

Εισαγωγή

Επεξεργασία Ερωτήσεων

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιος Πιτσιλιάς 1

Εισαγωγή

Η Δομή ενός ΣΒΔ

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιος Πιτσιλιάς 2

Εισαγωγή

Γενική Εικόνα του Μαθήματος

ΜΕΡΟΣ 1 (Χρήση ΣΔΒΔ)

1. Μοντελοποίηση (Μοντέλο Ο/Σ, Σχεσιακό, Λογικός Σχεδιασμός)
2. Κανονικοποίηση
3. Προγραμματισμός (Σχεσιακή Άλγεβρα, Σχεσιακό Λογισμό, SQL)

Δημιουργία/Κατασκευή, Εισαγωγή Δεδομένων, Επεξεργασία Δεδομένων

ΜΕΡΟΣ 2 (Υλοποίηση ΣΔΒΔ)

Αποθήκευση
Ευρετήρια

Εκτέλεση Ερωτήσεων

Το εσωτερικό ενός ΣΔΒΔ

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιος Πιτσιλιάς 3

Εισαγωγή

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιος Πιτσιλιάς 4

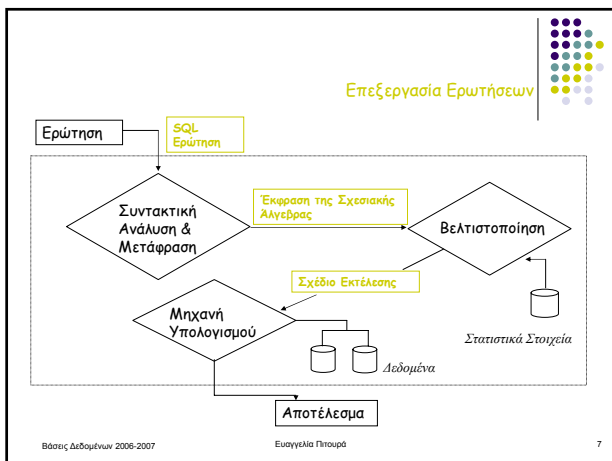
Εισαγωγή

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιος Πιτσιλιάς 5

Επεξεργασία Ερωτήσεων

Θα δούμε την «πορεία» μιας SQL ερώτησης (πως εκτελείται)

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιος Πιτσιλιάς 6



Επεξεργασία Ερωτήσεων

Τα βασικά βήματα στην επεξεργασία μιας ερώτησης είναι

1. Συντακτική Ανάλυση & Μετάφραση
2. Βελτιστοποίηση
3. Υπολογισμός

Βάσεις Δεδομένων 2006-2007 8

Συντακτική Ανάλυση & Μετάφραση

1. Συντακτική Ανάλυση (Parsing) & Μετάφραση

Η SQL ερώτηση μεταφράζεται σε μια εσωτερική μορφή αφού γίνει ο απαραίτητος συντακτικός και σημασιολογικός έλεγχος (π.χ., τα ονόματα που αναφέρονται είναι ονόματα σχέσεων που υπάρχουν)

Αντικατάσταση των όψεων από τον ορισμό τους

Σε ποια εσωτερική μορφή; Έκφραση της σχεσιακής άλγεβρας

select A_1, A_2, \dots, A_n
from R_1, R_2, \dots, R_m $\pi_{A_1, A_2, \dots, A_n} (\sigma_P (R_1 \times R_2 \times \dots \times R_m))$
where P

Βάσεις Δεδομένων 2006-2007 9

Βελτιστοποίηση

2. Βελτιστοποίηση

Μια SQL ερώτηση μπορεί να μεταφραστεί σε διαφορετικές (ισοδύναμες) εκφράσεις της σχεσιακής άλγεβρας

select balance • $\pi_{balance} (\sigma_{balance < 2500} (account))$
from account
where balance < 25000 • $\sigma_{balance < 2500} (\pi_{balance}(account))$

Με ποιο κριτήριο γίνεται η επιλογή της έκφρασης;

Βάσεις Δεδομένων 2006-2007 10

Βελτιστοποίηση

Κάθε πράξη της σχεσιακής άλγεβρας μπορεί να υλοποιηθεί με διαφορετικούς αλγόριθμους.

π.χ., για την υλοποίηση της επιλογής μπορεί είτε να σαρώσουμε (scan) όλο το αρχείο ελέγχοντας κάθε εγγραφή αν ικανοποιεί τη συνθήκη είτε αν υπάρχει π.χ., ένα Β* ευρετήριο στο γνώρισμα balance να χρησιμοποιήσουμε το ευρετήριο

Άρα δεν αρκεί ο προσδιορισμός της πράξης - πρέπει να προσδιορίζεται και ο αλγόριθμος που θα χρησιμοποιηθεί για την υλοποίησή της

Βάσεις Δεδομένων 2006-2007 11

Βελτιστοποίηση

βασικές (primitive) πράξεις: πράξη + αλγόριθμος

Σχέδιο εκτέλεσης (execution plan): μια ακολουθία από βασικές πράξεις

$\pi_{balance}$
 $\sigma_{balance < 2500}$, χρησιμοποίησε το ευρετήριο 1
 account

Βάσεις Δεδομένων 2006-2007 12



- Τα διαφορετικά σχέδια εκτέλεσης έχουν και διαφορεικό κόστος
- **Βελτιστοποίηση:** η διαδικασία επιλογής του σχεδίου εκτέλεσης που έχει το μικρότερο κόστος
- Εκτίμηση του κόστους (συνήθως χρήση στατιστικών στοιχείων)



3. Εκτέλεση

Μηχανή εκτέλεσης που εκτελεί τις βασικές πράξεις

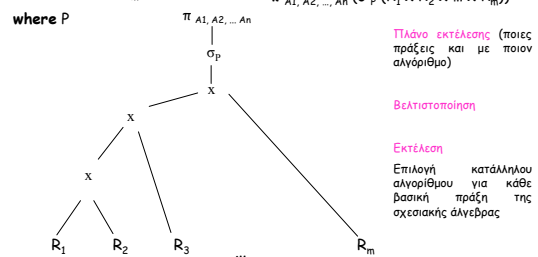


Υπάρχουν υλοποιημένοι μια σειρά από αλγόριθμοι για κάθε βασική πράξη (π.χ., που χρησιμοποιούν ή όχι ευρετήρια κλπ)

Γενικά, το ΣΔΒΔ με βάση κάποια *στατιστικά στοιχεία* κάνει μια *εκτίμηση του κόστους* και *επιλέγει τον αλγόριθμο* με τον μικρότερο (με βάση την εκτίμηση) κόστος



$\text{select } A_1, A_2, \dots, A_n$ **Μετάφραση** $\pi_{A_1, A_2, \dots, A_n}(\sigma_P(R_1 \times R_2 \times \dots \times R_m))$
 $\text{from } R_1, R_2, \dots, R_m$
 $\text{where } P$



Πλάνο εκτέλεσης (ποιες πράξεις και με ποιον αλγόριθμο)

Βελτιστοποίηση

Εκτέλεση

Επιλογή κατάλληλου αλγορίθμου για κάθε βασική πράξη της σχεσιακής άλγεβρας



Τι είδους στατιστικά στοιχεία;

Για ένα αρχείο δεδομένων μιας σχέσης R:

- n_R : αριθμός πλειάδων της σχέσης R
- b_R : αριθμός blocks της σχέσης R
- s_R : μέγεθος σε bytes κάθε πλειάδας της σχέσης R
- f_R : παράγοντας ομαδοποίησης (αριθμός εγγραφών ανά block) αν μη εκτεινόμενη, $f_R = \lfloor B / s_R \rfloor$ και $b_R = \lceil n_R / f_R \rceil$

Ενημέρωση στατιστικών στοιχείων;



Άλλα στατιστικά στοιχεία;

Π.χ., για μια πράξη επιλογής στο γνώρισμα A

- $V(A, R)$: αριθμός διαφορετικών τιμών του A
 $|\pi_A(R)|$ -- αν το A κλειδί;
- $SC(A, R)$: μέσος αριθμός πλειάδων που ικανοποιεί μια συνθήκη (δεδομένου ότι υπάρχει μια τουλάχιστον που την ικανοποιεί)
1 αν κλειδί, αν ομοιόμορφη;

Στατιστικά στοιχεία επίσης για το *αρχείο ευρετηρίου* (αν υπάρχει)

- f_i : παράγοντας διακλάδωσης, πολυεπίπεδο f_0 , Β: δέντρο ~ τάξη
- H_i : αριθμός επιπέδων
- LB_i : αριθμός block φύλλων

Με βάση τα στατιστικά επιλέγεται ο αλγόριθμος με το μικρότερο κόστος
Κόστος: Αριθμό blocks που μεταφέρονται

Πιθανοί *αλγόριθμοι εκτέλεσης* για την επιλογή:

- E1: Σειριακή αναζήτηση
- E2: Δυαδική αναζήτηση
- E3: Χρήση πρωτεύοντος ευρετηρίου/κατακερματισμού
- E4: Χρήση δευτερεύοντος ευρετηρίου/κατακερματισμού

Για τον E2 πρέπει το αρχείο να είναι ταξινομημένο

Για τους E3 και E4 λέμε ότι έχουμε μονοπάτι προσπέλασης (access path)

Τι θα καλύψουμε στο μάθημα:

1. Αλγόριθμοι εκτέλεσης βασικών πράξεων

- Επιλογή
- Προβολή
- Πράξεις συνόλων
- Συνένωση

2. Βασικές Ιδέες για Βελτιστοποίηση

Επιλογή

Θα εξετάσουμε:

- Επιλογή με συνθήκη ισότητας ($\sigma_{A = a} (R)$)
- Επιλογή με συνθήκη σύγκρισης - διαστήματος/περιοχής (range query) ($\sigma_{A \leq u} (R)$) ή ($\sigma_{A \geq u} (R)$)
- Επιλογή με σύζευξη ($\sigma_{\theta_1 \text{ AND } \theta_2 \dots \text{ AND } \theta_n} (R)$)
- Επιλογή με διάζευξη ($\sigma_{\theta_1 \text{ OR } \theta_2 \dots \text{ OR } \theta_n} (R)$)

Θα δούμε ποιος είναι ο καλύτερος (με το μικρότερο κόστος σε blocks) αλγόριθμος για την εκτέλεση της πράξης

Επιλεκτικότητα επιλογής:

το πλήθος των εγγραφών (πλειάδων) που επιλέγονται (δηλ. ικανοποιούν την συνθήκη) το πλήθος των εγγραφών (πλειάδων) του αρχείου (σχέσης)

• Έστω $s_i = | \sigma_{\theta_i} (R) |$

επιλεκτικότητα: s_i / n_R

Αν θ_i **συνθήκη ισότητας** σε ένα γνώρισμα υποψήφιο κλειδί $s_i = 1 / n_R$

Αν θ_i **συνθήκη ισότητας** σε ένα γνώρισμα, ομοιόμορφη κατανομή, k διακριτές τιμές, $s_i = k / n_R$

Επιλογή - συνθήκη ισότητας $\sigma_{A = a} (R)$

E1 Σειριακή αναζήτηση

Διάβασμα (scan) όλου του αρχείου

b_R : αριθμός blocks της σχέσης R

b_R

$b_R/2$ αν το A υποψήφιο κλειδί (οπότε το αποτέλεσμα έχει μόνο μία πλειάδα, σταματάμε την αναζήτηση μόλις τη βρούμε)

Μπορεί να χρησιμοποιηθεί σε οποιοδήποτε αρχείο

Επιλογή: Συνθήκη Ισότητας

b_R : αριθμός blocks της σχέσης R
 $SC(A, R)$: μέσος αριθμός πλειάδων που ικανοποιεί μια συνθήκη
 f_R : παράγοντας ομαδοποίησης

E2 Διαδική αναζήτηση

Μπορεί να χρησιμοποιηθεί μόνο αν το αρχείο είναι διατεταγμένο με βάση το γνώρισμα της επιλογής

$$\lceil \log(b_R) \rceil \leftarrow \text{Εύρεση της πρώτης}$$

$$+ \lceil SC(A, r)/f_R \rceil - 1 \leftarrow \text{Εύρεση των υπόλοιπων}$$

Αν το A υποψήφιο κλειδί:

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιο Πλιουρά 25

Επιλογή: Συνθήκη Ισότητας

b_R : αριθμός blocks της σχέσης R
 $SC(A, R)$: μέσος αριθμός πλειάδων που ικανοποιεί μια συνθήκη
 f_R : παράγοντας ομαδοποίησης
 HT : αριθμός επιπέδων

E3 Χρήση πρωτεύοντος (πολυεπίπεδου) ευρετηρίου

Πρωτεύον ευρετήριο σημαίνει ταξινομημένο αρχείο

Μπορεί να χρησιμοποιηθεί μόνο αν υπάρχει τέτοιο ευρετήριο στο A

$$HT_i + 1 \leftarrow \text{Εύρεση και μεταφορά της πρώτης}$$

Αν το A δεν είναι υποψήφιο κλειδί -- ευρετήριου συστάδων

$$HT_i + \lceil SC(A, R)/f_R \rceil \leftarrow \text{Εύρεση και των υπόλοιπων}$$

ΣΗΜΕΙΩΣΗ: Πρωτεύον ευρετήριο στο A, σημαίνει ότι οι εγγραφές του αρχείου δεδομένων είναι ταξινομημένες (διατεταγμένες) ως προς A άρα οι υπόλοιπες εγγραφές με την ίδια τιμή (αν υπάρχουν) βρίσκονται σε γειτονικά blocks του αρχείου δεδομένων

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιο Πλιουρά 26

Επιλογή: Συνθήκη Ισότητας

b_R : αριθμός blocks της σχέσης R
 $SC(A, R)$: μέσος αριθμός πλειάδων που ικανοποιεί μια συνθήκη
 f_R : παράγοντας ομαδοποίησης
 HT : αριθμός επιπέδων

E4 Χρήση δευτερεύοντος (πολυεπίπεδου) ευρετηρίου

Μπορεί να χρησιμοποιηθεί μόνο αν υπάρχει τέτοιο ευρετήριο στο A

Αν το A είναι υποψήφιο κλειδί

$$HT_i + 1 \leftarrow \text{Εύρεση και μεταφορά της πρώτης}$$

Αν το A δεν είναι υποψήφιο κλειδί ± κόστος για την εύρεση των υπολοίπων

$$HT_i + \text{ενδιάμεσο επίπεδο} + SC(A, R) \leftarrow \text{Εύρεση και των υπόλοιπων}$$

Στη χειρότερη περίπτωση κάθε εγγραφή που ικανοποιεί τη συνθήκη σε διαφορετικό block

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιο Πλιουρά 27

Επιλογή: Συνθήκη με Σύγκριση

Επιλογή - συνθήκη με σύγκριση

$$\sigma_{A \leq u} (R) \text{ ή } \sigma_{A \geq u} (R)$$

Έστω ότι c πλειάδες ικανοποιούν τη συνθήκη

Γενικά $c = n_R/2$ (δηλαδή, οι μισές)

Έστω min, max (μικρότερη, μεγαλύτερη τιμή του A), αν ομοιόμορφη κατανομή και $\sigma_{A \leq u} (R)$

$$c = \begin{cases} 0 & \text{αν } u < \min \\ n_R & \text{αν } u \geq \max \\ n_R * [(u - \min) / (\max - \min)] & \end{cases}$$

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιο Πλιουρά 28

Επιλογή: Συνθήκη με Σύγκριση

$$\sigma_{A \leq u} (R)$$

Θα δούμε δύο αλγόριθμους:

E5 Χρήση πρωτεύοντος πολυ-επίπεδου ευρετηρίου

E6 Χρήση δευτερεύοντος πολυ-επίπεδου ευρετηρίου

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιο Πλιουρά 29

Επιλογή: Συνθήκη με Σύγκριση

E5 Χρήση πρωτεύοντος (πολυεπίπεδου) ευρετηρίου

$A \geq u$

Πρωτεύον, σημαίνει ταξινομημένο αρχείο, έτσι σε αύξουσα διάταξη

- Χρήση ευρετηρίου για την εύρεση της πρώτης εγγραφής $A \geq u$
- Σάρωση όλου του αρχείου ξεκινώντας από αυτήν την εγγραφή

$$HT_i + \lceil c / f_R \rceil$$

$A \leq u$

Δε χρειάζεται ευρετήριο, γιατί:

c : επιλεξιμότητα (πλειάδες που ικανοποιούν την συνθήκη)
 f_R : παράγοντας ομαδοποίησης
 HT : αριθμός επιπέδων

Βάσεις Δεδομένων 2006-2007 Ευαγγέλιο Πλιουρά 30



E6 Χρήση δευτερεύοντος (πολυεπίπεδου) ευρετηρίου

- Εύρεση του πρώτου φύλλου του ευρετηρίου *Μη ταξινομημένο αρχείο*
- Για κάθε block (φύλλο) του ευρετηρίου διάβασε το αντίστοιχο block δεδομένων *(σημείωση, αν B+ δέντρο χρησιμοποιούμε το δείκτη ανάμεσα στα φύλλα)*

Σάρωση των φύλλων του δέντρου

$A \leq u$ από την αρχή έως το u

$A \geq u$ από το u έως το τέλος

*c: επιλεξιμότητα (πλειάδες που ικανοποιούν την συνθήκη)
 n_R : αριθμός εγγραφών
 LB : αριθμός block φύλλων
 HT : αριθμός επιπέδων*

Π.χ., αν $c = n_R/2$, και A κλειδί τότε (αν κάθε εγγραφή σε διαφορετικό block)

$$HT_i + LB_i/2 + n_R/2$$



Επιλογή - συνθήκη σύζευξης

$$\sigma \theta_1 \text{ AND } \theta_2 \dots \text{ AND } \theta_n (R)$$

Επιλεκτικότητα μιας συνθήκης:

το πλήθος των εγγραφών (πλειάδων) που την ικανοποιούν
το πλήθος των εγγραφών (πλειάδων) του αρχείου (σχέσης)

- Αν οι συνθήκες είναι ανεξάρτητες, το μέγεθος του αποτελέσματος:

$$\frac{n_R * s_1 * s_2 * \dots * s_n}{n_R^n}$$

s_i επιλεκτικότητα της θ_i



E7 Συζευκτική επιλογή με χρήση ενός απλού ευρετηρίου

Υπάρχει διαδρομή προσπέλασης για ένα από τα γνωρίσματα που εμφανίζονται σε οποιαδήποτε απλή συνθήκη

Επιλογή του γνωρίσματος στην απλή συνθήκη με τη *μικρότερη* επιλεκτικότητα (γιατί:)

Χρήση μιας από τις προηγούμενες μεθόδους για την ανάκτηση των εγγραφών που ικανοποιούν αυτήν την συνθήκη και

έλεγχος για κάθε επιλεγμένη εγγραφή αν ικανοποιεί και τις υπόλοιπες συνθήκες



E8 Συζευκτική επιλογή με χρήση σύνθετου ευρετηρίου

Αν υπάρχει ευρετήριο στο συνδυασμό δύο ή περισσότερων γνωρισμάτων που εμφανίζονται σε οποιαδήποτε απλές συνθήκες



E9 Συζευκτική επιλογή με τομή δεικτών

Αν υπάρχουν ευρετήρια σε περισσότερα από ένα από τα γνωρίσματα

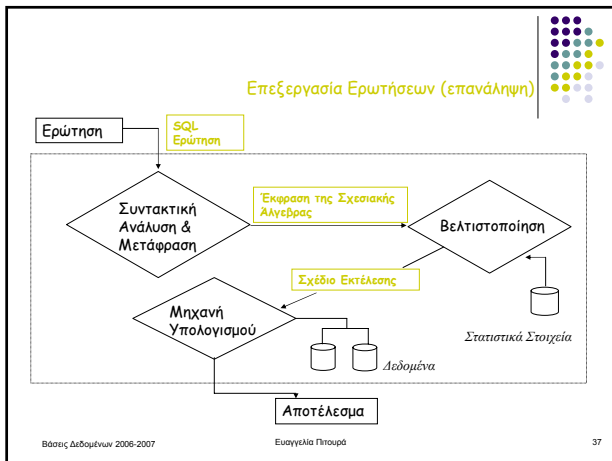
Τότε διαβάζουμε τα blocks του αρχείου δεδομένων που δίνονται από όλα τα ευρετήρια



Επιλογή - συνθήκη διάζευξης

$$\sigma \theta_1 \text{ OR } \theta_2 \dots \text{ OR } \theta_n (R)$$

Αν έστω και μία από τις συνθήκες δεν έχει διαδρομή προσπέλασης -> σάρωση όλου του αρχείου



Επεξεργασία Ερωτήσεων

Τι θα καλύψουμε στο μάθημα:

- Αλγόριθμους εκτέλεσης βασικών πράξεων
 - Επιλογή
 - Πράξεις συνόλων
 - Συνένωση
- Βασικές Ιδέες για Βελτιστοποίηση

Βάσεις Δεδομένων 2006-2007 Ευαγγέλια Πτιουρά 38

Συνένωση

Συνένωση

$$R \bowtie R.A \text{ op } S.B \quad S$$

Σ1 Εμφωλευμένος (εσωτερικός - εξωτερικός) βρόγχος
 Σ2 Χρήση μιας δομής προσπέλασης
 Σ3 Ταξινόμηση-Συγχώνευση
 Σ4 Συνένωση με κατακερματισμό

Έχει σημασία πόσο χώρο μνήμης κάθε χρονική στιγμή (buffers) μπορούμε να χρησιμοποιήσουμε για τις σχέσεις
 Αρχικά, ας υποθέσουμε ότι έχουμε **μόνο 2 blocks**

Βάσεις Δεδομένων 2006-2007 Ευαγγέλια Πτιουρά 39

Συνένωση

Επιλεκτικότητα συνένωσης μιας σχέσης:

$$\frac{\text{το πλήθος των εγγραφών (πλειάδων) που επιλέγονται}}{\text{το πλήθος των εγγραφών (πλειάδων) του αρχείου (σχέσης)}}$$

- Σε ορισμένες περιπτώσεις μπορεί να δημιουργηθεί ένα ευρετήριο ειδικά για τη συνένωση

Βάσεις Δεδομένων 2006-2007 Ευαγγέλια Πτιουρά 40

Συνένωση (εμφωλευμένος βρόγχος)

Σ1 Εμφωλευμένος (εσωτερικός-εξωτερικός) βρόγχος

Για κάθε εγγραφή t της R
 Για κάθε εγγραφή s της S
 Αν $t[A]$ op $s[B]$ πρόσθεσε το t s στο αποτέλεσμα

Αγνοώντας την εγγραφή των blocks του αποτελέσματος

$$b_r + n_r * b_s$$

Βάσεις Δεδομένων 2006-2007 Ευαγγέλια Πτιουρά 41

Συνένωση (εμφωλευμένος βρόγχος)

Για κάθε block B_r της R
 Για κάθε block B_s της S
 Για κάθε εγγραφή t του B_r
 Για κάθε εγγραφή s του B_s
 Αν $t[A]$ op $s[B]$ πρόσθεσε το t s στο αποτέλεσμα

Αγνοώντας την εγγραφή των blocks του αποτελέσματος

$$b_r + b_r * b_s$$

Συμφέρει η τοποθέτηση της μικρότερης σχέσης στον εξωτερικό βρόγχο

Βάσεις Δεδομένων 2006-2007 Ευαγγέλια Πτιουρά 42

Συνένωση (εμφωλευμένος βρόγχος)

Πριν θεωρήσουμε ότι έχουμε 2 block στη μνήμη (buffers) διαθέσιμους
 Αν υπάρχουν $n_B > 2$ blocks στη μνήμη που μπορεί να χρησιμοποιηθούν για τον υπολογισμό της συνένωσης συμφέρει να διαβάζουμε τα blocks της σχέσης του εξωτερικού βρόγχου ανά $n_B - 1$

Για κάθε $n_B - 1$ block B_r της R
 Για κάθε block B_s της S
 Για κάθε εγγραφή t του B_r
 Για κάθε εγγραφή s του B_s
 Αν $t[A]$ or $s[B]$ πρόσθεσε το t s στο αποτέλεσμα

$$b_R + \lceil (b_R / (n_B - 1)) \rceil * b_S$$

Συνένωση (χρήση ευρετηρίου)

Σ2 Χρήση μιας δομής προσπέλασης

Η σχέση για την οποία υπάρχει ευρετήριο τοποθετείται στον εσωτερικό βρόγχο. Έστω ότι υπάρχει ευρετήριο για το γνώρισμα B της σχέσης S

Για κάθε block B_r της R
 Για κάθε εγγραφή t του B_r
 Χρησιμοποίησε το ευρετήριο στο B για να βρεις τις εγγραφές s της S τέτοιες ώστε $t[A]$ or $s[B]$

$b_R + n_R * C$ όπου C το κόστος μιας επιλογής στο S (δηλαδή της εύρεσης της εγγραφής (εγγραφών) του S που ικανοποιούν τη συνθήκη)

Συνένωση (ταξινόμηση-συγχώνευση)

Σ3 Ταξινόμηση - Συγχώνευση

Ταξινόμησε τις πλειάδες της R στο γνώρισμα A
 Ταξινόμησε τις πλειάδες της S στο γνώρισμα B
 $i := 1; j := 1;$
 while ($i \leq n_R$ and $j \leq n_S$)
 if ($R_i[A] < S_j[B]$)
 $i := i + 1$; (*προχώρησε το δείκτη στην R *)
 if ($R_i[A] > S_j[B]$)
 $j := j + 1$; (* προχώρησε το δείκτη στην S *)

Συνένωση (ταξινόμηση-συγχώνευση)

else (* $R_i[A] = S_j[B]$ *)
 πρόσθεσε το R_i, S_j στο αποτέλεσμα
 $k := j + 1$; (* γράψε και τις άλλες πλειάδες της S που ταιριάζουν, αν υπάρχουν *)
 while (($k \leq n_S$) and ($R_i[A] = S_k[B]$))
 πρόσθεσε το R_i, S_k στο αποτέλεσμα
 $k := k + 1$;
 $m := i + 1$; (* γράψε και τις άλλες πλειάδες της R που ταιριάζουν, αν υπάρχουν *)
 while (($m \leq n_R$) and ($R_m[A] = S_j[B]$))
 πρόσθεσε το R_m, S_j στο αποτέλεσμα
 $m := m + 1$;
 $i := m; j := k;$

Συνένωση (ταξινόμηση-συγχώνευση)

Αν αγνοήσουμε τη ταξινόμηση για τη συγχώνευση (merge) απλή σάρωση των δύο αρχείων:

$$b_R + b_S$$

Ταξινόμηση: $b_R * \log(b_R) + b_S * \log(b_S)$

Συνένωση (με κατακερματισμό)

Σ4 Συνένωση με κατακερματισμό

- χωρίζουμε με βάση μια συνάρτηση κατακερματισμού h τις πλειάδες της S και της R σε κάδους -- στον ίδιο κάδο αν $h(t_R[A]) = h(t_S[B])$
- δηλαδή οι πλειάδες με $t_R[A] = t_S[B]$ πέφτουν στον ίδιο κάδο άρα αρκεί να ελέγξουμε μεταξύ τους τις πλειάδες που πέφτουν στον ίδιο κάδο

Συνένωση (με κατακερματισμό)

Κατακερμάτισε τις εγγραφές της R χρησιμοποιώντας την $h(t_R[A])$

Για κάθε εγγραφή t_S της S

$k := h(t_S[B])$

σύγκρινε το $t_S[B]$ με $t_{R_i}[A]$ για όλες τις εγγραφές t_{R_i} της R στον κάδο k

- Χρησιμοποιούμε την μικρότερη σχέση για το πρώτο πέρασμα.
- Αν όλοι οι κάδοι που προκύπτουν χωράνε στη μνήμη, κόστος $b_R + b_S$

Συνένωση (με κατακερματισμό)

Αν δεν χωρούν όλοι οι κάδοι - τροποποίηση

Πράξεις Συνόλων

Ταξινομήσε τις πλειάδες της R σε ένα γνώρισμα (έστω A)

Ταξινομήσε τις πλειάδες της S στο ίδιο γνώρισμα

$i := 1; j := 1;$

while ($i \leq n_R$ and $j \leq n_S$)

if ($R_i[A] > S_j[A]$)

Τομή

τίποτα

Ένωση

γράψε το S_j στο αποτέλεσμα

Διαφορά

τίποτα

$j := j + 1$

Πράξεις Συνόλων

else if ($R_i[A] < S_j[A]$)

Τομή

τίποτα

Ένωση

γράψε το R_i στο αποτέλεσμα

Διαφορά

γράψε το R_i στο αποτέλεσμα

$i := i + 1$

else ($* R_i[A] = S_j[A] *$)

Τομή

γράψε το R_i στο αποτέλεσμα

$i := i + 1;$

$j := j + 1;$

Ένωση

$i := i + 1;$

Διαφορά

$i := i + 1;$

$j := j + 1;$

Πράξεις Συνόλων

Αν υπάρχουν ακόμα εγγραφές για κάποιο αρχείο:

Ένωση

while ($i \leq n_R$)

γράψε το R_i στο αποτέλεσμα

$i := i + 1;$

while ($j \leq n_S$)

γράψε το S_j στο αποτέλεσμα

$j := j + 1;$

Διαφορά

while ($i \leq n_R$)

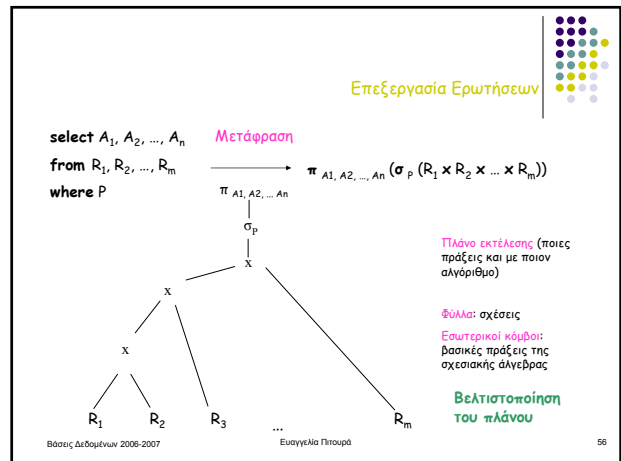
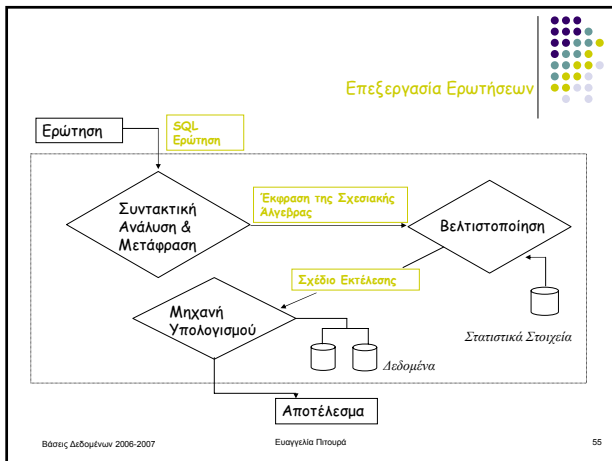
γράψε το R_i στο αποτέλεσμα

$i := i + 1;$

Επεξεργασία Ερωτήσεων

Τι θα καλύψουμε στο μάθημα:

- Βασικές Ιδέες για Βελτιστοποίηση
- Αλγόριθμους εκτέλεσης βασικών πράξεων
 - Επιλογή
 - Προβολή
 - Πράξεις συνόλων
 - Συνένωση



Βελτιστοποίηση Ερωτήσεων

Δέντρο ερώτησης

Φύλλα: σχέσεις εισόδου

Εσωτερικοί κόμβοι: πράξεις της σχεσιακής άλγεβρας

Εκτέλεση δέντρου ερώτησης - χρήση των αλγορίθμων που περιγράψαμε

Βάσεις Δεδομένων 2006-2007 Ευαγγέλια Πλουρά 57

Βελτιστοποίηση Ερωτήσεων

Μερικοί ευριστικοί κανόνες

Γενική ιδέα: εκτέλεση πρώτα των πράξεων με μικρή επιλεκτικότητα ώστε να περιοριστεί το μέγεθος των ενδιάμεσων αποτελεσμάτων

1. Διάσπαση των πράξεων επιλογής με συζευκτικές συνθήκες σε ακολουθίες πράξεων επιλογής
2. Μετατοπίζουμε την πράξη επιλογής όσο πιο κάτω επιτρέπεται από τα γνωρίσματα που περιλαμβάνονται στη συνθήκη
3. Επαναδιευθέτηση των φύλλων ώστε να εκτελούνται πρώτα οι σχέσεις που έχουν τις πιο περιοριστικές πράξεις επιλογής

Βάσεις Δεδομένων 2006-2007 Ευαγγέλια Πλουρά 58

Βελτιστοποίηση Ερωτήσεων

4. Συνδυασμός μιας πράξης καρτεσιανού γινομένου με μια πράξη επιλογής που ακολουθεί
5. Διάσπαση και μετακίνηση των λιστών προβολής όσο πιο κάτω γίνεται στο δέντρο
6. Εντοπισμός υποδέντρων με ομάδες πράξεων που μπορεί να εκτελεστούν με κοινό αλγόριθμο

Βάσεις Δεδομένων 2006-2007 Ευαγγέλια Πλουρά 59

Τέλος!

Βάσεις Δεδομένων 2006-2007 Ευαγγέλια Πλουρά 60