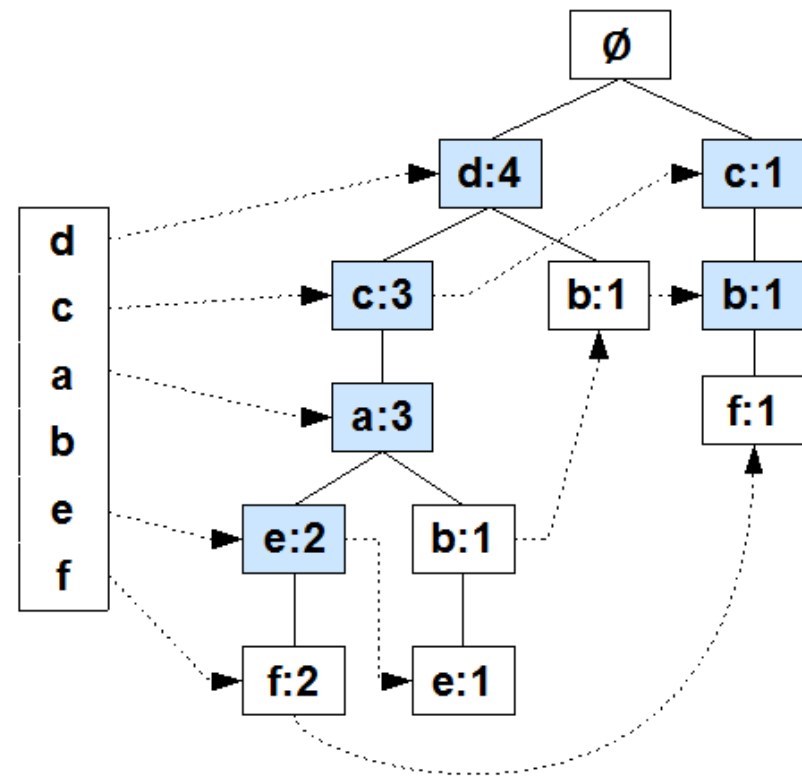


- {d, c, a, e, f}*
- {d, c, a, b, e}*
- {d, b}*
- {c, b, f}*
- {d, c, a, e, f}*



Άλλο ένα παράδειγμα (3/6)

Επίθεμα	Υ.Σ. Μονοπάτια
f	{(dcae:2), (cb:1)}
e	{(dca:2), (dcab:1)}
b	{(dca:1), (d:1), (c:1)}
a	{(dc:3)}
c	{(d:3), \emptyset :1}
d	{ \emptyset :4}

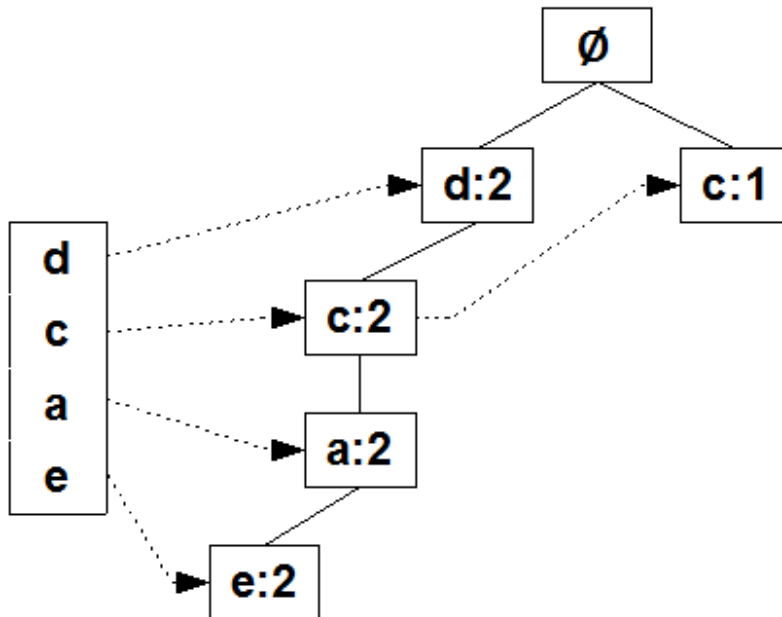


Υπό Συνθήκη FP-tree | f



Άλλο ένα παράδειγμα (4/6)

f	$\{(dcae:2), (cb:1)\}$
---	------------------------



Επίθεμα	Υ.Σ. Μονοπάτια f
ef	$\{(dca:2)\}$
af	$\{(dc:2)\}$
df	$\{\emptyset:2\}$
cf	$\{(d:2), \emptyset:1\}$



Άλλο ένα παράδειγμα (5/6)

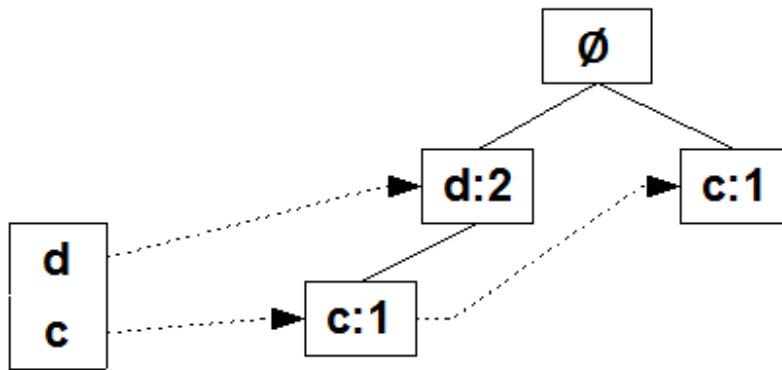
ef	{{dca:2}}
----	-----------

- Συνδυασμός ef με κάθε υποσύνολο του dca (και το κενό):
 - ef, def, cef, aef, dcef, daef, caef, dcaef
 - Όλα με υποστήριξη ίση με 2
 - Δηλ. όταν μένει μόνο ένα υπο συνθήκη μονοπάτι, σταματάμε την αναδρομική διαδικασία.
- (όμοια και για όλα τα άλλα ΥΣ μονοπάτια του f)



Άλλο ένα παράδειγμα (6/6)

b	$\{(dca:1), (d:1), (c:1)\}$
----------	-----------------------------



Επίθεμα	Υ.Σ. Μονοπάτια b
cb	$\{(d:1), (\emptyset:1)\}$
db	$\{\emptyset:2\}$



Μειονεκτήματα υποστήριξης, εμπιστοσύνης

	Κινηματογράφος	Όχι κινηματογράφος	
Θέατρο	15	5	20
Όχι Θέατρο	65	15	80
	80	20	100

- Θέατρο => Κινηματογράφος (15%, 75%)
- $P(\text{Κινηματογράφος}) = 80\% > 75\%$



Μέτρο ενδιαφέροντος

- Έστω ο κανόνας $A \rightarrow B$

—————
()

- \wedge : ανεξαρτησία.
- \vee : θετική συσχέτιση.
- \neg : αρνητική συσχέτιση.



Παράδειγμα

	Κινηματογράφος	Όχι κινηματογράφος	
Θέατρο	15	5	20
Όχι Θέατρο	65	15	80
	80	20	100

- $I = 0.15 / (0.2 * 0.8) = 0.9375 < 1$



Παράδοξο Simpson (1/3)

	Ραδιόφωνο	Όχι ραδιόφωνο	
Τηλεόραση	99	81	180
Όχι Τηλεόραση	54	66	120
	153	147	300

- Τηλεόραση \rightarrow Ραδιόφωνο (εμπ = $99/180 = 55\%$)
- Όχι Τηλεόραση \rightarrow Ραδιόφωνο (εμπ = $54/120 = 45\%$)
- Συσχέτιση(Τηλεόραση, Ραδιόφωνο) = 1.07
- Θετική συσχέτιση μεταξύ τηλεόρασης και ραδιοφώνου.



Παράδοξο Simpson (2/3)

	Ραδιόφωνο	Όχι ραδιόφωνο	
Τηλεόραση	1	9	10
Όχι Τηλεόραση	4	30	34
	5	13	44

ανήλικοι

- Τηλεόραση => Ραδιόφωνο (εμπ = $1/10 = 10\%$)
- Όχι Τηλεόραση => Ραδιόφωνο (εμπ = $4/34 = 11.8\%$)
- $I(\text{Τηλεόραση}, \text{Ραδιόφωνο}) = 0.88$



Παράδοξο Simpson (3/3)

	Ραδιόφωνο	Όχι ραδιόφωνο	
Τηλεόραση	98	72	170
Όχι Τηλεόραση	50	36	86
	148	108	156

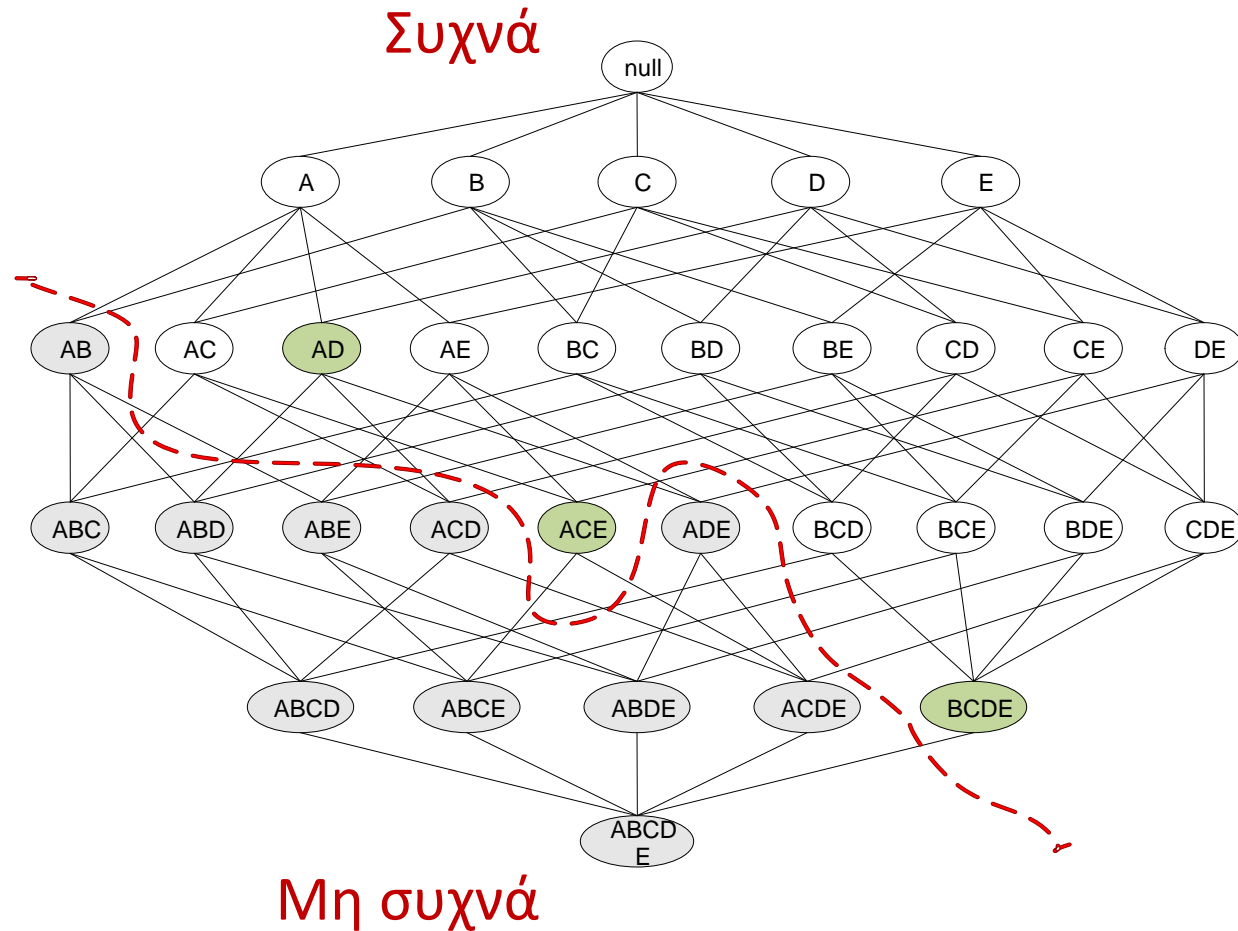
ανήλικοι

- Τηλεόραση => Ραδιόφωνο (εμπ = $98/170 = 57.7\%$)
- Όχι Τηλεόραση => Ραδιόφωνο (εμπ = $50/86 = 58.1\%$)
- $I(\text{Τηλεόραση}, \text{Ραδιόφωνο}) = 0.60$



Maximal συχνά στοιχειοσύνολα

- Ένα στοιχειοσύνολο είναι **maximal συχνό** αν κανένα από τα άμεσα υπερσύνολά του δεν είναι συχνό.
- Προσφέρουν μια συνοπτική αναπαράσταση των συχνών στοιχειοσυνόλων.
- Είναι το μικρότερο σύνολο στοιχειοσυνόλων από το οποίο μπορούμε να πάρουμε όλα τα συχνά στοιχειοσύνολα.
- ΟΜΩΣ: Δεν προσφέρουν καμιά πληροφορία για την υποστήριξη των υποσυνόλων τους.



Κλειστά συχνά στοιχειosύνολα

- Ένα στοιχειosύνολο είναι κλειστό (closed) αν κανένα από τα άμεσα υπερσύνολα του δεν έχει την ίδια υποστήριξη με αυτό (δηλαδή, έχει μικρότερη υποστήριξη).
- Ένα στοιχειosύνολο είναι κλειστό συχνό στοιχειosύνολο 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.
- Η υποστήριξη ενός μη κλειστού στοιχειosυνόλου πρέπει να είναι ίση με την μεγαλύτερη υποστήριξη ανάμεσα στα υπερσύνολά του.



Σημείωμα Αναφοράς

Copyright Αριστοτέλειο Πανεπιστήμιο Θεσσαλονίκης, Αναστάσιος Γούναρης.
«Αποθήκες Δεδομένων και Εξόρυξη Δεδομένων. Ενότητα 12. Κανόνες
Συσχέτισης – Μέρος Β΄». Έκδοση: 1.0. Θεσσαλονίκη 2014.

Διαθέσιμο από τη δικτυακή διεύθυνση:<http://eclass.auth.gr/courses/OCRS182/>



Σημείωμα Αδειοδότησης

Το παρόν υλικό διατίθεται με τους όρους της άδειας χρήσης Creative Commons Αναφορά - Μη Εμπορική Χρήση - Παρόμοια Διανομή 4.0 [1] ή μεταγενέστερη, Διεθνής Έκδοση. Εξαιρούνται τα αυτοτελή έργα τρίτων π.χ. φωτογραφίες, διαγράμματα κ.λ.π., τα οποία εμπεριέχονται σε αυτό και τα οποία αναφέρονται μαζί με τους όρους χρήσης τους στο «Σημείωμα Χρήσης Έργων Τρίτων».



Ο δικαιούχος μπορεί να παρέχει στον αδειοδόχο ξεχωριστή άδεια να χρησιμοποιεί το έργο για εμπορική χρήση, εφόσον αυτό του ζητηθεί.

Ως **Μη Εμπορική** ορίζεται η χρήση:

- που δεν περιλαμβάνει άμεσο ή έμμεσο οικονομικό όφελος από την χρήση του έργου, για το διανομέα του έργου και αδειοδόχο
- που δεν περιλαμβάνει οικονομική συναλλαγή ως προϋπόθεση για τη χρήση ή πρόσβαση στο έργο
- που δεν προσπορίζει στο διανομέα του έργου και αδειοδόχο έμμεσο οικονομικό όφελος (π.χ. διαφημίσεις) από την προβολή του έργου σε διαδικτυακό τόπο

[1] <http://creativecommons.org/licenses/by-nc-sa/4.0/>





Τέλος ενότητας

Επεξεργασία: Ανδρέας Κοσματόπουλος
Θεσσαλονίκη, Χειμερινό Εξάμηνο 2013-2014



Ευρωπαϊκή Ένωση
Ευρωπαϊκό Κοινωνικό Ταμείο



ΥΠΟΥΡΓΕΙΟ ΠΑΙΔΕΙΑΣ ΚΑΙ ΘΡΗΣΚΕΥΜΑΤΩΝ
ΕΙΔΙΚΗ ΥΠΗΡΕΣΙΑ ΔΙΑΧΕΙΡΙΣΗΣ

Με τη συγχρηματοδότηση της Ελλάδας και της Ευρωπαϊκής Ένωσης



ΕΥΡΩΠΑΪΚΟ ΚΟΙΝΩΝΙΚΟ ΤΑΜΕΙΟ



ΑΡΙΣΤΟΤΕΛΕΙΟ
ΠΑΝΕΠΙΣΤΗΜΙΟ
ΘΕΣΣΑΛΟΝΙΚΗΣ

Σημειώματα

Διατήρηση Σημειωμάτων

Οποιαδήποτε αναπαραγωγή ή διασκευή του υλικού θα πρέπει να συμπεριλαμβάνει:

- το Σημείωμα Αναφοράς
- το Σημείωμα Αδειοδότησης
- τη δήλωση Διατήρησης Σημειωμάτων
- το Σημείωμα Χρήσης Έργων Τρίτων (εφόσον υπάρχει)

μαζί με τους συνοδευόμενους υπερσυνδέσμους.

