OPTIMIZATION OF SELECTIONS IN RELATIONAL EXPRESSIONS

J. DEMETROVICS, HO THUAN

Computer and Automation Institute, Hungarian Academy of Sciences Budapest, Hungary

NGUEN THANH THUY

Hanoi Polytechnical Institute Vietnam

ABSTRACT

One operation met usually in the relational expressions is the selection of a relation R on a conditional expression $E = \bigwedge_{i=1}^{A} E_i$. In this paper, basing upon the estimations of probability of the tuples reR satisfying E_i , we shall show one simple O(nlogn) algorithm, where n is the length of E, rearranging the sub-expressions of E and so, the average probabilistic complexity of the algorithm for finding

 $\sigma_{\rm F}(R) = \{r \in R/r \text{ satisfies } E\}$

is minimal.

§O. INTRODUCTION

On operation met usually in expressions of relational algebra is the selection of a relation \mathcal{R} on a conditional expression \mathcal{E} . In general, it requires time O(N), where N is the number of the tuples in \mathcal{R} , to perform that selection. However, when it is discussed in the common situation with respect to other operations, for instance, projection, join, cartesian product ... the following principle is in priority: "Perform the selections and the projections as early as possible."

The transformation

$$\sigma_{\mathcal{E}}(\mathcal{R}) \Longrightarrow \sigma_{\mathcal{E}_1}(\sigma_{\mathcal{E}_2}(\ldots\sigma_{\mathcal{E}_m}(\mathcal{R})\ldots))$$

when & is of the form

$$\mathcal{E} = \bigwedge_{i=1}^{m} \mathcal{E}_i$$

is performed for the above principle. When an initial parse tree of a relational expression is reduced to a better form by the general optimization principles for relational expressions [1,3,4], it is possible that in the obtained parse tree there is a conjunctive selection

$$\sigma_{E_1} \wedge \ldots \wedge E_n^{(R)}$$

of a relation R on

 $E = E_1 \wedge \ldots \wedge E_n$.

Due to the commutativity and the associativity of the operation \wedge , the relational expression

$$\sigma$$
 n (R)
 $\bigwedge_{n=1}^{E}$ i

can be reduced to

$$\sigma$$
 n (R)
 $\bigwedge_{i=1}^{E} \tau_{i}$

where $\tau = \{\tau_1, \ldots, \tau_n\}$ is a permutation of $\{1, \ldots, n\}$. That is why we want to find the best permutation $\tau = \{\tau_1, \ldots, \tau_n\}$ such that the time complexity (cost) to find $\sigma_E(R)$ is minimal. In this paper, basing upon the estimations of probability of the tuples rER satisfying the logical expressions E, and the definition of the average probabilistic complexity of an algorithm, the best ordering $\tau = \{\tau_1, \ldots, \tau_n\}$ of the subexpressions E_i , $i=\overline{1,n}$ will be obtained such that the average probabilistic cost (complexity) of the algorithm finding $\sigma_E(R) = \{r R/r \text{ satisfies } E\}$ is minimal.

When E is an arbitrary a logical expression (as defined in §1) it can be reduced to the conjunctive - disjunctive normal form

Е	=	n V	n _i ∧	E.
		i=1		Ēj

where E_j^i is of the form either AOB or AOc or cOA, where A, B are attributes, c is a constant and O is a comparision operator $\Theta \in \{ \ge, >, <, \le, =, \neq \}$.

As the algorithm for a conjunctive selection can be used for a disjunctive selection with some modifications, so for a given arbitrary logical expression E, it is possible to find the best ordering of the subexpressions of E such that the average probabilistic cost of the algorithm finding $\sigma_{\rm F}({\rm R})$ is minimal.

It is interesting that when the cost to find the best ordering is added to the cost of the algorithm finding

$$\sigma_{\rm E}^{\rm (R)} = \sigma_{\rm n}^{\rm n} \varepsilon_{\tau_{\rm i}}^{\rm (R)}$$
,

the total cost remains desirable i.e. is less than the cost of the algorithm finding

$$\sigma_{E}(R) = \sigma_{n}(R)$$
$$\bigwedge_{i=1}^{K} i$$

with large N. The algorithm shown here for finding the best ordering τ of the subexpressions of E can be implemented in the computers as a subroutine without any access to the secondary memory devices containing the file R. Its time complexity is of O(nlogn) where n is the length of E.

§1. BASIC DEFINITIONS

Definition 1: A relation R with a set of attributes $U=\alpha(R)=\{A_1,\ldots,A_k\}$ and the corresponding ranges D_1,\ldots,D_k is defined as follows

 $R \stackrel{\text{def}}{=} \{r: \cup \to \bigcup_{i=1}^{k} D_i \mid \forall i \ l \leq i \leq k \ r(A_i) \in D_i \}$

or

$$\mathbb{R} \stackrel{\text{def}}{=} \{ r = \langle t_1, t_2, \dots, t_k \rangle | \forall i \ l \leq i \leq k \ t_i \in D_i \}$$

Eack rER is called a tuple of R.

Definition 2: A logical expression E in R with the set of attributes $U=\alpha(R)$ and the ranges D_1, \ldots, D_k can be defined recursively as follows:

- 1) An expression of the form AΘB, AΘC, CΘA, A,B∈U, k c ∈ U D, Θ∈{>, >, ≤, <, =, ≠}, is a simple logical i=1 expression.
- 2) If E_1 , E_2 are logical expressions, then $E_1 \lor E_2$, $E_1 \land E_2$, $\neg E_1$ are also logical expressions.

Definition 3: A logical expression E in R which is of the form $E=E_1 \wedge \dots \wedge E_n$ is called conjunctive logical expression.

<u>Definition 4:</u> Given a logical expression E in a relation R with the set of attributes v and the ranges D_i , $i=\overline{1,k}$, and a tuple r of R.

We say that \mathbf{r} satisfies E if when subtituting the names of the attributes A in E by the value r. $A \in \bigcup_{i=1}^{k} D_i$ of the tuple $r \in R$,

the obtained logical expression has the value "true".

<u>Definition 5:</u> Given a relation R and a logical expression E. The selection of the relation R on the condition E, denoted by $\sigma_{_{\rm F}}({\rm R})$ is defined as follows:

$$\sigma_{E}(R) = \{r \in R | r \text{ satisfies } E\}$$
.

If E is a conjunctive logical expression, the selection $\sigma_{E}^{(R)}$ is called a conjunctive selection.

<u>Definition 6:</u> Let Ω be a probability space of finite cardinality, i.e. in Ω is defined a probability measure $\rho: 2^{\Omega} \longrightarrow [0,1]$ satisfying the probability axioms. Put

and

$$\Omega = \{\omega_1, \omega_2, \dots, \omega_g\}$$

$$\rho_i = \rho(\{\omega_i\}), i = \overline{i, g}$$
.

Then, the average probabilistic value of the real valued function f defined on Ω corresponding to the probability measure ρ is defined as

$$\bar{f}_{\rho} = \sum_{i=1}^{g} f(\omega_i) \rho_i .$$

§2. AN APPROACH TO THE PROBABILITY ESTIMATIONS

Let a relation R be given with the set of attributes υ and the ranges D_i , $i=\overline{1,k}$. As defined in definition 2, a logical expression can be constructed by the simple logical expressions and the logical operators { \land, \lor, \neg }.

The following statistical parameters obtained and collected during the manipulation of R can be estimated:

1) The distribution of values of the attributes in v.

1

- The exact upper bound and lower bound for the attributes in v during the manipulation.
- 3) The number of distinct values of an attribute. One of the simple assumptions is that the values of every attribute X € v are uniformly distributed in the segment [X_m, X_M] where X_m=min R [X], X_M=max R[X]. (H1).

With the assumption (H1), it is easy to give the probability estimations Pr(E) for the tuple $r \in R$ satisfying E.

For instance, for X, $Y \in v$; $a, b \in R[X]$, we have

$$Pr(X=a) = \frac{1}{card R[X]}; \qquad (2.1)$$

$$Pr(X \ge a) = \begin{cases} \frac{X_{M}^{-a}}{X_{M}^{-X_{m}}}, & X_{M} \le a < X_{M} \\ \frac{1}{card R[X]}, & a = X_{M} \end{cases}$$
(2.2)

$$Pr(X > a) = Pr(X \ge a) - \frac{1}{card R[X]}$$
(2.3)

$$\Pr(a \leq X \leq b) = \begin{cases} \frac{b-a}{X_{M}-X_{m}}, & X_{m} \leq a \leq b \leq X_{M} \\ \\ \frac{X_{M}-a}{X_{M}-X_{m}}, & \frac{1}{card R[X]}, & X_{m} \leq a \leq b = X_{M} \end{cases}$$
(2.4)

$$Pr(X=Y) = 0 \qquad \text{if } X_{M} < X_{M} < Y_{M} < Y_{M} \qquad (2.5)$$

$$or \ Y_{m} < Y_{M} < X_{M} < X_{M} < X_{M} \qquad (2.5)$$

$$Pr(X=Y) = \frac{1}{card R[X] \cdot card R[Y]}$$
(2.6)
if $X_m < X_M = Y_m < Y_M \cdot$
or $Y_m < Y_M = X_m < X_M$

$$Pr(X=Y) = \frac{1}{\max\left(\frac{X_{M}-Y_{m}}{X_{M}-X_{m}} \operatorname{card} R[X], \frac{X_{M}-Y_{m}}{Y_{M}-Y_{m}} \operatorname{card} R[Y]\right)}$$
(2.7)

$$if \quad X_{m} < Y_{m} < X_{M} < Y_{M}$$

$$or \quad Y_{m} < X_{m} < Y_{M} < X_{M}$$

$$Pr(X=Y) = \frac{1}{\max(\operatorname{card} R[X], \operatorname{card} R[Y])}$$
(2.8)

if $X_m = Y_m < X_M = Y_M$

(a special case of (2.7)).

<u>Remark 1.</u> To compute Pr(E) for a non simple logical expression, we use the following rules:

a)
$$Pr(E_1 \wedge E_2) = Pr(E_1) \cdot Pr(E_2)$$
 where $E_1 \neq E_2$ are independent.
b) $Pr(E_1 \vee E_2) = Pr(E_1) + Pr(E_2) - Pr(E_1 \wedge E_2)$;
c) $Pr(\neg E_1) = 1 - Pr(E_1)$.

<u>Remark 2.</u> In [2], for the computation of Pr(X=Y), the authors gave only formula (2.8), with the assumption (H1). Of course, when taking attention to the relative positions of the segments $[X_m, X_M]$ and $[Y_m, Y_M]$, this simple formula is not complete.

<u>Remark 3.</u> For other distribution types, the idea of this paper is also useful. One must only considered an appropriate method for computing the probability estimations.

§3. THE MAIN PROBLEM

Given a relation R with the set of attributes $U=\alpha(R)$ and the ranges D_i , $i=\overline{1,k}$. E is a logical expression in R. The selection $T=\sigma_E(R)$ can be obtained by the following algorithm:

$$T:=\phi$$
for each rER do
if test (r,E) then add r to T
$$\left. \begin{array}{c} (1) \end{array} \right.$$

The function test (r,E) performs two operations:

- i) Replaces the names of attributes $A \in U$ by the values r.A of the tuple reR .
- ii) Computes the obtained logical expression and assigns the result to the function test.

Given a tuple r $\in R$. Denote cost (r,E) the cost paid to perform the function test (r,E). In this section, we always consider that E is a conjunctive logical expression

$$E = \bigwedge_{i=1}^{n} E_{i}$$

The function test (r,E) can be computed by two different ways:

Way 1: Compute all functions test (r,E,) and then set

test(r,E) =
$$\bigwedge_{i=1}^{n}$$
 test (r,E_i).

We have the algorithm:

test (r,E): = <u>true</u>; (2) <u>for</u> i : = 1 <u>to</u> n <u>do</u> test (r,E)=test(r,E) \land \land test (r,E_i) .

The one-pass compilers compute often the conjunctive logical expressions in this way, for instance the compiler on the

hypothetical computer P-code for language PASCAL-S.

<u>Way 2:</u> It is obvious that if there is some i_0 during the computation of test (r, E_i) such that $test(r, E_{i0}) = \underline{false}$ then it is possible to conclude immediately that $test(r, E) = \underline{false}$ and to halt the calculation of test (r, E). This is expressed in the algorithm as follows.

Test
$$(r,E)$$
: = false
for i : = 1 to n do
if \neg test (r,E_i) then goto L;
test (r,E) : = true; (3)

L:

Intuitively, it is easy to see that the method in the algorithm (3) is very natural. (Of course it is better than the algorithm (2).) However, it is very interesting if we know the probabilities of the tuples $r \in R$ satisfying the expressions E_i in R and so we can expect that there exist a best ordering of the subexpressions E_i such that the average probabilistic cost of the algorithm (3) is minimal.

To make clear this idea we do as follows:

At first, basing on the probability estimations $s_i = Pr(E_i)$, $i = \overline{1, n}$ of E_i in R and the costs $c_i = cost(r, E_i)$ $i = \overline{1, n}$ paid to compute the functions test (r, E_i) , we can compute the average probabilistic cost of the algorithm (3). Then, analyzing the mathematical expression cost (r, E) represented by c_i , s_i $i = \overline{1, n}$, we try to find the best ordering $\tau = \{\tau_1, \dots, \tau_n\}$ -a permutation of $\{1, \dots, n\}$, such that the value of the expression cost (r, E) is minimum.

Note that to compute test (r,E) for a given E and rER, the replacements in step i) are necessary and the time cost is the same for every rER.

It is obvious that:

$$cost (r, E_i) = c_i > 0 (constant) \forall r \in \mathbb{R}$$
 (4)

Assume that for each E_i , by the rules as in §2 we can def $s_i = Pr(E_i)$, $i = \overline{1, n}$ and the logical subexpressions are independent of each other.

The relation R is partitioned into T and T_i , $i=\overline{1,n}$ as follows:

R = T II (U T)

$$T = \sigma_{E}(R) = \{r \in R/test (r, E) = true\}$$

$$T_{i} = \{r \in R/\forall j, j=\overline{1,i-1} test(r, E_{j}) = true, test(r, E_{i}) = false$$

n

it is evident that

$$T \cap T_{j} = \phi, j = \overline{1, n}$$
$$T_{i} \cap T_{j} = \phi, j \neq i.$$

Define

$$\rho^{*}(T) = \prod_{i=1}^{n} s_{i}, \rho^{*}(T_{i}) = \prod_{j=1}^{i-1} s_{j}(1-s_{i}), i=\overline{2,n}$$

$$\rho^*(T_1) = 1 - s_1$$

We have

$$\rho^{*}(T) + \sum_{i=1}^{n} \rho^{*}(T_{i}) = \sum_{i=1}^{n} \prod_{j=1}^{n} s_{j} (1-s_{i}) + \prod_{j=1}^{n} s_{j} = 1$$

Indeed, set $h_i = \prod_{j=1}^{i} s_j, h_o = 1$

$$\rho^{*}(T_{i}) = \prod_{j=1}^{i-1} s_{j}(1-s_{i}) = h_{i-1} - h_{i}, \quad i=2,n$$

$$\rho^{*}(T_{1}) = 1-s_{1} = h_{0} - h_{1}$$

$$\sum_{i=1}^{n} \rho^{*}(T_{i}) + \rho^{*}(T) = \sum_{i=1}^{n} (h_{i-1} - h_{i}) + h_{n} = h_{o} = 1$$

Moreover in T i.e. for the tuples rER for which test(r,E)=true, it is necessary to compute all test (r,E_i), $i=\overline{1,n}$ for the final result of test (r,E), therefore it requires $\sum_{i=1}^{n} c_{i}$.

In T_i , because test $(r, E_i) = \underline{false}$, the computation of test(r, E)halts and it takes $\sum_{j=1}^{\Sigma} c_j$. By the definition 6, the average probabilistic cost of the algorithm (3) is

 $\overline{\operatorname{cost}_{3}}(r,E) = \rho^{*}(T)\operatorname{cost}(r\in T,E) + \sum_{i=1}^{n} \rho^{*}(T_{i})\operatorname{cost}(r\in T_{i},E) =$

 $= \begin{pmatrix} n & n & n & i-1 & i \\ (\Sigma & c_{i}) & \Pi & s_{i} + \Sigma & \Pi & s_{j} (1-s_{i}) (\Sigma & c_{j}) \\ i=1 & i=1 & i=1 & j=2 & j & j=1 \end{pmatrix}$ (5)

Return to the algorithm' (2) computing the function test(r,E), the worst-case cost and the average probabilistic cost are the same. We have:

$$\operatorname{cost}_{2}(\mathbf{r}, \mathbf{E}) = \sum_{i=1}^{n} c_{i}$$
(6)

The following result is obvious.

Proposition 1.

$$cost_{2}(r,E) \leq cost_{2}(r,E)$$

where $\overline{\text{cost}_3}(r, E)$ and $\text{cost}_2(r, E)$ are expressed by the formulae (5), (6) respectively.

Proof. It is not difficult to see that:

$$\overline{\operatorname{cost}_{3}}(\mathbf{r}, \mathbf{E}) \leq (\sum_{i=1}^{n} c_{i}) \begin{bmatrix} n & n & i-1 \\ \Pi & s_{i} + \sum_{i=1}^{n} \Pi & s_{i} \end{bmatrix} = \sum_{i=1}^{n} c_{i} = \sum$$

$$= cost_2(r, E)$$

The above proof shows that the algorithm (3) has always the cost less than the cost of algorithm (2).

Now, (5) can be transformed in the following way:

$$\overline{\operatorname{cost}}_{3}(\mathbf{r}, \mathbf{E}) = (\sum_{i=1}^{n} c_{i}) \prod_{i=1}^{n} s_{i} + \sum_{i=1}^{n} (1-s_{i}) \prod_{j=1}^{n} s_{j} (\sum_{j=1}^{i} c_{j})$$

Set

$$g_{i} = \sum_{j=1}^{i} c_{j}, g_{o}=0$$

 $\overline{\operatorname{cost}_{3}}(r,E) = g_{n}h_{n} + \sum_{i=1}^{n} (h_{i-1}-h_{i})g_{i} =$

=

$$g_{n}h_{n} + \sum_{i=1}^{n} h_{i-1}g_{i} - \sum_{i=1}^{n} h_{i}g_{i} =$$

$$= \sum_{i=1}^{n} h_{i-1}g_i - \sum_{i=1}^{n-1} h_ig_i = \sum_{i=1}^{n} h_{i-1}g_i - \sum_{i=1}^{n} h_{i-1}g_{i-1} =$$

 $= \sum_{i=1}^{n} h_{i-1} (g_i - g_{i-1}) = \sum_{i=1}^{n} c_i \prod_{j=1}^{i-1} s_j$

From here the following problem can be formulated:

Given the numbers
$$c_i > 0, i=\overline{1,n}$$

 $1 \ge S_i \ge 0$

Find the best permutation $\tau = \{\tau_1, \ldots, \tau_n\}$ of $\{1, \ldots, n\}$ such that

$$A(\tau) = \sum_{i=1}^{n} c_{\tau_{i}} \prod_{j=1}^{n-j} s_{\tau_{j}} \longrightarrow \min .$$

If it is possible to find the best permutation τ such that $A(\tau) \longrightarrow \min$, then by the commutativity and the associativity of

the logical operator \wedge , we have

$$\sigma_{E}(R) = \sigma_{A} \alpha_{E}(R) = \sigma_{A} \alpha_{E}(R)$$
$$\stackrel{(R)}{\underset{i=1}{\wedge}} \sigma_{i} \alpha_{i} \alpha_$$

This transformation should allow us to calculate the function test (r,E) by the algorithm (3) not with the ordering $\{1,\ldots,n\}$ of E_i but with the ordering $\{\tau_1,\ldots,\tau_n\}$. And so, the algorithm (3) becomes:

Test $(r,E) := \underline{false};$ <u>for</u> $i := l \underline{to} n \underline{do}$ <u>if</u> \neg test $(r,E_{\tau}) \underline{then} \underline{go to} L;$ test $(r,E) := \underline{true}$. (3.1)

L :

For this algorithm, the function test (r,E) can be computed with the average cost

$$\begin{array}{cccc} n & i-1 \\ \Sigma & C_{\tau} & \Pi & s_{\tau} & \longrightarrow & \text{min} \\ i=1 & i & j=1 & j \end{array}$$

The following proposition will show the way to find the best T.

Proposition 2.

Let $s_i > 0, i = \overline{1, n}$, τ , τ' be two permutations of $\{1, \ldots, n\}$ whose i_0 -th and (i_0+1) -th elements are changed with each other, i.e.

$$\tau'_{io} = \tau_{i_0+1}, \tau'_{i_0+1} = \tau_{i_0}.$$

$$t_{\tau_{i_{o}}} = \frac{u_{\tau_{i_{o}}}}{c_{\tau_{i_{o}}}} < t_{\tau_{i_{o}}+1} = \frac{u_{\tau_{i_{o}}+1}}{c_{\tau_{i_{o}}+1}}$$

If

 $u_i = 1-s_i, i=\overline{1,n}$.

Then

 $A(\tau') < A(\tau)$

Proof.

$$A(\tau) = \sum_{i=1}^{n} c_{\tau_{i}} \prod_{j=1}^{n-1} s_{\tau_{j}}$$
$$A(\tau') = \sum_{i=1}^{n} c_{\tau, \prod_{i=1}^{i-1}} s_{\tau'_{j}}$$

when i < io

$$\begin{array}{cccc} \Pi & \mathbf{s}_{\tau'} = \Pi & \mathbf{s}_{\tau} \\ \mathbf{j} = \mathbf{l} & \mathbf{j} & \mathbf{j} = \mathbf{l} & \mathbf{j} \end{array}$$

i-1 i-1

when i=io

$$c_{\tau'_{i_o}} = c_{\tau_{i_o+1}}$$

c_{t'} = c_t

 $i_0 - 1$ $i_0 - 1$ Π $s_{\tau'} = \Pi$ s_{τ} j = 1 j j = 1 j

when i=i_+1

$$i_{0} = \pi i_{0} + 1 = \pi i_{0}$$

$$i_{0} = \pi i_{0} + 1 = \pi i_{0}$$

$$i_{0} = \pi i_{0} + 1 = \pi i_{0} + 1$$

$$j = 1 = \pi j = \pi i_{0} + 1$$

when i>i_o+1

$$c_{\tau'} = c_{\tau_i}$$

$$= \frac{i_{0}^{-1}}{j_{1}} s_{\tau_{j}} \cdot s_{\tau_{i_{0}+1}} \cdot s_{\tau_{i_{0}}} \cdot \frac{i_{-1}}{j_{1}} s_{\tau_{j}} =$$

$$= \frac{i_{-1}}{j_{1}} s_{\tau_{j}} \cdot$$

$$A(\tau') - A(\tau) = c_{\tau_{i_{0}+1}} \frac{i_{0}^{-1}}{j_{1}} s_{\tau_{j}} + c_{\tau_{i_{0}}} \frac{i_{0}^{-1}}{j_{1}} s_{\tau_{j}} \cdot s_{\tau_{i_{0}+1}}$$

$$- c_{\tau_{i_{0}}} \frac{i_{0}^{-1}}{j_{1}} s_{\tau_{j}} - c_{\tau_{i_{0}+1}} \frac{i_{0}^{-1}}{j_{1}} s_{\tau_{j}} \cdot s_{\tau_{i_{0}}}$$

$$= \frac{i_{0}^{-1}}{j_{1}} s_{\tau_{j}} (c_{\tau_{i_{0}+1}} + c_{\tau_{i_{0}}} s_{\tau_{i_{0}+1}} - c_{\tau_{i_{0}} - c_{\tau_{i_{0}+1}} s_{\tau_{i_{0}}})$$

$$= \frac{i_{0}^{-1}}{j_{1}} s_{\tau_{j}} (c_{\tau_{i_{0}+1}} + c_{\tau_{i_{0}}} s_{\tau_{i_{0}+1}} - c_{\tau_{i_{0}} - c_{\tau_{i_{0}+1}} s_{\tau_{i_{0}}})$$

$$= \frac{i_{0}^{-1}}{j_{1}} s_{\tau_{j}} (c_{\tau_{i_{0}+1}} (1 - s_{\tau_{i_{0}}}) - c_{\tau_{i_{0}}} (1 - s_{\tau_{i_{0}+1}}))$$

$$= \frac{i_{0}^{-1}}{j_{1}} s_{\tau_{j}} \cdot c_{\tau_{i_{0}}} \cdot c_{\tau_{i_{0}+1}} (\frac{1 - s_{\tau_{i_{0}}}}{c_{\tau_{i_{0}}}} - \frac{1 - s_{\tau_{i_{0}+1}}}{c_{\tau_{i_{0}+1}}}) < 0$$

$$\rightarrow A(\tau') < A(\tau) .$$

- 41 -

Proposition 3.

Given the numbers

 $\begin{array}{c} 0 \leqslant s_{i} \leqslant 1 \\ 0 \leqslant c_{i} \end{array} \right\} i = \overline{1, n}$

Set

$$t_i = \frac{1 - s_i}{c_i}$$

If $\tau = \{1, ..., n\}$ satisfies $t_i \ge t_{i+1}$, i=1, n-1 then $A(\tau_0)$ is the minimum.

Proof.

Let $\tau = \{\tau_1, \ldots, \tau_n\}$ be a permutation of $\{1, \ldots, n\} = \tau_0$. We have to prove that $A(\tau_0) \leq A(\tau)$.

First we remark that from $\{t_{\tau_i}\} \otimes$ the sequence $\{t_i\}, i=\overline{1, n}$ (corresponding to τ_0) can be obtained by

- (i) the bubble sorting algorithm permuting sequentially the adjacent elements t_{τ} , t_{τ} satisfying $t_{\tau} < t_{\tau}$ and i_{0} i_{0} +1 i_{0} i_{0} +1
- (ii) the permutations (if necessary) of the elements with equal values in the obtained sequence.

By proposition 2, if (i) should be carried out then we should obtain τ^1 satisfying $A(\tau^1) < A(\tau)$.

Basing upon the proof of the proposition 2, we have: The permutations of the elements with equal values in the sequence $\{t_{\tau_1}\}\$ do not change the value of A, i.e. $A(\tau_0)=A(\tau^1)$. i i=1,n

From here follows $A(\tau_0) = A(\tau^1) < A(\tau)$. When the step (i) does not take places, it is not difficult to see that $A(\tau_0) = A(\tau)$.

The following algorithm will give the best result for any conjunctive logical expression.

Algorithm Al.

Input: $E = E_1 \land \dots \land E_n$

$$R = \{r : U \longrightarrow \bigcup_{i=1}^{k} D_i / \forall i r(A_i) \in D_i\}$$

Output:

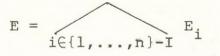
 $\tau = \{\tau_1, \dots, \tau_n\}$ is the permutation of $\{1, \dots, n\}$ such that

$$A(\tau) = \sum_{i=1}^{n} c_{\tau i} \prod_{j=1}^{i-1} s_{\tau} \longrightarrow \min_{i=1}^{n}$$

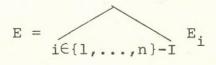
Method

- For each i, estimate the probability Pr(E_i) by the formulae in §2 or by the formulae given by the system programmers basing upon the statistical parameters during the manipulation of R.
- 2) If there exists i such that $Pr(E_i) = 0$, then inform $\sigma_F(R) = \phi$.
- 3) If there exist i such that $Pr(E_i) = 1$ then delete E_i from E.

More generally, denote $I = \{i_0 / S(E_{i_0}) = 1\}$. Consider



Renumber the expressions E, in



n := n - card I

4) Define c_i, i=1,n (In practice, in order to define c_i, we compute the function test (r,E_i) for any tuple and give c_i = cost (r,E_i).

For instance: if

 $E_i = A\Theta B$ then $c_i = 2a + b$ $E_i = A\Theta c$, $c_i = a + b$ $E_i = c\Theta A$, $c_i = a + b$

where

a is the cost to bustitute A by r.A;

b is the cost paid to compare two elements in U D_i . i=1

k

5)
$$u_i = 1 - s_i$$
, $i = \overline{1, n}$.

6)
$$t_i = u_i / c_i$$
, $i = \overline{1, n}$.

7) Sort {t_i} such that t_i ≥ t_{i+1} i=1,n-1 Step 7 can be performed by one of the sorting algorithms, in general, of complexity O(nlogn).

8) Print the best ordering obtained $\tau = \{\tau_1, \ldots, \tau_n\}$.

9) Print the value $A(\tau) = \sum_{i=1}^{n} c_{\tau} \prod_{j=1}^{n-1} s_{\tau}$.

Costing the algorithm Al.

Steps	1,2,3	are	pert	formed	with	the	cost	Kl	n
step	4	has	the	comple	exity			к2	n
step	5,6							к3	n
step	7							К4	nlogn
step	8,9							К5	n

Kn + H n log n
$$\approx O(nlogn)$$

K = Kl + K2 + K3 + K5
H = K4 .

Remark 4.

The proof of the proposition 3 is based on the bubble sorting algorithm of complexity $O(n^2)$ but the step 7 of the algorithm Al uses any sorting algorithm of complexity O(nlogn). However, there is no matters about the correctness of the algorithm Al.

The algorithm Al can be implemented without any access to the secondary memory devices containing the file R.

Theorem 1.

Let R be a relation, card R = N, E = $\bigwedge_{i=1}^{n} E_i, \tau_0 = \{1, \dots, n\}$ is the best ordering of E_i 's i.e.

$$\begin{array}{cccc}n & i-1\\ A(\tau_{O}) &= & \Sigma & C & \Pi & s \\ & i=1 & j=1 & j \end{array} \longrightarrow \min$$

Then, the cost of the algorithm (1) finding $\sigma_{\rm E}^{}({\rm R})$ with the function test (r,E) computed by:

- 1) Algorithm (2) is $C^1 = N$. $\sum_{i=1}^{n} C_i$
- 2) Algorithm (3) with the best ordering τ_0 of E's is

$$C^{2} = N \cdot \sum_{i=1}^{n} C_{i} \prod_{j=1}^{i-1} S_{j} + F(n)$$

where F(n) = Kn + Hnlogn is the cost paid to perform the algorithm Al.

3) Algorithm (3) with an arbitrary ordering $\tau = \{\tau_1, ..., \tau_n\}$ is

$$C^{3} = N \Sigma C_{\tau} I S_{\tau}$$

$$i=1 i j=1 j$$

we have the inequalities:

$$C^{1} \ge C^{3}$$

 $C^{1} \ge C^{2}$ with large N
 $C^{3} \ge C^{2}$ with large N.

§4 Extensions

Extension 1. If E is of the form $E=E_1 \vee \ldots \vee VE_n$ then using the symbols as above and the De Morgan's law

$$\neg (E_1 V \dots VE_n) = \neg E_1 \land \dots \land \neg E_n$$

we have: the cost payed to compute test (r,E) is

 $\overline{\operatorname{cost}}(\mathbf{r}, \mathbf{E}) = \sum_{\substack{\Sigma \\ i=1}}^{n} c_{i} \prod_{\substack{j=1 \\ j=1}}^{i-1} s'_{j} \text{ where } s'_{j} = 1 - s_{j}, j = \overline{1, n}$

Proposition 4.

 $s_i = Pr(E_i)$ $0 \leq s_i \leq 1$ $c_i > 0$ $i=\overline{1,n}$ $t_i = s_i/c_i$ Given

If $\tau_0 = \{1, \dots, n\}$ satisfies $t_i > t_{i+1}$, $i=\overline{1, n-1}$ then

cost $(r,E) = \sum_{\substack{z \in C \\ i=1}}^{n} c_{i} \prod_{j=1}^{i-1} s'_{j} \longrightarrow min.$

Extension 2. Algorithm A2.

Input: An arbitrary logical expression E (as defined by def.2).

: ;

Output The best ordering of the simple logical subexpressions of E.

Method.

1) Reduce E to the conjunctive disjunctive normal form

$$r(E) = \bigvee_{i=1}^{n} \bigvee_{j=1}^{n_{i}} E_{j}^{i} \cdot$$

2) Apply algorithm Al to

$$E_{i} = \bigwedge_{j=1}^{n_{i}} E_{j}^{i}$$

to give the best ordering τ^{i} of $E^{i}_{j},\;j{=}\overline{1,n_{i}}$.

- 3) Apply the modified algorithm to $V E_i$ with $c_i = A(\tau^i)$ and $s_i = Pr(E_i)$ defined by the estimating formulae analogous. to one's in §2.
- 4) Print the best ordering $E = V \land E_{i}^{i}$. $i=1 j=1 \tau_{j}^{i}$
- 5) Print the value $C = \sum_{\substack{\Sigma \\ i=1 \\ i=1 \\ j=1 \\ j=$

CONCLUSION

Independently, our approach is quite near to the Hanani's one [5]. However, our approach seems to be more straightforward, easy for extensions and the complexity analysis of the algorithm proposed is much elaborate.

ACKNOWLEDGEMENT

The authors wish to thank Dr Béla Uhrin and Katalin Fridl for their valuable comments and suggestions.

REFERENCES

- [1] Ulmann, J.; Principles of database systems. Computer Science Press. Second edition, 1982.
- [2] Adiba, M. and Delobel, C.: Bases de données et systèmes relationnels. Paris, Dunod, 1982.
- [3] Steve Talbot: An investigation into logical optimization of relational query languages. The comp. J. Vol. 27, No.4, November 1984.
- [4] Matthias Jarke and Jürgen Koch: Query optimization in Database systems. ACM Computing Surveys, Vol.16. No.2, June 1984.
- [5] Hanani, M.Z.: An optimal evaluation of Boolean expression in an on line query system. Comm. ACM 20, 5, May 1977.

- 49 -

J. DEMETROVICS, HO THUAN, NGUEN THANH THUY

Összefoglaló

Legyen $E = \bigwedge_{i=1}^{n} E_i$ egy feltételes kifejezés és R egy

reláció. A cikkben egy O(n log n) algoritmust mutatnak be a szerzők, amely a $\sigma_{E}(R) := \{r \in R/r \text{ kielégiti az E-t}\}$ mennyiséget /átlagban/ minimális lépésszámban határozza meg /azaz, amely komplexitásának várható értéke minimális/.

Оптимизация выборок из реляциённых выражений.

Й. Деметрович, Хо Тхуан, Нгуен Тханх Тхуи

Резюме

Пусть $E = \bigwedge_{i=1}^{n} E_i$ есть условное выражение и R реляция. В статье показывается O/n logn/ алгоритм который /в среднем/ минимизирует число шагов для нахождения $\sigma/R/: = \{r \in R/r\}$