CSE 544 Theory of Query Languages

Announcements

- Project Milestone: was due on Friday
 Expect feedback by the end of this week
- Project presentations: Friday, 6/7
 Reserve all day, stay tuned for announcements
- Homework 3: was due yesterday
 Did you remember to turn off your servers!?!
- Homework 4: will be posted in 2-3 days
- Next paper review: next Wednesday, 5/29

Complexity Classes

A decision problem:

- We have a property (a.k.a. problem)
- Given an input X of size n, decide if it satisfies the property

In other words, define have to compute a function f(X) = 0 or 1

We are interested in these classes



The Class AC⁰

What is AC⁰?

The Class AC⁰

A problem f is in AC⁰ if, for every n, there exists a Boolean circuit s.t.:

- It consists of unbounded fan-in AND, unbounded fan-in OR, and NOT gates
- If the inputs X₁, ..., X_n encode an input X, then the circuit's output is f(X)
- The circuit size is $n^{O(1)}$
- The circuit depth is O(1)

Example in AC⁰

<u>Problem</u>: given an input string $X_1, ..., X_n$ in $\{0,1\}^n$, check if it has at least two 1's

0100101 - yes

00001000 - no

Example in AC⁰

<u>Problem</u>: given an input string $X_1, ..., X_n$ in $\{0,1\}^n$, check if it has at least two 1's

0100101 - yes

00001000 - no



Are these in AC⁰ or not?

Carry bit: The sum $(X_nX_{n-1}...X_1) + (Y_nY_{n-1}...Y_1)$ has a carry bit

(1001)+(0101) = (1110) NO (1001)+(0111) = (10000) YES

Triangle: A graph given by the n × n adjacency matrix contains a directed triangle



Parity: X_1, \dots, X_n has an even number of 1's

10010101101 YES 10010101100 NO

s-t Reachability (GAP): A graph given by the n × n adjacency matrix contains a path from node s to node t



Are these in AC⁰ or not?



Relational Queries

<u>Theorem</u>: Every Boolean relational query defines a property that is in AC⁰



select distinct R.A, S.C from R, S where R.B=S.B

Example

R	А	В
	а	а
	b	b
	а	b

S	В	С
	С	с
	b	с



























Another Example

 $Q=\exists y.R(a',y) \land (\forall z.S(y,z) \rightarrow \exists u.R(z,u))$

Practice at home: Show that Q is in AC⁰ by showing how to construct a circuit for computing Q. What is the depth? What fanouts have your OR and AND gates ?

Discussion

In class: make sure you understand very well why every relational query is in AC⁰

- Consequence 1 (for theoreticians and their friends):
 - SELECT-DISTINCT-FROM-WHERE queries cannot express PARITY, GAP
- Consequence 2 (for fans of Big Data)
 "SQL is embarrassingly parallel"

The Classes L and NL



LOGSPACE and NLOGSPACE

• What is LOGSPACE (or L)?

• What is NLOGSPACE (or NL)?

LOGSPACE and NLOGSPACE

- A problem is in LOGSPACE (or L) if it can be computed by a deterministic Turing machine using O(log n) space
- A problem is in NLOGSPACE (or NL) if it can be computed by a non-deterministic Turing machine using O(log n) space

O(log n) space refers to the working tape: the input is on a separate tape of size n

Examples

- GAP is in NLOGSPACE (why?)
- 1-GAP (each node has outdegree ≤ 1) is in LOGSPACE (why?)
- Recall that none of these problems is in AC⁰
- Theorem: GAP is complete for NLOGSPACE
- Theorem: 1-GAP is complete for LOGSPACE

It can express GAP

In many ways!!

T(x,y) := R(x,y)T(x,y) :- R(x,z), T(z,y) Answer :- T('s','t')

This proves that datalog can express strictly more queries than the relational calculus

The Classes PTIME, NP, PSPACE



The Classes PTIME, NP, PSPACE

• What is **PTIME**?

• What is NP?

In class

• What is **PSPACE**?

What is the Complexity of Datalog?



The Same-Generation Problem

Problem: We have a database of microbes, where each microbe x may have several children y:

Parent(x,y)

Find all microbes in the same generation with "M62251"

Can we solve it in datalog?

Is this problem in NLOGSPACE?

The Same-Generation Problem

Problem: We have a database of microbes, where each microbe x may have several children y:

Parent(x,y)

Find all microbes in the same generation with "M62251"



30

The Same-Generation Problem

Problem: We have a database of microbes, where each microbe x may have several children y:

Parent(x,y)

Find all microbes in the same generation with "M62251"



Discussion

- The same-generation problem was a trap:
 SG is no more complex than GAP!
 - Lesson: GAP is more than meets the eyes

 But datalog <u>is</u> more expressive than NLOGSPACE: it captures all of PTIME, in ways we discuss next, in three steps

Step 1: Complexity of Datalog

Theorem. Datalog is in PTIME.

More precisely, fix any Boolean datalog program P. The problem: given D, check if P(D) = true is in PTIME

Proof: ... [discuss in class]

Step 1: Complexity of Datalog

Theorem. Datalog is in PTIME.

More precisely, fix any Boolean datalog program P. The problem: given D, check if P(D) = true is in PTIME

Proof: ... [discuss in class]

Which of the following are in PTIME? Stratified, inflationary-fixpoint, partial-fixpoint datalog[¬].

Circuit Value Problem

Input = a rooted DAG; leaves labeled 0/1, internal nodes labeled AND/OR **Output** = check if the value of the root is 1

Note: NOT nodes could be added w.l.o.g. (why?)



Circuit Value Problem

Input = a rooted DAG; leaves labeled 0/1, internal nodes labeled AND/OR **Output** = check if the value of the root is 1

Note: NOT nodes could be added w.l.o.g. (why?)

Theorem.

The Circuit Value Problem is complete for PTIME

In class:

- 1. How can we compute it in PTIME?
- 2. Why isn't it in NLOGSPACE?



Theorem. Datalog can express the Circuit Value Problem

EDBs: root(x) and(x,y1,y2) or(x,y1,y2) zeroLeaf(x)

oneLeaf(x)

Theorem. Datalog can express the Circuit Value Problem



Discussion

- Step 1: datalog is in PTIME
 - Stratified, and inflationary datalog[¬] are in PTIME
- Step 2: datalog can express a PTIME complete problem
- Step 3: can datalog express <u>all</u> PTIME problems?

Step 3

<u>Theorem</u>. For every problem in PTIME there exists a program in inflationary-fixpoint datalog[¬] that expresses that problem

Caveat: the program must have access to a total order on the active domain. Otherwise inflationary-fixpoint datalog[¬] cannot even express parity!

> Make sure you understand why pure datalog (without negation) cannot express all of PTIME

Finally: partial-fixpoint Datalog[¬]

<u>Theorem</u>. Partial-fixpoint datalog[¬] can express precisely the problems that are in PSPACE

Same caveat: for completeness we need access to an order relation

Which are "Easy" to Parallelize?

Relational calculus = AC⁰

• Add transitive closure = NLOGSPACE

Inflationary datalog = PTIME

Which are "Easy" to Parallelize?

- Relational calculus = AC⁰
- YES!! "embarrassingly parallel"

- Add transitive closure = NLOGSPACE
- MAYBE: the path-doubling program

- Inflationary datalog = PTIME
- NO: circuit value problem

Descriptive Complexity

- In computational complexity one describes complexity classes in terms of a computational model
 - Turing Machine, circuit, etc
- In descriptive complexity one describes complexity classes in terms of the logic ("query language") that captures that class

Descriptive Complexity

[Immerman, Vardi] Assume we have access to an order relation < (and to a BIT relation for AC⁰)

 $RC = AC^0$

RC + Transitive Closure = NL

Inflationary datalog[¬] = PTIME

Partial fixpoint datalog[¬] = PSPACE