

Query Evaluation on Probabilistic Databases

CSE 544: Wednesday, May 24, 2006

Problem Setting

Tables:

Review

name	rating	p
Monkey Love	good	.5
	fair	.2
	fair	.6
	poor	.9

Movie

title	year	p
Twelve Monkeys	1995	.8
Monkey Love	1997	.4
Monkey Love	1935	.9
Monkey Love Pl	2005	.7

Problem: complexity of query evaluation

Queries:

```
A(x,y) :- Review(x,y),  
         Movie(x,z), z > 1991
```

Answers:

title	rating	p
Twelve Monkeys	fair	.53
Monkey Love	good	.42
Monkey Love Pl	fair	.15

Top k



Two Problems

Fixed schema S , conjunctive query $Q(x,y)$

Query evaluation problem

Fix answer tuple (a,b)

Given database I , compute $\Pr(Q(a,b))$

Top- k answering problem

Fix $k > 0$

Given database I , find k answer tuples with highest probabilities

Related Work: DB

- Cavallo&Pitarelli:1987
- Barbara,Garcia-Molina, Porter:1992
- Lakshmanan,Leone,Ross&Subrahmanian:1997
- Fuhr&Roellke:1997
- Dalvi&S:2004
- Widom:2005

Related Work: Logic

- Query reliability [Gradel,Gurevitch,Hirsch'98]
- Degrees of belief
[Bacchus,Grove,Halpern,Koller'96]
- Probabilistic Logic [Nielson]
- Probabilistic model checking [Kwiatkowska'02]
- Probabilistic Relational Model
[Taskar,Abbeel,Koller'02]

Probabilistic Database

Schema S , Domain D , Set of instances $Inst$

Definition

Probabilistic database is a probability distribution

$$Pr : Inst \rightarrow [0,1], \quad \sum_I Pr[I] = 1$$

If $Pr[I] > 0$ then I is called "possible world"

Probabilistic Database

Representation:

- Independent tuples:
I-database DB over some schema S^i
- Independent and disjoint tuples:
ID-database DB over some schema S^{id}



Semantics:

- DB "means" probability distribution Pr over schema S



Independent Events

- A tuple is in the database with probability p
- Any two tuples are independent events

Representation

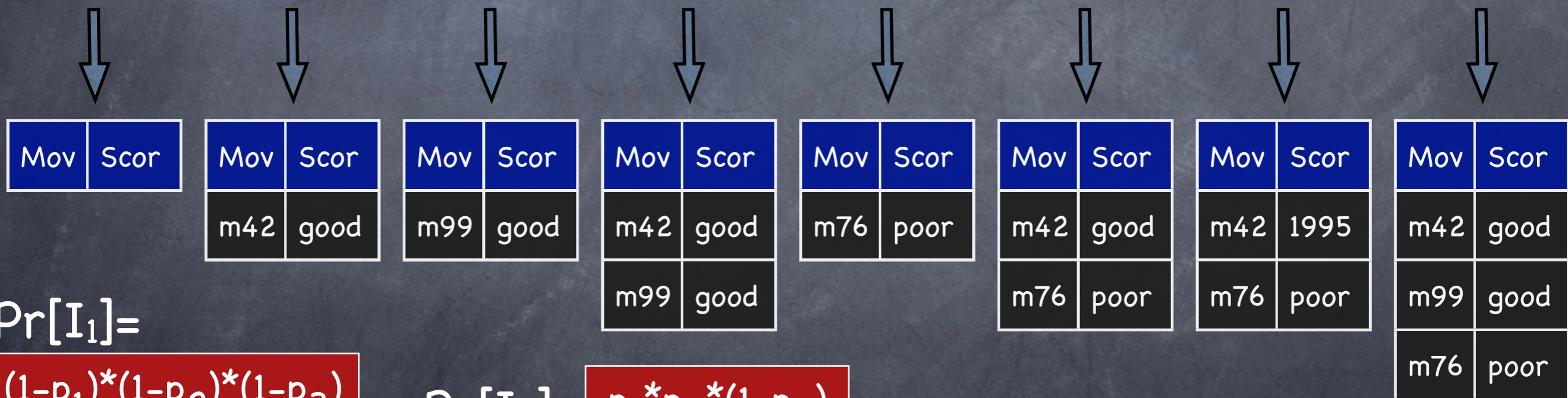
I-Databases

Reviewsⁱ(M,S,p)

Movie	Score	p
m42	good	p ₁
m99	good	p ₂
m76	poor	p ₃



Reviews(M,S)



Pr[I₁]=

$$(1-p_1)*(1-p_2)*(1-p_3)$$

$$\text{Pr}[I_4]= p_1 * p_2 * (1-p_3)$$

$$\text{Pr}[I_8]= p_1 * p_2 * p_3$$

$$\text{Pr}[I_1] + \text{Pr}[I_2] + \dots + \text{Pr}[I_8] = 1$$

Possible worlds semantics

Disjoint Events

Needed in

- Many-to-1 matchings
- Possible values for attributes [Barbara'92]

Name	Age
John	34 (0.3)
	43 (0.7)
Mary	25



Name	Age	p
John	34	0.3
John	43	0.7
Mary	25	1.0

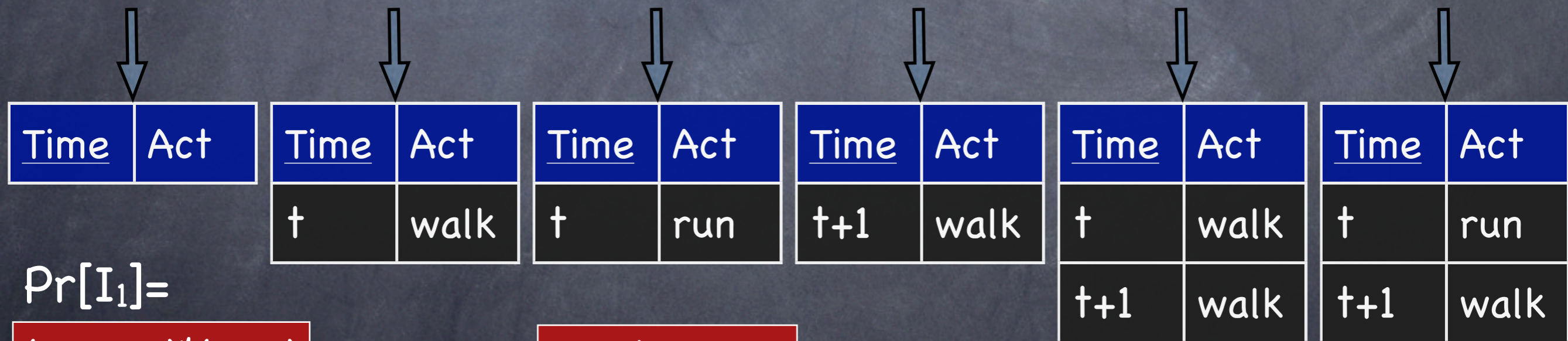
ID-Databases

Activities^{id}

Time ^d	Activity	p
t	walk	p ₁
t	run	p ₂
t+1	walk	p ₃



Activities



Pr[I₁]=

$$(1-p_1-p_2)*(1-p_3)$$

$$\text{Pr}[I_3]= p_2*(1-p_3)$$

$$\text{Pr}[I_5]= p_1*p_3$$

$$\text{Pr}[I_1] + \text{Pr}[I_2] + \dots + \text{Pr}[I_6] = 1$$

ID subsumes I

Reviews^{id}

Movie ^d	Score ^d	p
m42	good	p ₁
m99	good	p ₂
m76	poor	p ₃

=

Reviewsⁱ

Movie	Score	p
m42	good	p ₁
m99	good	p ₂
m76	poor	p ₃

Note:

Reviews^{id}

Movie	Score	p
m42	good	p ₁
m99	good	p ₂
m76	poor	p ₃

means all
tuples are
disjoint

Queries

Syntax: conjunctive queries over schema S

$Q(y) :- \text{Movie}(x,y), \text{Review}(x,z), z \geq 3$

Movieⁱ

id	year	p
m42	1995	0.95
m99	2002	0.65
m76	2002	0.1
m05	2005	0.7

Reviewⁱ

mid	rating	p
m42	4	0.7
m42	5	0.45
m99	5	0.82
m99	4	0.68
m05	5	0.79

Two Query Semantics

Possible answer sets

- Given set A:

$$\Pr[\{t \mid I \models Q(t)\} = A]$$

- Used for views

Possible tuples

- Given tuple t:

$$\Pr[I \models Q(t)]$$

- Used for query evaluation and top-k



This
talk

Query Semantics

id	year	p_1
m42	2004	
m99	1901	
m76	1902	

id	year	p_2
m99	1935	
m05	1903	

id	year	p_3
m76	1995	
m99	1935	
m05	2004	

id	year	p_4
m87	1934	
m44	1904	

$Q(y) :- \text{Movie}(x,y), \text{Review}(x,z)$

Tuple probabilities

year	p
1935	$p_2 + p_3 = 0.6$
2004	$p_1 + p_3 = 0.5$
1995	$p_3 = 0.2$
...	...

↓ top k

Summary on Data Model

- Data Model:

 - Semantics = possible worlds

 - Syntax = I-databases or ID-databases

- Queries:

 - Syntax = unchanged (conjunctive queries)

 - Semantics = tuple probabilities

Problem Definition

Fix schema S , query Q , answer tuple t

Problem: given I/ID-database DB , compute $\Pr[I \models Q(t)]$

notation: $\Pr[Q(t)]$

Conventions:

For upper bounds (P or $\#P$): probabilities are rationals

For lower bounds ($\#P$): probabilities are $1/2$

Query Evaluation on I-Databases

Outline

- Intuition
- Extensional plans: PTIME case
- Hard queries: #P-complete case
- Dichotomy Theorem

$Q(y) :- \text{Movie}(x,y), \text{Review}(x,z)$

Intuition

Movieⁱ

id	year	p
m42	1995	p ₁
m99	2002	p ₂
m76	2002	p ₃
m05	2005	p ₄

Reviewⁱ

mid	rate	p
m42	4	q ₁
m42	2	q ₂
m42	3	q ₃
m99	1	q ₄
m99	3	q ₅
m76	5	q ₆

Answer

Year	p
1995	$p_1 \times (1 - (1 - q_1) \times (1 - q_2) \times (1 - q_3))$
2002	$1 - (1 - p_2 \times (1 - (1 - q_4) \times (1 - q_5))) \times (1 - p_3 \times q_6)$

I-Extensional Plans

[Barbara92,Lakshmanan97]

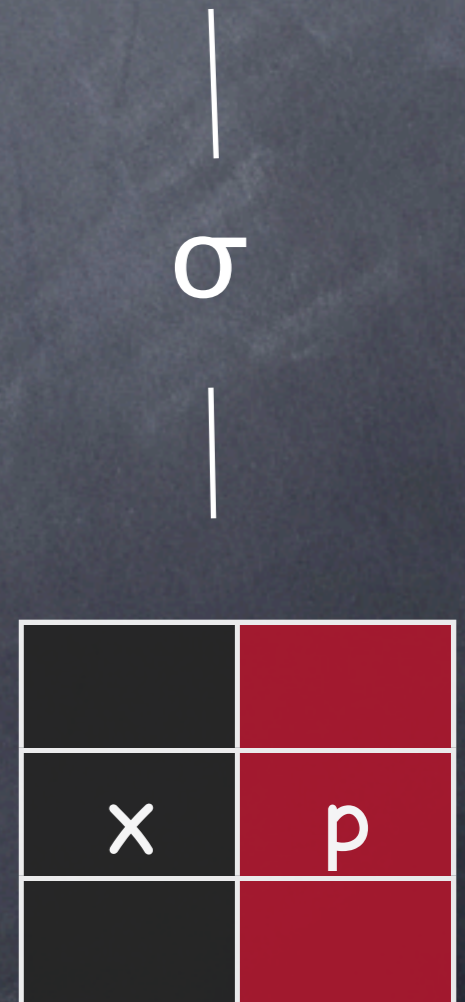
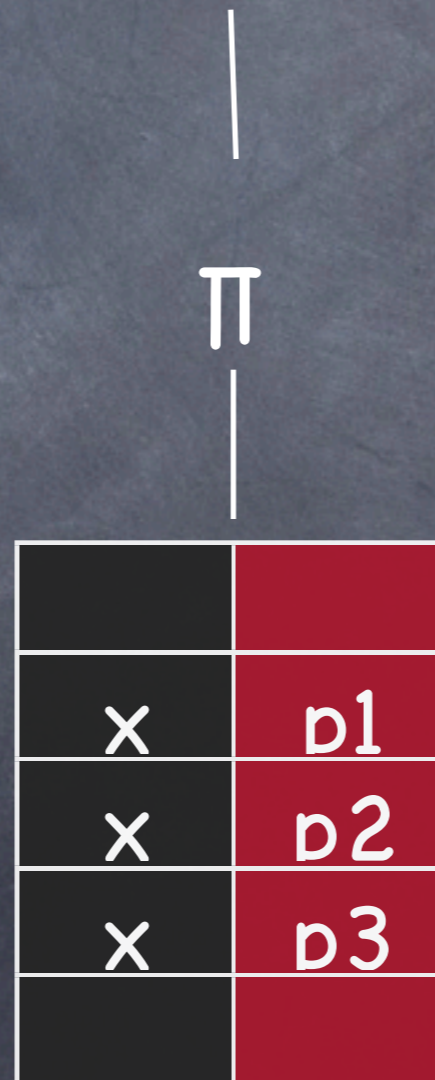
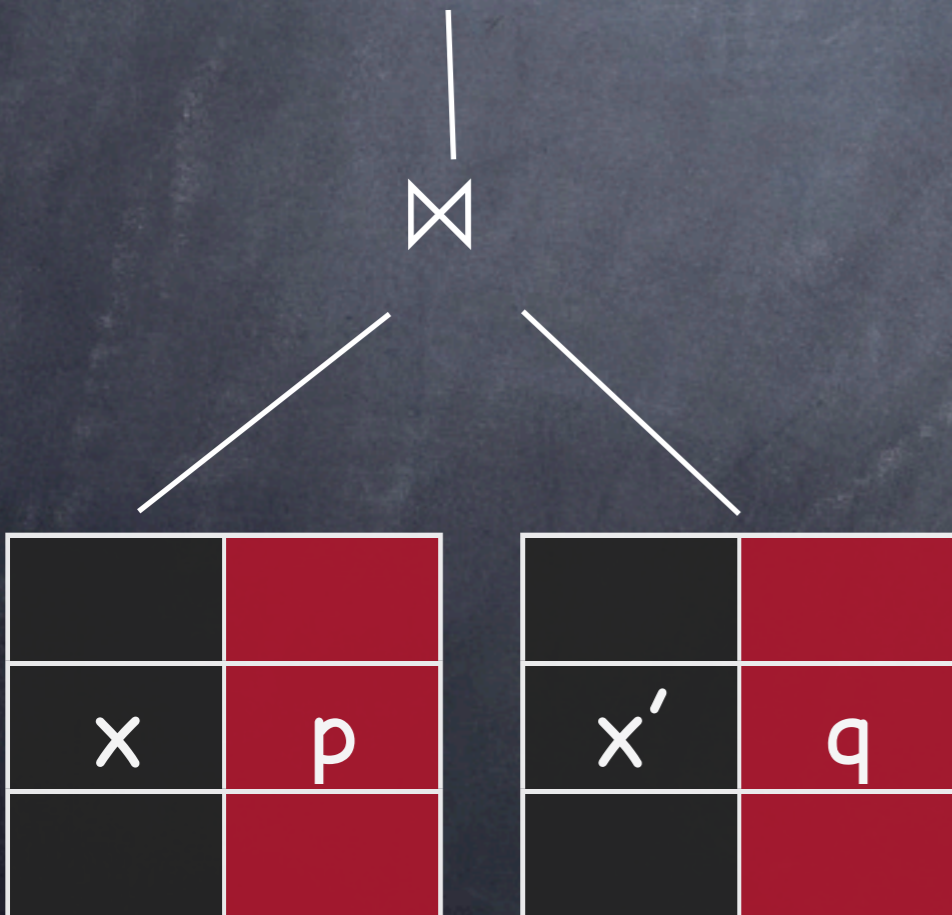
- Add p
- Join \bowtie $p = p_1 * p_2$
- Projection Π $p = 1 - (1 - p_1)(1 - p_2) \dots (1 - p_n)$
- Selection σ $p = p$
- Note: data complexity is PTIME

Extensional Query Plans

x	x'	pq

x	$1-(1-p1)(1-p2)(1-p3)$

x	p

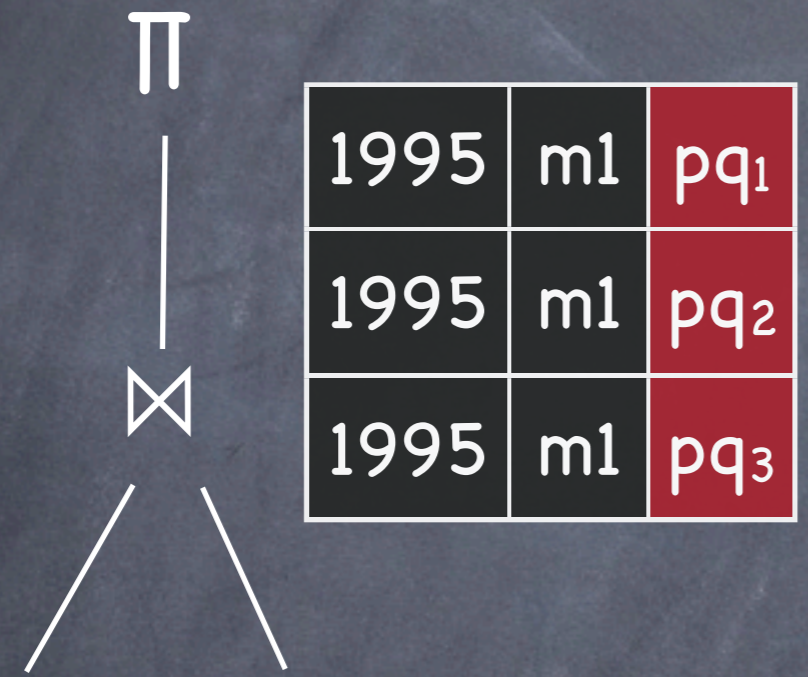


Extensional Query Plans

- Each tuple t has a probability $t.P$
- Algebra operators compute $t.P$
- Data complexity: PTIME

$Q(y) :- \text{Movie}(x,y), \text{Review}(x,z)$

1995 $1-(1-pq_1)(1-pq_2)(1-pq_3)$



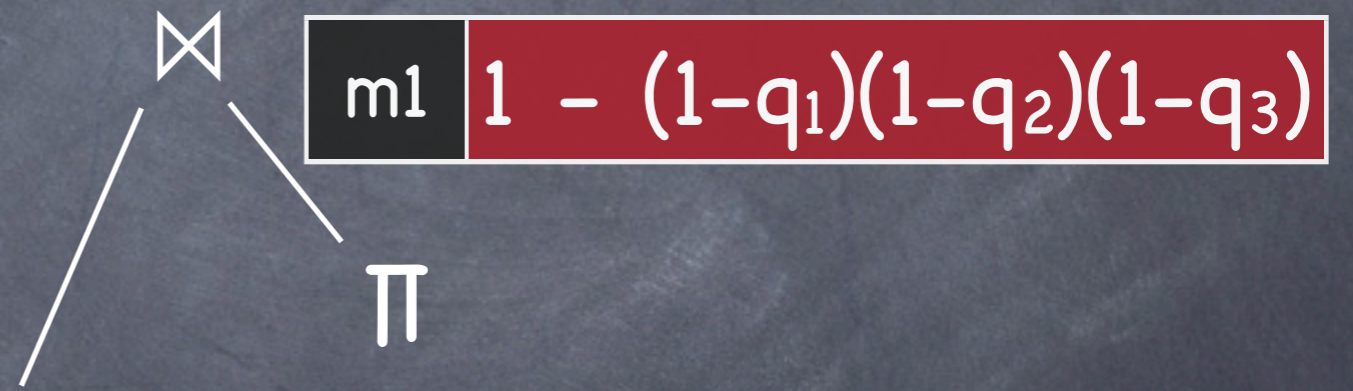
Movie **Review**

1995	p	m1	q1
		m1	q2
		m1	q3

INCORRECT!

Π

1995 m1 $p(1-(1-q_1)(1-q_2)(1-q_3))$



Movie **Review**

1995	p	m1	q1
		m1	q2
		m1	q3



CORRECT

#P-Complete Queries

R^i

A	p
	p_1
	p_2
	p_3
	p_4

S

A	B

T^i

B	p
	q_1
	q_2
	q_3
	q_4

$Q_{\text{bad}} :- R^i(x), S(x,y), T^i(y)$

Theorem: Data complexity is #P-complete

Proof:

Theorem [Provan&Ball83] Counting the number of satisfying assignments for bipartite DNF is #P-complete

Reduction:

$$x_2 y_3 \vee x_1 y_2 \vee x_4 y_3 \vee x_3 y_1$$

R^i

A	p
x_1	1/2
x_2	1/2
x_3	1/2
x_4	1/2

S

A	B
x_2	y_3
x_1	y_2
x_4	y_3
x_3	y_1

T^i

B	p
y_1	1/2
y_2	1/2
y_3	1/2

$$Q_{\text{bad}} := R^i(x), S(x,y), T^i(y)$$

I-Dichotomy

Q = boolean conjunctive query

Definition 1. For each variable x :

$\text{goals}(x)$ = set of goals that contain x

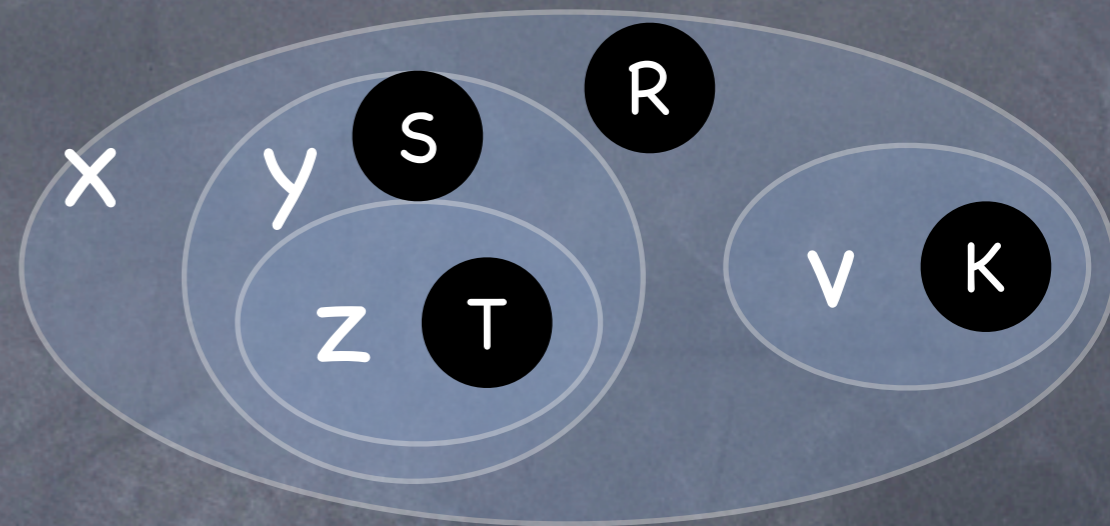
Definition 2. Q is hierarchical if for all x, y :

(a) $\text{goals}(x) \cap \text{goals}(y) = \emptyset$, or

(b) $\text{goals}(x) \subseteq \text{goals}(y)$, or

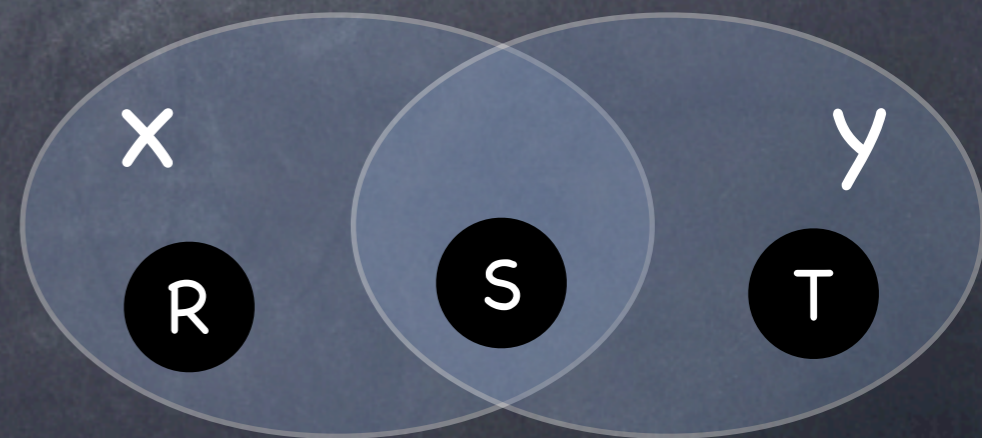
(c) $\text{goals}(y) \subseteq \text{goals}(x)$

$Q :- R(x), S(x,y), T(x,y,z), K(x,v)$



“hierarchical”

$Q :- R(x), S(x,y), T(y)$



“non-hierarchical”

[Dalvi&S.'04]

I-Dichotomy

Schema $S^i = \{R_1^i, R_2^i, \dots, R_m^i\}$

Theorem Let Q = conjunctive query w/o self-joins.
Then one of the following holds:

Q is in PTIME

Q has a correct extensional plan

Q is hierarchical

or:

Q is #P-complete

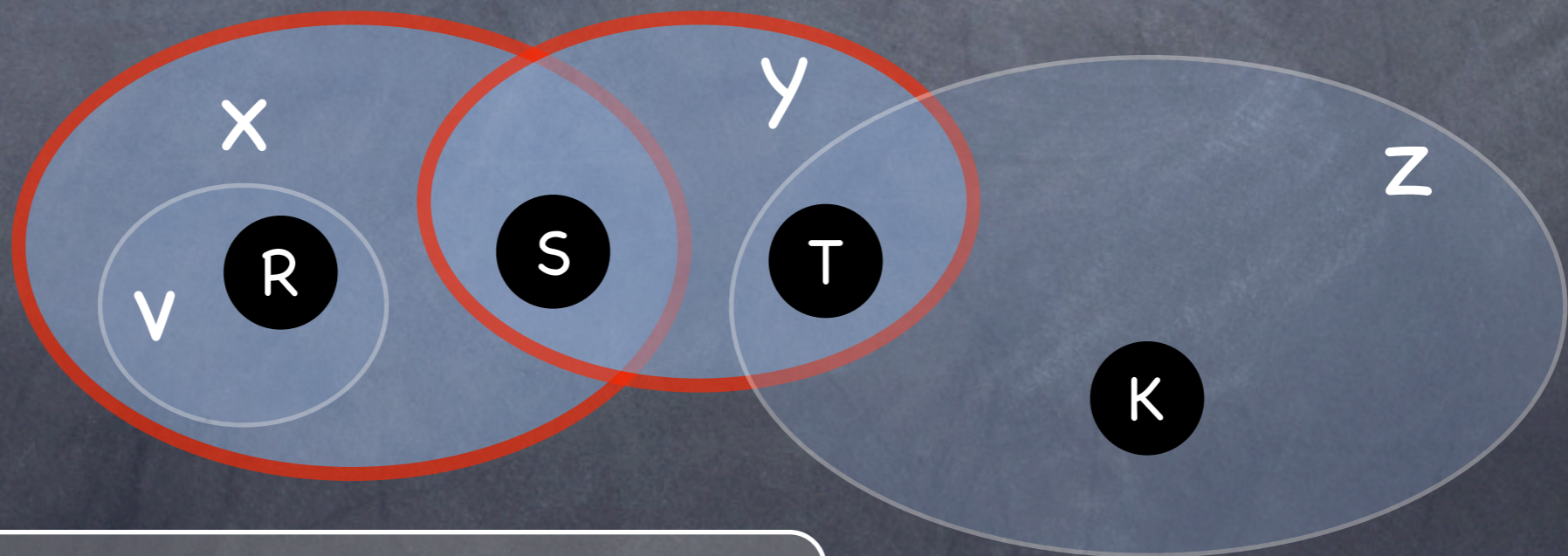
Q has subgoals $R(x, \dots), S(x, y, \dots), T(y, \dots)$

Proof

Lemma 1.

If Q is non-hierarchical, then $\#P$ -complete

Proof:



$Q := R^i(v, \underline{x}), S^i(\underline{x}, y), T^i(y, z), K^i(z)$

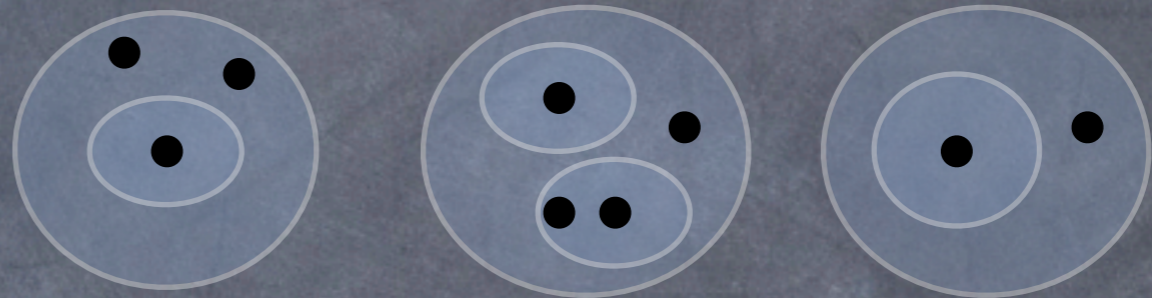
rest is like for Q_{bad}

Proof

Lemma 2. If Q is hierarchical, then P_{TIME}

Proof:

Case 1: has no root



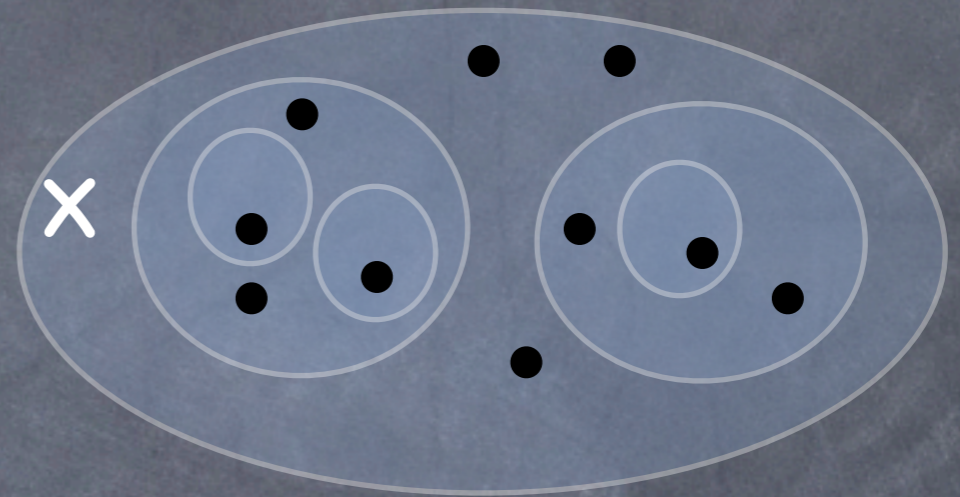
$$\Pr(Q) = \Pr(Q_1) \Pr(Q_2) \Pr(Q_3)$$

This is extensional join \bowtie

Proof

Case 2: has root x

$\text{Dom} = \{a_1, a_2, \dots, a_n\}$



$\text{Pr}(Q) =$

$$1 - (1 - \text{Pr}(Q(a_1/x))(1 - \text{Pr}(Q(a_2/x)) \dots (1 - \text{Pr}(Q(a_n/x))))$$

This is an extensional projection: Π

QED

Query Evaluation on ID-Databases

- ID-extensional plans
- #P-complete queries
- Dichotomy Theorem

Extensional Plans for ID-DBs

- Only difference: two kinds of projections:
independent $1 - (1 - p_1) \dots (1 - p_n)$
disjoint $p_1 + \dots + p_n$

#P-Complete Queries

$$Q_1 :- R^i(x), S^i(x,y), T^i(y)$$

$$Q_2 :- R^d(x^d,y), S^d(y^d)$$

$$Q_3 :- R^d(x^d,y), S^d(z^d,y)$$

[Dalvi&S.'04]

I-DB Dichotomy

Schema S^{id} s.t. each table is either R^i or R^{id}

Theorem Let Q = conjunctive query w/o self-joins.
Then one of the following holds:

Q is in PTIME

Q has a correct extensional plan

or:

Q is #P-complete

Q has one of Q_1, Q_2, Q_3 as subqueries

Extensions

Extensions of the dichotomy theorem exists for:

- Mixed schemas (some relations are deterministic)
- Functional dependencies

Summary on Query Evaluation

Extensional plans: popular, efficient, BUT

- “Equivalent” plans lead to different results
- Some queries admit “correct” plans

Some simple queries: #P-complete complexity

Dichotomy theorem

Future work: remove ‘no-self-join’ restriction

Conclusions

- Strong motivation from practical applications
Merge query and search technologies
- Probabilistic DB's are hard !
Hacks don't work (yet). Need principled approach.

Thank you !

Questions ?