CSE 444: Database Internals

Lecture 12 Query Optimization (part 3)

Reminders

- I'm not Magda
- Lab 2 is due on Friday by 11pm
- Lab 3 released this Friday (transactions, yay!)
- HW 5 due next week
- Quiz on 4/25 (next Monday)

Acknowledgments

Today's lecture focuses on how to actually implement the Selinger optimizer

(Many slides from Sam Madden at MIT)

Goal:

• How to order a series joins over N tables A,B,C,...

E.g. A.a = B.b AND A.c = D.d AND B.e = C.f

Goal:

• How to order a series joins over N tables A,B,C,...

E.g. A.a = B.b AND A.c = D.d AND B.e = C.f

Problem:

• ... too ... many ... plans ...

Goal:

- How to order a series joins over N tables A,B,C,...
- E.g. A.a = B.b AND A.c = D.d AND B.e = C.f

Problem:

• N! ways to order joins; e.g. ABCD, ACBD,

• $C_{N-1} = \frac{1}{N} \binom{2(N-1)}{N-1}$ plans/ordering; e.g. (((AB)C)D), ((AB)(CD)))

- Multiple implementations (hash, nested loops)
- Naïve approach does not scale
 - E.g. N = 20, #join orders $20! = 2.4 \times 10^{18}$; many more plans

- Only left-deep plans: (((AB)C)D) eliminate C_{N-1} .
- Push down selections
- Don't consider Cartesian products
- Dynamic programming algorithm

OrderJoins:

R = set of relations to join
For d = 1 to |R|:
For S in {all size-d subsets of R}:
Pick a∈S with lowest cost (S-a)⊠a

OrderJoins:

```
R = set of relations to join
For d = 1 to |R|:
  For S in {all size-d subsets of R}:
     Pick a \in S with lowest cost (S-a)\bowtiea
                   What is the cost?
                     * Cost to scan a
                     * Cost to produce S-a
                     * Cost to join (S-a) with a
```

OrderJoins:

```
R = set of relations to join
For d = 1 to |R|:
  For S in {all size-d subsets of R}:
     Pick a \in S with lowest cost (S-a)\bowtiea
                    What is the cost?
                      * Cost to scan a
                      * Cost to produce S-a ← Calculated in previous iteration
                      * Cost to join (S-a) with a
```

OrderJoins:

```
R = set of relations to join
For d = 1 to |R|:
For S in {all size-d subsets of R}:
optjoin(S) = (S - a) join a,
where a is the single relation that minimizes:
cost(optjoin(S - a)) +
min.cost to join (S - a) with a +
min.access cost for a
```

Note: **optjoin**(S-a) is cached from previous iterations

Subplan S	optJoin(S)	Cost(OptJoin(S))
А		

- orderJoins(A, B, C, D)
- Assume all joins are NL

- orderJoins(A, B, C, D)
- Assume all joins are NL

Subplan S	optJoin(S)	Cost(OptJoin(S))
А	Index scan	100
В	Seq. scan	50
С	Seq scan	120
D	B+tree scan	400

- d = 1
- A = best way to access A
 - (sequential scan, predicate-pushdown on index, etc)
 - B = best way to access B
 - C = best way to access C
 - D = best way to access D
- Total number of steps: choose(N, 1)

- orderJoins(A, B, C, D)
- d = 2
- {A,B} = AB or BA
- use previously computed
 - best way to access A and B

Subplan S	optJoin(S)	Cost(OptJoin(S))
A	Index scan	100
В	Seq. scan	50

- orderJoins(A, B, C, D)
- d = 2
- {A,B} = AB or BA
- use previously computed
 - best way to access A and B

-		
Subplan S	optJoin(S)	Cost(OptJoin(S))
А	Index scan	100
В	Seq. scan	50
{A, B}	BA	156

- orderJoins(A, B, C, D)
- d = 2
- {A,B} = AB or BA
- use previously computed
 - best way to access A and B
 - $\{B,C\} = BC \text{ or } CB$

Subplan S	optJoin(S)	Cost(OptJoin(S))
А	Index scan	100
В	Seq. scan	50
{A, B}	BA	156
{B, C}	BC	98

• orderJoins(A, B, C, D)

	Subplan S	optJoin(S)	Cost(OptJoin(S))
)	А	Index scan	100
	В	Seq. scan	50
	{A, B}	BA	156
	{B, C}	BC	98
nd B			

- d = 2
- {A,B} = AB or BA
- use previously computed
 - best way to access A and B
 - $\{B,C\} = BC \text{ or } CB$

• orderJoins(A, B, C, D)

	Subplan S	optJoin(S)	Cost(OptJoin(S))
	A	Index scan	100
	В	Seq. scan	50
	{A, B}	BA	156
	{B, C}	BC	98
3			

- d = 2
- {A,B} = AB or BA
- use previously computed
 - best way to access A and B
 - $\{B,C\} = BC \text{ or } CB$
 - $\{C,D\} = CD \text{ or } DC$
 - {A,C} = AC or CA
 - $\{B,D\} = BD \text{ or } DB$
 - $\{A,D\} = AD \text{ or } DA$

- orderJoins(A, B, C, D)
- Subplan S Cost(OptJoin(S)) optJoin(S) Α Index scan 100 Β 50 Seq. scan . . . **{A**, B} BA 156 {B, C} BC 98

- d = 2
- {A,B} = AB or BA
- use previously computed
 - best way to access A and B
 - $\{B,C\} = BC \text{ or } CB$
 - $\{C,D\} = CD \text{ or } DC$
 - {A,C} = AC or CA
 - {B,D} = BD or DB
 - $\{A,D\} = AD \text{ or } DA$
- Total number of steps: choose(N, 2) × 2

• orderJoins(A, B, C, D)

• d = 3

 $\{A,B,C\} =$ Remove A: compare A($\{B,C\}$) to ($\{B,C\}$)A

Subplan S	optJoin(S)	Cost(OptJoin(S))
А	Index scan	100
В	Seq. scan	50
{A, B}	BA	156
{B, C}	BC	98

• orderJoins(A, B, C, D)

• d = 3

 $\{A,B,C\} =$ Remove A: compare A($\{B,C\}$) to ($\{B,C\}$)A

Subplan S		optJoin(S)	Cost(OptJoin(S))
A		Index scan	100
В		Seq. scan	50
{A, B}		BA	156
{B, C}		BC	98

optJoin(B,C) and its cost are already cached in table

Example Subplan S

А

В

. . . .

orderJoins(A, B, C, D)

• d = 3

{A, B} BA {B, C} BC {A, B, C} BAC ${A,B,C} =$ Remove A: compare A($\{B,C\}$) to ($\{B,C\}$)

Remove B: compare $B(\{A,C\})$ to $(\{A,C\})$ B

Remove C: compare C({A,B}) to ({A,B})C

optJoin(B,C) and its cost are already cached in table

Cost(OptJoin(S))

100

50

156

98

500

optJoin(S)

Index scan

Seq. scan





24

Α

B

{A, B}

{B, C}

{A, B, C}

{B, C, D}

Subplan S

optJoin(S)

Seq. scan

Index

scan

BA

BC

BAC

DBC

Cost(OptJoin(S))

optJoin(B, C, D)

and its cost are

already cached

in table

100

50

156

98

500

150

• orderJoins(A, B, C, D)

- d = 4
- {A,B,C,D} =
- •
- ullet
- Remove A: compare A({B, C, D}) to ({B, C, D})
- Remove B: compare B({A,C,D}) to ({A,C,D})B
- Remove C: compare C({A,B,D}) to ({A,B,D})C
 - Remove D: compare D({A,B,C}) to ({A,B,C})D
- Total number of steps: choose(N, 4) \times 4 \times 2

Complexity

• Total #subsets considered

- Choose(N, 1) + Choose(N, 2) + + Choose (N, N)
- All nonempty subsets of a size N set: $2^{N} 1$
- Equivalently: number of binary strings of size N, except 00...0:
- 000, 001, 010, 011, 100, 101, 110, 111

Complexity

• Total #subsets considered

- Choose(N, 1) + Choose(N, 2) + + Choose (N, N)
- All nonempty subsets of a size N set: $2^{N} 1$
- Equivalently: number of binary strings of size N, except 00...0:
- 000, 001, 010, 011, 100, 101, 110, 111
- For each subset of size d:
 - d ways to remove one element
 - 2 ways for compute AB or BA

Complexity

• Total #subsets considered

- Choose(N, 1) + Choose(N, 2) + + Choose (N, N)
- All nonempty subsets of a size N set: $2^{N} 1$
- Equivalently: number of binary strings of size N, except 00...0:
- 000, 001, 010, 011, 100, 101, 110, 111
- For each subset of size d:
 - d ways to remove one element
 - 2 ways for compute AB or BA
- Total #plans considered
 - Choose(N, 1) + 2 Choose(N, 2) + + N Choose (N, N)
 - Equivalently: total number of 1's in all strings of size N
 - N 2^{N-1} because every 1 occurs 2^{N-1} times
 - Need to further multiply by 2, to account for AB or BA

Why Left-Deep

Asymmetric, cost depends on the order

- Left: Outer relation Right: Inner relation
- For nested-loop-join, we try to load the outer (typically smaller) relation in memory, then read the inner relation one page at a time B(R) + B(R)*B(S) or B(R) + B(R)/M * B(S)
- For index-join, we assume right (inner) relation has index

Why Left-Deep

- Advantages of left-deep trees?
 - 1. Fits well with standard join algorithms (nested loop, one-pass), more efficient
 - 2. One pass join: Uses smaller memory
 - 1. ((R, S), T), can reuse the space for R while joining (R, S) with T
 - 2. (R, (S, T)): Need to hold R, compute (S, T), then join with R, worse if more relations
 - 3. Nested loop join, consider top-down iterator next()
 - 1. ((R, S), T), Reads the chunks of (R, S) once, reads stored base relation T multiple times
 - (R, (S, T)): Reads the chunks of R once, reads computed relation (S, T) multiple times, either more time or more space

Excruciatingly Detailed Optimization Example

Interesting Orders

- Some query plans produce data in sorted order
 - E.g scan over a primary index, merge-join
 - Called *interesting order*
- Next operator may use this order
 - E.g. can be another merge-join
- For each subset of relations, compute multiple optimal plans, one for each interesting order
- Increases complexity by factor k+1, where k=number of interesting orders