

CSE 544: Lecture 11 Theory

Monday, May 3, 2004

Query Minimization

Definition A conjunctive query q is minimal if for every other conjunctive query q' s.t. $q \equiv q'$, q' has at least as many predicates ('subgoals') as q

Are these queries minimal ?

```
q(x) :- R(x,y), R(y,z), R(x,x)
```

```
q(x) :- R(x,y), R(y,z), R(x,'Alice')
```

Query Minimization

- Query minimization algorithm

Choose a subgoal g of q

Remove g : let q' be the new query

We already know $q \subseteq q'$ (why ?)

If $q' \subseteq q$ then permanently remove g

- Notice: the order in which we inspect subgoals doesn't matter

Query Minimization In Practice

- No database system today performs minimization !!!
- Reason:
 - It's hard (NP-complete)
 - Users don't write non-minimal queries
- However, non-minimal queries arise when using views intensively

Query Minimization for Views

```
CREATE VIEW HappyBoaters
SELECT DISTINCT E1.name, E1.manager
FROM Employee E1, Employee E2
WHERE E1.manager = E2.name
and E1.boater='YES'
and E2.boater='YES'
```

This query is minimal

Query Minimization for Views

Now compute the Very-Happy-Boaters

```
SELECT DISTINCT H1.name
FROM HappyBoaters H1, HappyBoaters H2
WHERE H1.manager = H2.name
```

This query is also minimal

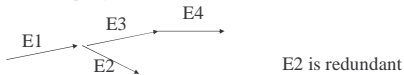
What happens in SQL when we run a query on a view ?

Query Minimization for Views

View Expansion

```
SELECT DISTINCT E1.name
FROM Employee E1, Employee E2, Employee E3, Employee E4
WHERE E1.manager = E2.name and E1.boater = 'YES' and E2.boater = 'YES'
and E3.manager = E4.name and E3.boater = 'YES' and E4.boater = 'YES'
and E1.manager = E3.name
```

This query is no longer minimal !



Monotone Queries

Definition A query q is monotone if:
For every two databases D, D'
if $D \subseteq D'$ then $q(D) \subseteq q(D')$

Which queries below are monotone ?

$\phi \equiv \exists x.R(x,x)$

$\phi \equiv \exists x.\exists y.\exists z.\exists u.(R(x,y) \wedge R(y,z) \wedge R(z,u))$

$\phi \equiv \exists x.\forall y.R(x,y)$

Monotone Queries

- **Theorem.** Every conjunctive query is monotone
- Stronger: every UCQ query is monotone

How To Impress Your Students Or Your Boss

- Find all drinkers that like some beer that is not served by the bar "Black Cat"

```
SELECT L.drinker
FROM Likes L
WHERE L.beer not in (SELECT S.beer
                     FROM Serves S
                     WHERE S.bar = 'Black Cat')
```

- Can you write as a simple SELECT-FROM-WHERE (I.e. without a subquery) ?

Expressive Power of FO

- The following queries cannot be expressed in FO:
 - Transitive closure:
 - $\forall x.\forall y$, there exists x_1, \dots, x_n s.t.
 $R(x,x_1) \wedge R(x_1,x_2) \wedge \dots \wedge R(x_{n-1},x_n) \wedge R(x_n,y)$
 - Parity: the number of edges in R is even

Datalog

- Adds recursion, so we can compute transitive closure
- A datalog program (query) consists of several datalog rules:

```
P1(t1) :- body1
P2(t2) :- body2
⋮
Pn(tn) :- bodyn
```

Datalog

Terminology:

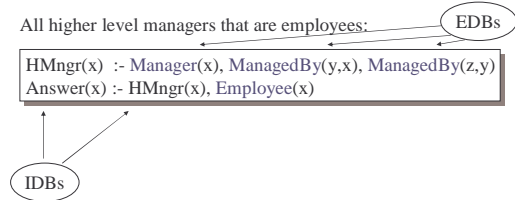
- EDB = extensional database predicates
 - The database predicates
- IDB = intentional database predicates
 - The new predicates constructed by the program

Datalog

Employee(x), ManagedBy(x,y), Manager(y)

All higher level managers that are employees:

$HMngr(x) :- Manager(x), ManagedBy(y,x), ManagedBy(z,y)$
 $Answer(x) :- HMngr(x), Employee(x)$



Datalog

Employee(x), ManagedBy(x,y), Manager(y)

All persons:

$Person(x) :- Manager(x)$
 $Person(x) :- Employee(x)$

Manager \cup Employee

Unfolding non-recursive rules

Graph: $R(x,y)$

$P(x,y) :- R(x,u), R(u,v), R(v,y)$
 $A(x,y) :- P(x,u), P(u,y)$

Can “unfold” it into:

$A(x,y) :- R(x,u), R(u,v), R(v,w), R(w,m), R(m,n), R(n,y)$

Unfolding non-recursive rules

Graph: $R(x,y)$

$P(x,y) :- R(x,y)$
 $P(x,y) :- R(x,u), R(u,y)$
 $A(x,y) :- P(x,y)$

Now the unfolding has a union:

$A(x,y) :- R(x,y) \vee \exists u(R(x,u) \wedge R(u,y))$

Recursion in Datalog

Graph: $R(x,y)$

Transitive closure:

$P(x,y) :- R(x,y)$
 $P(x,y) :- P(x,u), R(u,y)$

Transitive closure:

$P(x,y) :- R(x,y)$
 $P(x,y) :- P(x,u), P(u,y)$

Recursion in Datalog

Boolean trees:

Leaf0(x), Leaf1(x),
AND(x, y₁, y₂), OR(x, y₁, y₂),
Root(x)

- Write a program that computes:
Answer() :- true if the root node is 1

Recursion in Datalog

One(x)	:- Leaf1(x)
One(x)	:- AND(x, y ₁ , y ₂), One(y ₁), One(y ₂)
One(x)	:- OR(x, y ₁ , y ₂), One(y ₁)
One(x)	:- OR(x, y ₁ , y ₂), One(y ₂)
Answer()	:- Root(x), One(x)

Exercise

Boolean trees:

Leaf0(x), Leaf1(x),
AND(x, y₁, y₂), OR(x, y₁, y₂), Not(x,y),
Root(x)

- **Hint:** compute both One(x) and Zero(x) here you need to use Leaf0

Variants of Datalog

	without recursion	with recursion
without \neg	Non-recursive Datalog = UCQ (why ?)	Datalog
with \neg	Non-recursive Datalog [¬] = FO	Datalog [¬]

Non-recursive Datalog

- Union of Conjunctive Queries = UCQ
 - Containment is decidable, and NP-complete
- Non-recursive Datalog
 - Is equivalent to UCQ
 - Hence containment is decidable here too
 - Is it still NP-complete ?

Non-recursive Datalog

- A non-recursive datalog:

T ₁ (x,y)	:- R(x,u), R(u,y)
T ₂ (x,y)	:- T ₁ (x,u), T ₁ (u,y)
...	...
T _n (x,y)	:- T _{n-1} (x,u), T _{n-1} (u,y)
Answer(x,y)	:- T _n (x,y)

- Its unfolding as a CQ:

Answer(x,y)	:- R(x,u ₁), R(u ₁ , u ₂), R(u ₂ , u ₃), ... R(u _n , y)
-------------	--

- How big is this query ?

Query Complexity

- Given a query ϕ in FO
- And given a model $\mathbf{D} = (D, R_1^D, \dots, R_k^D)$
- What is the complexity of computing the answer $\phi(\mathbf{D})$

Query Complexity

Vardi's classification:

Data Complexity:

- Fix ϕ . Compute $\phi(\mathbf{D})$ as a function of $|\mathbf{D}|$

Query Complexity:

- Fix \mathbf{D} . Compute $\phi(\mathbf{D})$ as a function of $|\phi|$

Combined Complexity:

- Compute $\phi(\mathbf{D})$ as a function of $|\mathbf{D}|$ and $|\phi|$

Which is the most important in databases ?

Example

$$\phi(x) \equiv \exists u.(R(u,x) \wedge \forall y.(S(y,v) \Rightarrow \neg R(x,y)))$$

R =

3	8
7	5
0	8
09	7
6	9
7	6
89	8
98	7
4	0

S =

43	4
5	58
8	6
9	79
6	67
4	7
6	8

How do we proceed ?

General Evaluation Algorithm

for every subexpression ϕ_i of ϕ , ($i = 1, \dots, m$)
 compute the answer to ϕ_i as a table $T_i(x_1, \dots, x_n)$
 return T_m

Theorem. If ϕ has k variables then one can compute $\phi(\mathbf{D})$ in time $O(|\phi|^*|\mathbf{D}|^k)$

Data Complexity = $O(|\mathbf{D}|^k)$ = in PTIME
 Query Complexity = $O(|\phi|^*c^k)$ = in EXPTIME

General Evaluation Algorithm

Example: $\phi(x) \equiv \exists u.(R(u,x) \wedge \forall y.(S(y,v) \Rightarrow \neg R(x,y)))$

$\phi_1(u,x)$	$\equiv R(u,x)$
$\phi_2(y,v)$	$\equiv S(y,v)$
$\phi_3(x,y)$	$\equiv \neg R(x,y)$
$\phi_4(y)$	$\equiv \exists v.\phi_2(y,v)$
$\phi_5(x,y)$	$\equiv \phi_4(y) \Rightarrow \phi_3(x,y)$
$\phi_6(x)$	$\equiv \forall y.\phi_5(x,y)$
$\phi_7(u,x)$	$\equiv \phi_1(u,x) \wedge \phi_6(x)$
$\phi_8(x)$	$\equiv \exists u.\phi_7(u,x) \equiv \phi(x)$

Complexity

Theorem. If ϕ has k variables then one can compute $\phi(\mathbf{D})$ in time $O(|\phi|^*|\mathbf{D}|^k)$

Remark. The number of variables matters !

Paying Attention to Variables

- Compute all chains of length m

$$\text{Chain}_m(x,y) \text{ :- } R(x,u_1), R(u_1, u_2), R(u_2, u_3), \dots, R(u_{m-1}, y)$$

- We used m+1 variables
- Can you rewrite it with fewer variables ?

Counting Variables

- $\text{FO}^k = \text{FO}$ restricted to variables x_1, \dots, x_k
- Write Chain_m in FO^3 :

$$\text{Chain}_m(x,y) \text{ :- } \exists u.R(x,u) \wedge (\exists x.R(u, x) \wedge (\exists u.R(x,u) \dots \wedge (\exists u.R(u, y) \dots)))$$

Query Complexity

- Note: it suffices to investigate boolean queries only
 - If non-boolean, do this:

```

for  $a_1$  in D, ...,  $a_k$  in D
if ( $a_1, \dots, a_k$ ) in  $\phi(D)$  /* this is a boolean query */
then output ( $a_1, \dots, a_k$ )
    
```

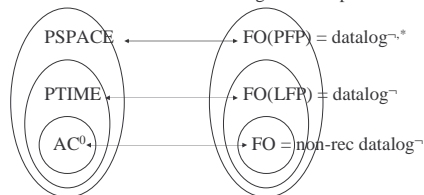
Computational Complexity Classes

Recall computational complexity classes:

- AC^0
- LOGSPACE
- NLOGSPACE
- PTIME
- NP
- PSPACE
- EXPTIME
- EXPSPACE
- (Kalmar) Elementary Functions
- Turing Computable functions

Data Complexity of Query Languages

Paper: On the Unusual Effectiveness of Logic in Computer Science



Important: the more complex a QL, the harder it is to optimize

Views

Employee(x), ManagedBy(x,y), Manager(y)

Views {

$$L(x,y) \text{ :- } \text{ManagedBy}(x,u), \text{ManagedBy}(u,y)$$

$$E(x,y) \text{ :- } \text{ManagedBy}(x,y), \text{Employee}(y)$$

Query {

$$Q(x,y) \text{ :- } \text{ManagedBy}(x,u), \text{ManagedBy}(u,v), \text{ManagedBy}(v,w), \text{ManagedBy}(w,y), \text{Employee}(y)$$

How can we answer Q if we only have L and E ?

Views

- Query rewriting using views (when possible):

$Q(x,y) :- L(x,u), L(u,y), E(v,y)$

- Query answering:
 - Sometimes we cannot express it in CQ or FO, but we can still answer it

Views

Applications:

- Using advanced indexes
- Using replicated data
- Data integration [Ullman'99]