Updated 3/29/2022

Topic 3: Efficient query evaluation Unit 1: Acyclic query evaluation (continued) Lecture 19

Wolfgang Gatterbauer

CS7240 Principles of scalable data management (sp22)

https://northeastern-datalab.github.io/cs7240/sp22/

3/29/2022

Pre-class conversations

- Current topic: trees vs cycles
- Scribes, projects, keep on commenting on slides
- Today:
 - Full Semi-join reducers
 - Cycles: tree decompositions, Linear Programming Duality



Source: Moshe Y. Vardi. "A Logical Revolution", 2013. <u>https://www.cs.rice.edu/~vardi/comp409/logicrevo13.pdf</u>, <u>https://www.youtube.com/watch?v=7pu-ZJddxJQ&t=59m20s</u> Wolfgang Gatterbauer. Principles of scalable data management: <u>https://northeastern-datalab.github.io/cs7240/</u>

Outline: T3-1: Acyclic conjunctive queries

• T3-1: Acyclic conjunctive queries

- The semijoin operator
- Join trees & Yannakakis algorithm
- Query hypergraphs & GYO reduction
- A detailed Yannakakis example
- Full semijoin reductions
- T3-2: Cyclic conjunctive queries

Semijoin Reducer



 $Q(x,y,z) = R(x,y) \bowtie S(y,z) \bowtie T(z,w)$

DEFINITION: A full semijoin reducer (semijoin program) is:

a sequence of semijoins on the join tree $R'(x,y) = R(y,z) \ltimes ...$ $S'(y,z) = S(y,z) \ltimes ...$

...

s.t. there no more "dangling tuples" in the reduced relations

Then you can rewrite the query over the reduced relations $Q(x,y,z) = R'(x,y) \bowtie S'(y,z) \bowtie T'(z,w)$



A full reducer is

?

A full reducer is





1. Find a join tree



A full reducer is



3. collect at root (= bottom-up)

$$S_1(y, z) = S(y, z) \ltimes R(x, y)$$

$$T_1(z, y) = T(z, y) \ltimes S_1(y, z)$$



1. Find a join tree



A full reducer is

$$[R(x,y)] \xrightarrow{\bowtie} S(y,z) \xrightarrow{\bowtie} T(z,w)$$

3. collect at root (= bottom-up)

$$S_{1}(y,z) = S(y,z) \ltimes R(x,y)$$
4. distribute

$$T_{1}(z,y) = T(z,y) \ltimes S_{1}(y,z)$$

$$S_{2}(z,y) = S_{1}(y,z) \ltimes T_{1}(z,y)$$

$$R_{1}(x,y) = R(x,y) \ltimes S_{2}(y,z)$$

1. Find a join tree



5. The rewritten query is



Semijoin Reducer

$$Q(x,y,z) = R(x,y) \bowtie S(y,z) \bowtie T(z,w)$$

A full reducer is $R(x,y) \xrightarrow{\times} S(y,z) \xrightarrow{\times} T(z,w)$ 3. collect at root (= bottom-up) $S_1(y,z) = S(y,z) \ltimes R(x,y)$ 4. distribute $T_1(z,y) = T(z,y) \ltimes S_1(y,z)$ to leaves $S_2(z,y) = S_1(y,z) \ltimes T_1(z,y)$ $R_1(x,y) = R(x,y) \ltimes S_2(y,z)$

5. The rewritten query is $Q(x,y,z) = R_1(x,y) \bowtie S_2(y,z) \bowtie T_1(z,w)$



1. Find a join tree



Semijoin Reduction



The intuition of the sequence of semi-join reductions is as follows:

- a message is sent across an edge only if the subtree is already fully reduced
- a message that is sent across an edge contains *all the information* from the subtree rooted at the sender.

It follows that the reduction always needs to start at the leaves!

What happens if one uses the wrong sequence? (1) $\gtrsim \times \le$

T N2

RKS

GYO ear removal

- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

Q(x,y,z) := R(x,y), S(y,z), T(x,z).





GYO ear removal

- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

Q(x,y,z) := R(x,y), S(y,z), T(x,z).







GYO ear removal

- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

Q(x,y,z) := R(x,y), S(y,z), T(x,z).





There is no join tree! You can't fulfill the running intersection property...



GYO ear removal

- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

Q(x,y,z) := R(x,y), S(y,z), T(x,z), W(x,y,z).





GYO ear removal

- remove isolated nodes (variables)
- remove consumed or empty edges (atoms)

Q(x,y,z) := R(x,y), S(y,z), T(x,z), W(x,y,z).

Join tree







Q(x,y,z) := R(x,y), S(y,z), T(x,z), W(x,y,z).



Write 1) a full reducer and then 2) the new join expression in RA





Q(x,y,z) := R(x,y), S(y,z), T(x,z), W(x,y,z).



$$W_{1}(x, y, z) = W(x, y, z) \ltimes R(x, y)$$
$$W_{2}(x, y, z) = W_{1}(x, y, z) \ltimes S(y, z)$$
$$W_{3}(x, y, z) = W_{2}(x, y, z) \ltimes T(x, z)$$
$$R_{1}(x, y) = R(x, y) \ltimes W_{3}(x, y, z)$$
$$S_{1}(y, z) = S(y, z) \ltimes W_{3}(x, y, z)$$
$$T_{1}(x, z) = T(x, z) \ltimes W_{3}(x, y, z)$$

 $Q(x,y,z) = \left(\left(R_1(x,y) \bowtie S_1(y,z) \right) \bowtie T_1(x,z) \right) \bowtie W_3(x,y,z)$



Q(x,y,z) := R(x,y), S(y,z), T(x,z), W(x,y,z).





Q(x,y,z) := R(x,y), S(y,z), T(x,z), W(x,y,z).



$$W_{1}(x, y, z) = W(x, y, z) \ltimes R(x, y)$$
$$W_{2}(x, y, z) = W_{1}(x, y, z) \ltimes S(y, z)$$
$$W_{3}(x, y, z) = W_{2}(x, y, z) \ltimes T(x, z)$$
$$R_{1}(x, y) = R(x, y) \ltimes W_{3}(x, y, z)$$
$$S_{1}(y, z) = S(y, z) \ltimes W_{3}(x, y, z)$$
$$T_{1}(x, z) = T(x, z) \ltimes W_{3}(x, y, z)$$

 $Q(x,y,z) = \left(\left(R_1(x,y) \bowtie S_1(y,z) \right) \bowtie T_1(x,z) \right) \bowtie W_3(x,y,z)$

O

0



Q(x,y,z) := R(x,y), S(y,z), T(x,z), W(x,y,z).



Semi-join reductions can be extremely powerful

5040



(h) Combining (a)-(c)

Optimization is speculative (i.e. depends on the inputs

Opt. 3. Deterministic semi-join reduction

The most expensive operations in probabilistic query plans are the group-bys for the probabilistic project operations. These are often applied early in the plans to tuples which are later pruned and do not contribute to the final query result. Our third optimization is to first apply a *full semi-join reduction on the input relations* before starting the probabilistic evaluation from these reduced input relations.

We like to draw here an important connection to [54], which introduces the idea of "lazy plans" and shows orders of magnitude performance improvements for safe plans by computing confidences not after each join and projection, but rather at the very end of the plan. We note that our semijoin reduction serves the same purpose with similar performance improvements and also apply for safe queries. The advantage of semi-join reductions, however, is that we do not require any modifications to the query engine.

From: Gatterbauer, Suciu. "Dissociation and propagation for approximate lifted inference with standard relational database management systems", VLDBJ 2017. https://arxiv.org/pdf/1310.6257 Reference [54]: Olteanu, Huang, Koch. "Sprout: Lazy vs. eager query plans for tuple-independent probabilistic databases", ICDE 2009. https://doi.org/10.1109/ICDE.2009.123 179 Wolfgang Gatterbauer. Principles of scalable data management: https://northeastern-datalab.github.io/cs7240/