

Query Optimization

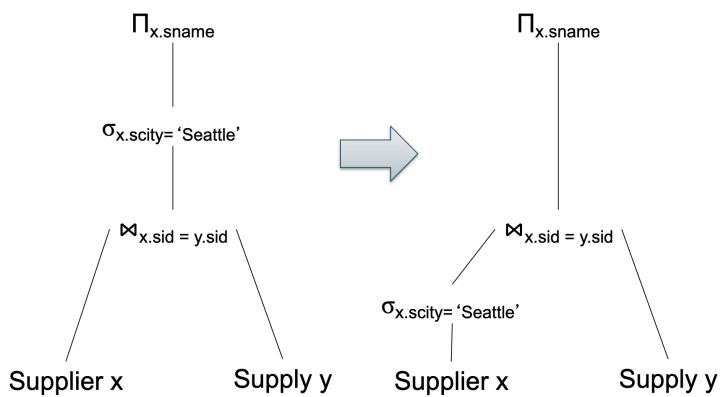
Friday, May 10, 2019 1:31 PM

- Main idea: replace a query plan with another one that is equivalent but cheaper.

- Push selections down

- Do selection before grouping or joining.

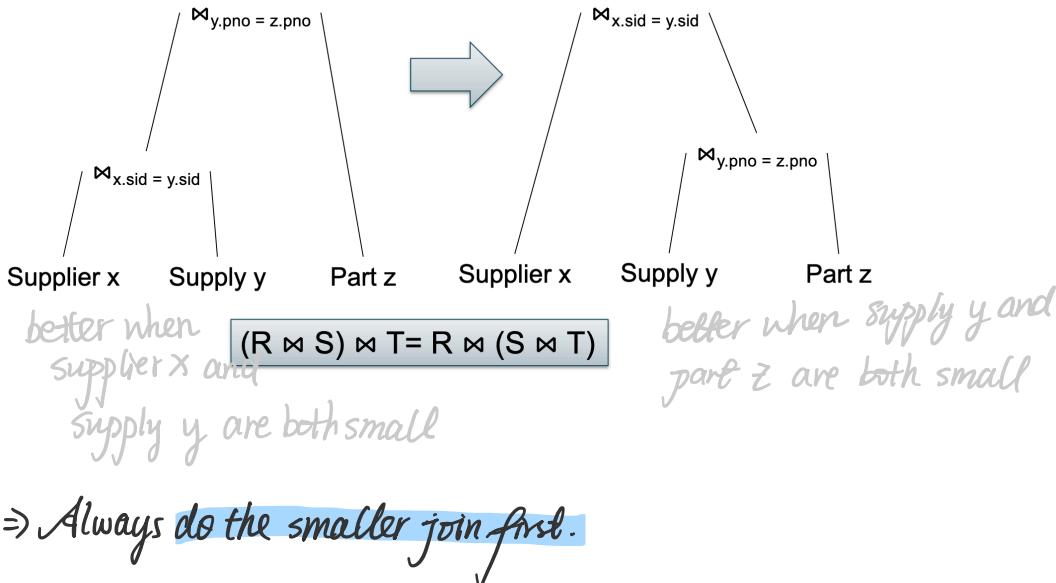
e.g. $\sigma_{C_1 \text{ and } C_2}(R \bowtie S) = \sigma_{C_1}(\sigma_{C_2}(R \bowtie S)) = \sigma_{C_1}(R \bowtie \sigma_{C_2}(S)) = \sigma_{C_1}(R) \bowtie \sigma_{C_2}(S)$



- Join Preorder: need to know the size of table.

- First join similar size tables.

- Cardinality problem



Why so many joins: schema normalize

- Size & Cost Estimation

(pages)

- $B(R)$ number of blocks for relation of R .
- $T(R)$ # of tuples in relation R
- $V(R,A)$ # of distinct values of attribute A If A is key, $V(R,A) = T(R)$

- Size estimation: estimate the size of a logical subplan

Cost estimation: estimate the cost of a physical subplan.

- Estimate size of a query plan $|P|$

▪ Worst case sizes:

$$|D_C(R)| \leq |R|$$

$$|R \bowtie S| \leq |R \times S|$$

- $f \times$ worst case f : selectivity factor

- estimating size of a relation D in columns $D(A_1), \dots, D(A_n)$

- Estimating size of a selection $R_{(A,D)}, S_{(C,U)}$
 - Assumption 1. Uniform distribution of values

$$|D_{A=v}(R)| \approx \frac{|T(R)|}{V(R,A)} \quad f = \frac{1}{V(R,A)}$$

- Assumption 2. Independence of attributes

$$|D_{A=v \text{ and } B=w}(R)| = \frac{|T(R)|}{(V(R,A) \cdot V(R,B))}$$

$\Rightarrow f = f_{A=v} \cdot f_{B=w}$

- Estimating Size of a join

- Assumption 3: Inclusion

if $V(R,B) \leq V(S,C)$ then $\Pi_B(R) \subseteq \Pi_C(S)$

$$|R \bowtie_{B=C} S| \approx \frac{|R| \cdot |S|}{V(S,C)}$$

- In general $|R \bowtie_{B=C} S| \approx \frac{|R| + |S|}{\max(V(R,B), V(S,C))}$

- I/O Cost of selection

- Sequential scan for relation costs $B(R)$

- Index-based selection

▪

- Clustered $B(R) \cdot \frac{1}{V(R,A)}$

- Unclustered $T(R) \cdot \frac{1}{V(R,A)}$

e.g. $B(R) = 2000$, $T(R) = 100,000$, $V(R, A) = 20$.