

Final Exam

DSE201, Winter 2016

Name: Ryan Riopelle

Brief Directions:

- Write clearly!
- Good luck!

DSE 201-Final-Ryan Röpelle

1.)

$$T(R) = 10^6$$

$$T(S) = 10^5$$

$$T(W) = 10^3$$

$$V(R,C) = 10^2$$

B is a key of w.B in non-null FK that references R.B.

A is Key of S and R.A is non null FK that references S.A.

Plans

$$1) \sigma_{R.C=1} ((R \bowtie S) \bowtie W)$$

$$2) \sigma_{R.C=1} ((R \bowtie W) \bowtie S)$$

$$3) (\sigma_{R.C=1} (R \times S)) \bowtie W$$

$$4) (\sigma_{R.C=1} (R \bowtie W)) \bowtie S$$

↓

$$a) T(a) = T(R \bowtie W) = T(W) = 10^3$$

$$b) T(\sigma_{R.C=1}(R \bowtie W)) = \frac{T(a)}{V(R,C)} = \frac{10^3}{10^2} = 10$$

$$\text{Intermediate} = 10^3 + 10$$

Plan 2 and 4 are less costly

$$T(a) = T(R \bowtie W) =$$

$$\frac{T(R)T(W)}{\max(V(R,B), V(W,B))} = 10^3$$

$$(b) = T(T(a) \bowtie S) = \frac{T(a)T(S)}{V(S,A)} =$$

$$\frac{(10^3)(10^5)}{10^5} = 10^3$$

$$\text{Intermediate} = 2 \times 10^3$$

Name: Ryan Riopelle

2 Algebra and Estimation

Produce an optimal algebraic expression for the following query over tables $R(A, B)$ and $S(A, C, D)$, where "optimal" means that it has the smallest total size of intermediate results, among all possible algebraic expressions that are equivalent to this query. Write the sizes of all intermediate results.

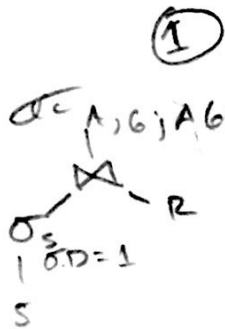
SELECT A, C, AGG(B) AS N
FROM R, S
WHERE S.A = R.A AND S.D = 1
GROUP BY A, C

given the following statistics

$$\begin{aligned} T(R) &= 10^9 \\ V(R, A) &= 10^6 \\ V(R, B) &= 10^9 \\ T(S) &= 10^{10} \\ V(S, A) &= 10^7 \\ V(S, C) &= 10^2 \\ V(S, D) &= 10 \end{aligned}$$

Assume (the common assumption) that

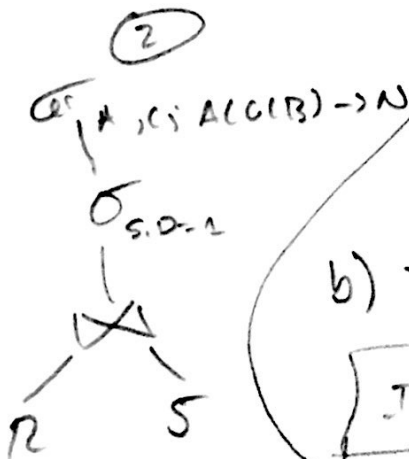
$$V(R, A) < V(S, A) \Rightarrow \pi_{AR} \subset \pi_{AS}$$



$$a) T(a) = T(\sigma_{S.D=1} S) = \frac{T(S)}{V(S, D)} = \frac{10^{10}}{10} = 10^9$$

$$b) T(R \bowtie T(a)) = \frac{T(R)T(a)}{\max(V(R, A), V(S, A))} = \frac{10^9 \cdot 10^9}{10^7} = 10^{11}$$

$$\text{Intermediate res 1} = 10^9 + 10^{11}$$



$$a) T(a) = T(R \bowtie S) = \frac{T(R)T(S)}{\max(V(R, A), V(S, A))} = \frac{10^9 \cdot 10^{10}}{10^7} = 10^{12}$$

$$b) T(\sigma_{S.D=1} T(a)) = \frac{T(a)}{V(S, D)} = \frac{10^{12}}{10} = 10^{11}$$

$$\text{Intermediate res 1} = 10^{12} + 10^{11}$$

Plan 1 is better!

Name: Ryan Riopelle

3)

$$R \bowtie_c S = \pi_A(R \bowtie_c S)$$

Except that it is wrong. Provide a correct one.

Equivalences Now, assume relations R , S and T . The notation $c(A)$ refers to a condition that refers to the list of attributes A only. Declare true or false each of the following. If the answer is "no", also provide counterexample.

True \rightarrow 1. $\sigma_{c(A)}(R \bowtie S) = (\sigma_{c(A)} R) \bowtie S$, where R has a list of attributes A and S has no attributes of A .

True \rightarrow 2. $\delta(R \bowtie S) = (\delta R) \bowtie S$

True \rightarrow 3. $(R \bowtie S) \bowtie T = (R \bowtie T) \bowtie S$, where all of R , S and T have a single common attribute A and no pair of R , S and T has a common attribute other than A .

Recall, δ is the duplicate elimination operator.

Queries

$C \bowtie_{cid = e, class} E$

Algebraic Definition

$$R \bowtie_c S = \pi_A(R \bowtie_c \delta_c S)$$

Name: Ryan Riopelle

4 Column Databases

Consider the table $R(D1, D2, D3, M)$. You already have a PostGres database and you wonder whether it is worthy to buy a column database in order to accomodate queries on R . Of course, the answer depends on knowing the queries that will be issued. For each of the following queries, declare whether a column database will be significantly better or whether PostGres is good enough (or even better). Just place a circle around the relevant system. "Significantly" means at least close to a multiple, say 2x. A 10% improvement is not significant.

When you consider column databases, assume they do not have indices. (This is not exactly accurate, as column databases also have indices, but we adopt it for the sake of the exercise.) Furthermore, assume that the table is in the order of terabytes and resides in hard disk. The main memory is only 32GB. Each one of the D attributes is an integer and M is a float. Assume that $V(D1, R) = V(D2, R) = V(D3, R) = 10^6$. Assume $V([D1, D2, D3], R) = 10^8$.

1. SELECT D1, SUM(M) FROM R GROUP BY D1: Column - Postgres
2. SELECT D1, D2, D3, SUM(M) FROM R: Column - Postgres or equal
3. SELECT D1, D2 SUM(M) FROM R WHERE D3=: Column - Postgres

As usual, "?" means that a constant will be given at query time.

- See work for more info.

Name: Ryan Piopelle

Database Final - Problem 4

4-1) Table $R(D_1, D_2, D_3, M)$, Size = Tera byte
 $V(D_1, R) = V(D_2, R) = V(D_3, R) = 10^6$
 $V([D_1, D_2, D_3], R) = 10^8$

Postgres: $O(\sigma_{D_1, D_2, D_3, M} R) = 10^8 \times [4+4+4+8] = 20 \times 10^2 \text{ bytes}$
 $= 1.86 \text{ GB}$

Column: $O(\sigma_{D_1, M} R) = 10^2 \times [4+8] = 12 \times 10^8 \text{ bytes} = 1.12 \text{ GB}$

Column DB = 40% faster

4-2) Select $D_1, D_2, D_3, \text{Sum}(M)$ From R

About equal, Postgres

4-3) Select $D_1, D_2, \text{Sum}(M)$ From R where
 $D_3 = ?$

Postgres: 1.86 GB

Column

$$T(\sigma_{D_3=R}) = \frac{T(R)}{V(D_3, R)} = \frac{V([D_1, D_2, D_3], R)}{V(D_3, R)}$$

$$O(\sigma_{D_1, D_2, M} R) = (4+4+8) \times 10^2 = 16 \times 10^2$$

$$O(\sigma_{D_3=R}) = 4 \times 10^8 = 381 \text{ MB}$$

$$O(\text{Total}) = (4 \times 10^8) + (16 \times 10^2) \approx 0.381 \text{ GB}$$

$$\frac{1.86 - 0.381}{1.86} =$$

80% faster - column