

Overview

Join algorithms

- (Block) Nested Loop Join - two-level for-loop
- Hash Join - compute a hash table of one input; probe the hash table with the other input
- Sort-Merge Join - sort both tables on one of the join conditions, then merge sorted lists

Indexes

- B-Tree index - supports point and range lookup.
- Hashtable - supports point lookup.

Clustered indexes are the way the data is stored in the "original" table; they are almost always B-trees. Unclustered indexes define secondary tables that reference the main table via pointers.

Cardinality estimation - the problem of estimating the number of tuples after an operation, such as a selection or join. Good estimates are critical to cost modeling; the larger the cardinality; the larger the cost. Here we use assumptions: that the values of a table are uniformly distributed among its distinct values, and that all joins are foreign key-primary key joins.

Cost Modeling

- $B(R)$ - the number of blocks used to store the relation R on disk
- $T(R)$ - the number of tuples in R (also known as R 's cardinality)
- $V(R, a)$ - the number of unique values of attribute a in relation R
- M - the number of pages that fit in memory

Cost-based Query Optimization compares plans by computing their estimated cost, then chooses the one with the cheapest estimated cost to execute.

Query execution - we learned about the iterator method (iterator interface) of executing the operators in a query plan. You may see it called the "pull-based model of query execution", because each operator "pulls" data from its child operators by calling `next()`. The three methods used are `open()`, `next()`, and `close()`.

Formula guide for cardinality estimation

In cost estimation, we assume data is uniformly distributed such that each distinct value has the same number of tuples.

Selectivity factor (X), assuming table $R(a, b)$ cartesian joined $S(a, c)$ and constants $x, x1, x2$:

- $R.a = x \Rightarrow X \cong \frac{1}{V(R,a)}$
- $R.a < x \Rightarrow X \cong \frac{x - \min(R.a)}{\max(R.a) - \min(R.a)}$
- $R.a > x \Rightarrow X \cong \frac{\max(R.a) - x}{\max(R.a) - \min(R.a)}$
- $x1 < R.a < x2 \Rightarrow X \cong \frac{x2 - x1}{\max(R.a) - \min(R.a)}$

- $R.a = S.a$ (equijoin) $\Rightarrow X \cong \frac{1}{\max(V(R,a), V(S,a))}$
- $\text{cond1 AND cond2} \Rightarrow X = X_1 * X_2$

On deriving the selectivity of an equijoin:

Why $R.a = S.a$ (equijoin) $\Rightarrow X \cong \frac{1}{\max(V(R,a), V(S,a))}$?

Let say x_0 a value such that $R.a = S.a = x_0$, that means when we do selection $R.a = x_0$ AND $S.a = x_0$, the selectivity is:

$$X \cong \frac{1}{V(R,a) * V(S,a)}$$

But there can be as many as $\min(V(R,a), V(S,a))$ distinct values of x_0 (for example R has 100 value of a, S has 1000 value of a, the number of value of a after join is 100 because $100 < 1000$, other S.a is filtered out. That means there can be 100 value of x_0 such that $R.a = S.a = x_0$)

Therefore, we multiply the above selectivity by $\min(V(R,a), V(S,a))$ which means the min value is crossed out of the denominator, leaving the maximum value. Thus

$$X \cong \frac{1}{\max(V(R,a), V(S,a))}$$

Note: this is the selectivity factor. To estimate the number of tuples in a join, multiply by the $T_1 T_2$, the number of tuples in a Cartesian product.

X *THE NUMBER OF TUPLES IN YOUR AGGREGATION/TABLE

Problems

1. (Adapted from 414 SP 17 Final)

Consider the relations $R(e, f)$, $S(f, g)$, and $X(g, h)$ in the query plan depicted above.

- Joins are natural joins.
- Every attribute is integer-valued.
- Assume that **every intermediate result is materialized** (i.e., written to disk).
- Assume that we are executing queries on a machine that has **11 memory pages** available.
- Assume **uniform distributions** on the attributes for the purpose of computing estimates.

Consider the following statistics:

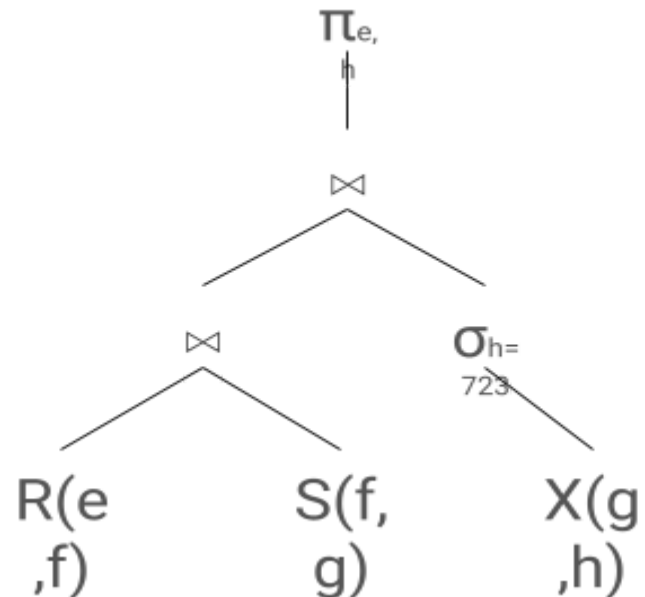


Table	#tuples	#blocks
R	1,000	100
S	5,000	200
X	100,000	10,000

Attribute	# distinct values	Minimum value	Maximum value
R.f	100	1	1,000
S.f	1,000	1	2,000
S.g	5,000	1	2,000
X.g	1,000	1	10,000
X.h	1,000	1	500,000

A. Estimate the number of tuples and blocks in the selection $\sigma_{h=723}(X)$.

Select * from X where h = 723

B. Estimate the number of tuples in the join $R \bowtie S$.

2. (Adapted from 414 AU 19 Final)

Assume we have relations R and S in a database, along with statistics on their attributes as shown below:

R(a integer,
b float
c integer)
S(d integer,
e float)

<u>R</u>	<u>S</u>
T(R) = 1,000	T(S) = 5,000
V(R, a) = 100	V(S, d) = 500
V(R, b) = 1,000	V(S, e) = 5,000
V(R, c) = 10	minimum value of e: 50.0
minimum value of b: 0.0	maximum value of e: 150.0
maximum value of b: 100.0	

$T(X)$ is the number of tuples in a relation X.

$V(X, y)$ is the number of distinct values for the attribute y in the relation X.

We assume the values of the attributes are uniformly distributed over their range.

For each of the queries below, estimate the number of tuples that will be returned in the output.

a)

```
SELECT *  
FROM R  
WHERE R.b > 25.0
```

b)

```
SELECT *  
FROM S  
WHERE S.d = 487
```

c)

```
SELECT *  
FROM S  
WHERE S.e > 70.0 AND S.e < 80.0
```

d)

```
SELECT *  
FROM R  
WHERE R.a = 15 AND R.c = 82
```

e)

```
SELECT *  
FROM R, S
```