

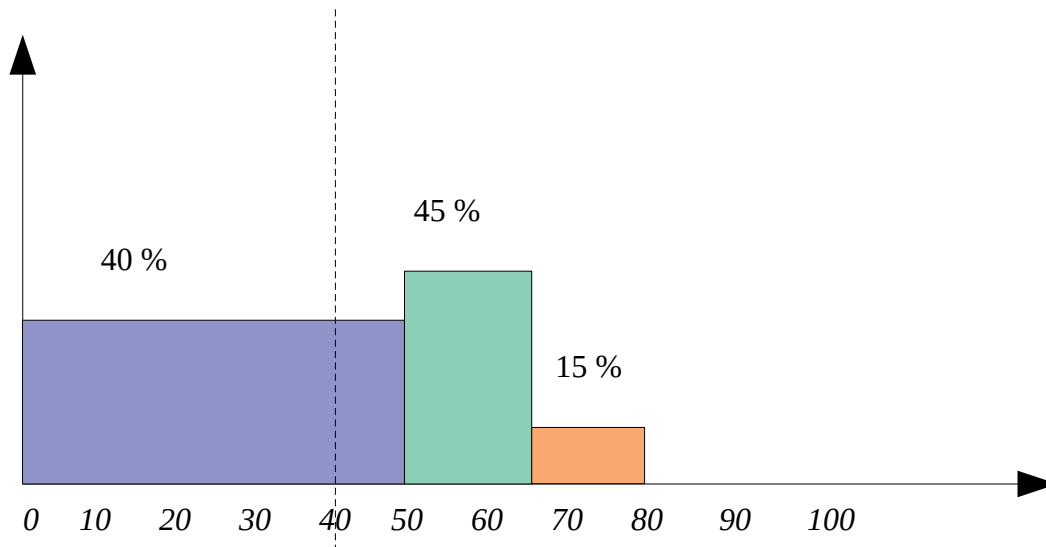
Προχωρημένα Θέματα Τεχνολογίας & Εφαρμογών Βάσεων Δεδομένων

Άσκηση 3

Όθωνας Γκαβαρδίνας ΑΜ: 2620

ΘΕΜΑ 1

(α)



-Απ' το ιστόγραμμα και την ερώτηση ισχύει:

$$\begin{aligned}\text{Ποσοστό} &= (50 - 40 / 50 - 0) * 40 \% + 45\% + 15 \% = \\ &= 1/5 * 40\% + 60\% \\ &= 8\% + 60\% = \underline{68\%}\end{aligned}$$

Κάθε πλειάδα της R αποτελεί το 1/30 μιας σελίδας.

Κάθε πλειάδα της S αποτελεί το 1/20 μιας σελίδας.

Συνολικά μια πλειάδα που διαθέτει τα στοιχεία και των δύο αποτελεί:

$$1/30 + 1/20 = (20 + 30) / 600 = 5 / 60 = \underline{1/12 \text{ μιας σελίδας στο δίσκο}}$$

1.

Το B είναι κλειδί, οπότε συμμετέχουν 100 διαφορετικές τιμές της S.

Το A είναι ξένο κλειδί στο B, οπότε, η R θα έχει μόνο τιμές του S.

Στο αποτέλεσμα, θα περιλαμβάνονται όλες οι πλειάδες της R.

Συνεπώς $150.000 * 100 = 15.000.000$ πλειάδες.

Δεν υπάρχουν εδώ ελάχιστα και μέγιστα. Μόνο μία περίπτωση.

Στην περίπτωση αυτή:

Μέγεθος αποτελέσματος:

$$15.000.000 * 68\% = 10200000 \text{ πλειάδες.}$$

Χώρος στο δίσκο:

$$10200000 * 1/12 = 850000 \text{ σελίδες.}$$

2.

Το B είναι κλειδί, οπότε συμμετέχουν 100 διαφορετικές τιμές της S.

-Ελάχιστο: Η R, στο A μπορεί να μην περιλαμβάνει καμία από αυτές τις τιμές.

Τότε το αποτέλεσμα θα έχει 0 πλειάδες και καθόλου χώρο στο δίσκο.

-Μέγιστο: Η R στο A μπορεί να περιλαμβάνει στη χειρότερη περίπτωση τις 100 διαφορετικές τιμές διασκορπισμένες σε όλες τις πλειάδες της. Τότε τα ταιριάσματα που θα γίνουν θα φέρουν ως αποτέλεσμα, $150.000 * 100 = 15.000.000$ πλειάδες.

Τότε το αποτέλεσμα θα είναι αυτό της 1. περίπτωσης.

3.

Ελάχιστο: Το A και το B μπορεί να μην έχουν κανένα κοινό.

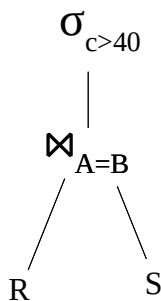
Τότε το αποτέλεσμα θα έχει 0 πλειάδες και καθόλου χώρο στο δίσκο.

Μέγιστο: Το A και το B μπορεί να έχουν όλες τις τιμές κοινές. Τότε $150.000 * 100 = 15.000.000$ πλειάδες.

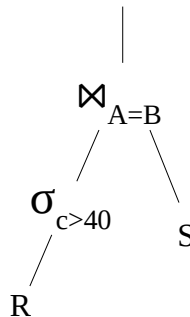
Τότε το αποτέλεσμα θα είναι αυτό της 1. περίπτωσης.

(β)

1.



2.



Στο πλάνο 1, μπορεί να γίνει pipelining, επειδή είναι αριστεροβαθές δένδρο, και οι πλειάδες μπορούν να εξαχθούν και από την R, και από την S ταυτόχρονα.

Στο πλάνο 2, δεν μπορεί να γίνει άμεσα pipelining, μπορεί όμως να εφαρμοστεί πρώτα το φίλτρο της επιλογής σε αυτό και στη συνέχεια να γίνει pipelining στα αποτελέσματα του φίλτρου με τις πλειάδες του S.

Επιλέγω το πλάνο 2.

(γ)

Με βάση το (α) ερώτημα, αν υποθέσουμε πως η γενική περίπτωση της άσκησης είναι το B να είναι πρωτεύον κλειδί και το A ξένο κλειδί, έχουμε range query με:

Selectivity: $15.000.000 * 68\% = 10200000$ πλειάδες.

Κόστος I/O:

Για κάθε μία περίπτωση εφαρμόζω την αντίστοιχη φόρμουλα.

M: σελίδες R ($150000 / 30 = 5000$)

m: πλειάδες R (150000)

N: σελίδες S ($100 / 20 = 5$)

n: πλειάδες S (100)

1. Κόστος = $M + M * N = 5000 + 5000 * 100 = 505000$ I/O
2. Κόστος = $M + \lceil M/B \rceil * N = 5000 + \lceil 150000/52 \rceil * 100 = 5000 + 288500 = 293500$ I/O
3. θεωρείται τυπικά, Κόστος = $M + N = 5000 + 5 = 5005$ I/O
ενώ μπορεί να φτάσει έως $M * N = 5000 * 5 = 25000$ I/O
4. Κόστος = $3 * (M+N) = 3 * 5005 = 15015$ I/O

(δ)

Πρέπει να ισχύουν οι εξής συνθήκες:

1. Πρέπει να σπρωχθούν προς τα φύλλα του δέντρου του πλάνου εκτέλεσης οι πράξεις επιλογής.
2. Πρέπει να τοποθετηθούν οι πράξεις προβολής από το βελτιστοποιητή σε κάθε τελεστή, έτσι ώστε να περνούν μόνο τα γνωρίσματα που χρειάζονται.
3. Δεν πρέπει να γίνεται χρήση καρτεσιανού γινομένου.
4. Το πλάνο να είναι ένα αριστεροβαθές δένδρο.

Το πλάνο εκτέλεσης που έχω επιλέξει στο (β) ερώτημα ικανοποιεί αυτές τις συνθήκες, και επομένως διαθέτει το ελάχιστο δυνατό κόστος, που έχει υπολογισθεί στο ερώτημα (γ).

ΘΕΜΑ 2

1.

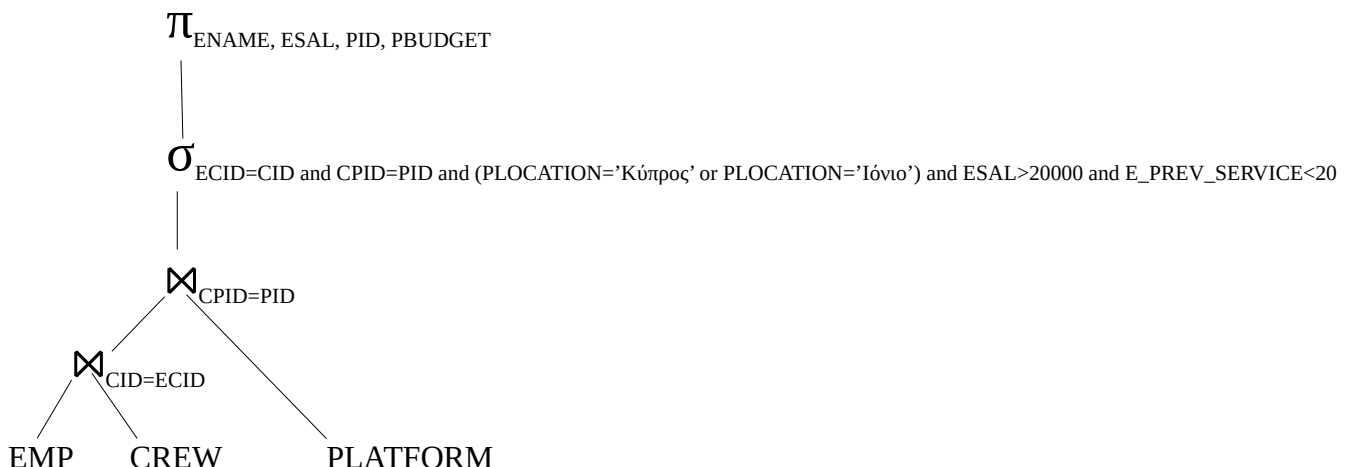
```
SELECT ENAME, ESAL, PID, PBUDGET
FROM EMP, CREW, PLATFORM
WHERE ECID=CID and CPID=PID
and (PLOCATION='Κύπρος' or PLOCATION='Ιόνιο')
and ESAL>20000 and E_PREV_SERVICE<20;
```

$\pi_{ENAME, ESAL, PID, PBUDGET} (\sigma_{PLOCATION='Κύπρος' \text{ or } PLOCATION='Ιόνιο'} (PLATFORM)) ..$

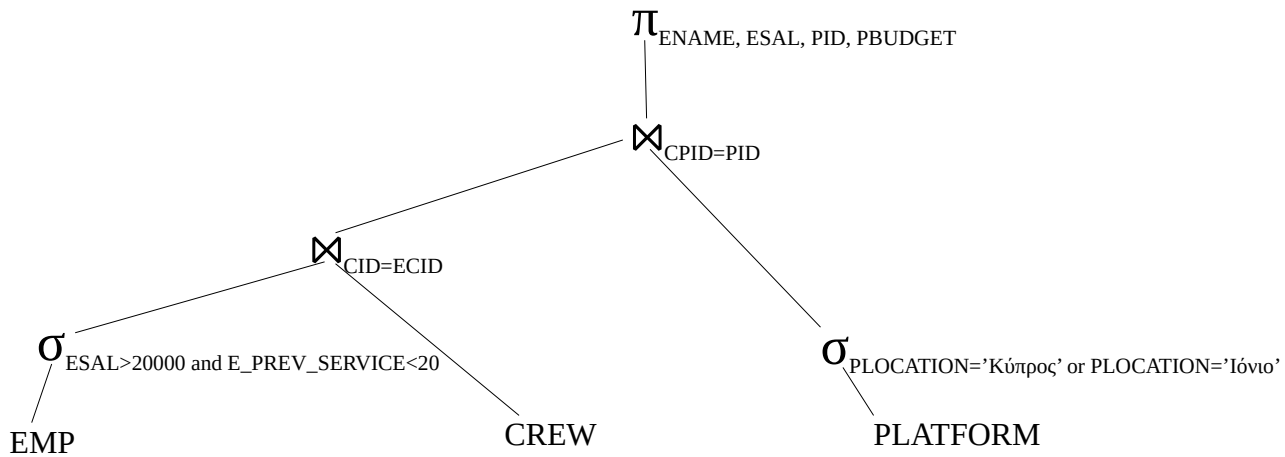
$.. \bowtie_{PID=CPID} (CREW) \bowtie_{CID=ECID} \sigma_{ESAL>20000 \text{ and } E_PREV_SERVICE<20} (EMP))$

2.

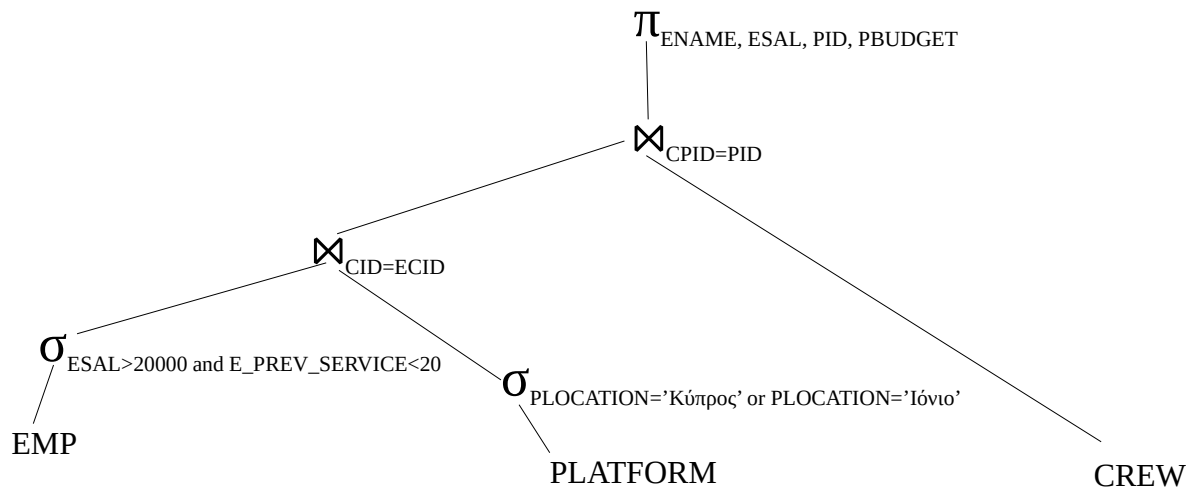
Τυφλή μετάφραση από SQL:



1ο εναλλακτικό πλάνο:



2ο εναλλακτικό πλάνο:



Στα δύο παραπάνω πλάνα πληρούνται οι εξής προϋποθέσεις:

- Τα selections έχουν τοποθετηθεί όσο πιο κοντά στα φύλλα γίνεται.
- Για κάθε τελεστή γίνονται projections με βάση τα πεδία που χρειάζεται ο DBMS.
- Δεν περιλαμβάνονται καρτεσιανά γινόμενα.
- Τα πλάνα αποτελούν αριστεροβαθή δέντρα, δηλαδή, σε κάθε join πράξη, εκτός της πρώτης, αριστερά βρίσκεται το output της προηγούμενης πράξης και δεξιά βρίσκεται πίνακας.

3.

Θα επιλέξω το 1ο εναλλακτικό πλάνο εκτέλεσης.

Γενικά, για nr εγγραφές της σχέσης r, και την πράξη $\sigma_{A \leq v}(r)$, με v σταθερά και A το γνώρισμα που με ενδιαφέρει ισχύει:

$$\text{εγγραφές} = nr * (v - \min(A, r)) / (\max(A, r) - \min(A, r))$$

και όταν, στην επιλογή υπάρχει σύζευξη, ισχύει ο τύπος:

$$nr * (\sum \pi\theta_i) / nr^n$$

$$\sigma_{\text{ESAL} > 20000 \text{ and } \text{E_PREV_SERVICE} < 20}$$

$$nr = 1000$$

$$\min(\text{EMP}, \text{ESAL}) = 20.000$$

$$\max(\text{EMP}, \text{ESAL}) = 600000$$

$$\min(\text{EMP}, \text{E_PREV_SERVICE}) = 0$$

$$\max(\text{EMP}, \text{E_PREV_SERVICE}) = 39$$

Επειδή το γνώρισμα E_PREV_SERVICE ξεκινά απ' το 0, προσθέτω +1 και +2, για να προκύψουν συνολικά 40 πλειάδες. (έστω 1 με 41)

$$\pi\theta_1 = 1000 * (60000 - 20001) / (60000 - 20000) = 1000 * 39999 / 40000 = 39999000 / 40000 = 999,975$$

$$\pi\theta_2 = 1000 * (21 - 1 - 1) / (41 - 1) = 950$$

$$\text{εγγραφές} = 1000 * \pi\theta_1 * \pi\theta_2 / 1000 * 1000 = 949976250 / 1000000 = 949,97625$$

(~950)

$$\sigma_{\text{PLOCATION} = \text{'Κύπρος'} \text{ or } \text{PLOCATION} = \text{'Ιόνιο'}}$$

Ο τύπος για ισότητα είναι: $\text{selectivity} = nr / V(A, r)$

όπου $V(A, r)$ είναι ο αριθμός των ξεχωριστών τιμών εμφάνισης της σχέσης r, για την ιδιότητα A.

και για τη διάζευξη ισχύει ο τύπος $\text{selectivity} = 1 - (1 - \text{selectivity}_1/nr) * (1 - \text{selectivity}_2/nr) * \dots$
και έτσι παίρνω μια πιθανότητα η οποία ονομάζεται

$$nr = 20$$

$$V(\text{PLOCATION}, \text{PLATFORM}) = 5$$

$$\pi\theta_1 = 20/5 = 4$$

$$\pi\theta_2 = 20/5 = 4$$

$$\text{selectivity} = 1 - (1 - 4/20) * (1 - 4/20) = 1 - 4/5 * 4/5 = 1 - 0.64 = 0.36$$

και

$$\text{εγγραφές} = nr * \text{selectivity} = 20 * 0.36 = 7,2 \quad (\sim 7)$$

$\boxtimes_{\text{CID=ECID}}$

Εγγραφές που έρχονται από αριστερά: 950

Εγγραφές που έρχονται από δεξιά: Δεν γνωρίζω

Αφού γίνεται συνένωση με γνώρισμα το κλειδί CID της σχέσης CREW, το οποίο έχει ξένο κλειδί το ECID, της σχέσης EMP, οι εγγραφές που θα προκύψουν θα είναι αυτές του αριστερού μέλους.

Επομένως: Εγγραφές = 950

$\boxtimes_{\text{CPID=PID}}$

Εγγραφές που έρχονται από αριστερά: 950

Εγγραφές που έρχονται από δεξιά: 7

Αφού γίνεται συνένωση με γνώρισμα το κλειδί PID της σχέσης PLATFORM, το οποίο έχει ξένο κλειδί το CPID, της σχέσης CREW, οι εγγραφές που θα προκύψουν θα είναι αυτές του αριστερού μέλους.

Επομένως: Εγγραφές = 950

$\Pi_{\text{ENAME, ESAL, PID, PBUDGET}}$

Με την προβολή από χώρο με αρκετές διαστάσεις θα μετατρέψω το μέγεθος των εγγραφών σε χώρο λιγότερων διαστάσεων.