

Topics in Probabilistic and Statistical Databases

Lecture 3: Representation and Query Evaluation

Dan Suciu
University of Washington

Review of Representation

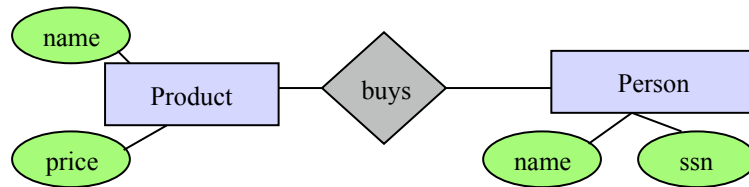
- Represent the possible worlds
- Represent the probability distribution

Review: Rule of Thumb #1

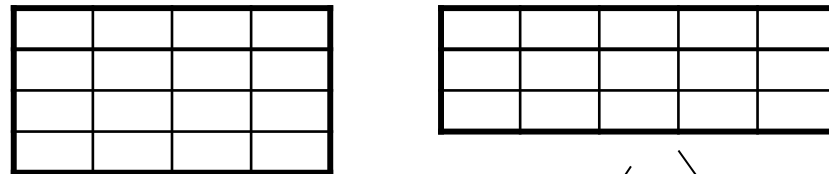
A ProbDB = an Incomplete DB + probabilities

DB 101: Database Design

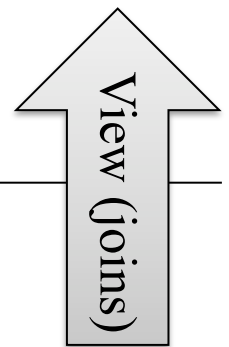
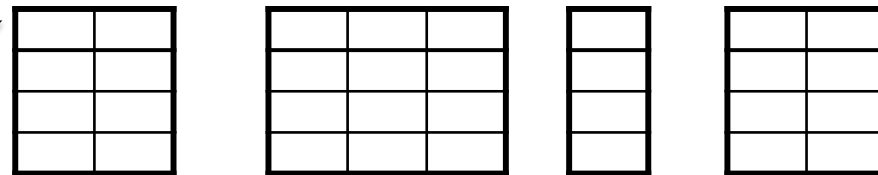
Conceptual Model:



Relational Model:
plus FD's



Normalization:
Eliminates *anomalies*



Review: Rule of Thumb #2

Probabilistic table =
Disjoint/independent tables + view (joins)

Discussion

Design

- How do we find the right decomposition ?
 - Normalization theory (n'th normal form...)
 - Probabilistic networks (same thing...)

Representation

- How to we recover the original table from the decomposed table(s) ?
 - At UW: we just use a SQL view
 - Elsewhere: equivalent formalism (e.g. U-relations)

[Antova'2008]

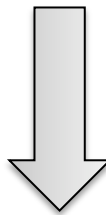
Representation: U-Relations

Two idea:

- Vertical decomposition
- Uniform lineage

Review: Vertical Decomposition

<u>TS</u>	Id	Type	Faction
a	1	Tank	Friend
b	3	Transport	Friend
c	3	Transport	Enemy
d	4	Transport	Enemy



<u>TS</u>	Id
a	1
b	3
c	3
d	4

<u>TS</u>	Type
a	Tank
b	Transport
c	Transport
d	Transport

<u>TS</u>	Faction
a	Friend
b	Friend
c	Enemy
d	Enemy

Review: Uniform Lineage

Saw

Witness	Car	X	V
Amy	Mazda	X1	1
Amy	Toyota	X1	2
Betty	Honda	X2	1

Drives

Person	Car	Y	W
Jimmy	Mazda	Y1	1
Jimmy	Toyota	Y1	2
Billy	Mazda	Y2	1
Billy	Honda	Y2	2

What is the lineage here ?

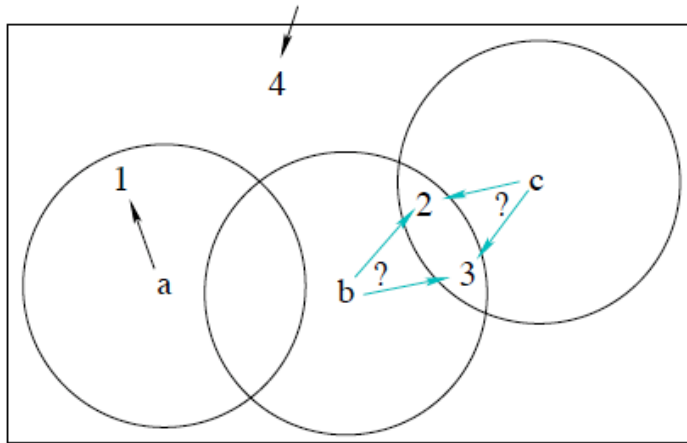
$q(w,p) :- \text{Saw}(w,c), \text{Drives}(c,p)$

Witness	Person	X	V	Y	W
Amy	Jimmy	X1	1	Y1	1
Amy	Jimmy	X1	2	Y1	2
Betty	Billy	X2	1	Y2	2

$X1=1 \wedge Y1=1 \vee$
 $X1=2 \wedge Y1=2$

[Antova'2008]

U-Relations



W	Var	Rng
	x	1
	x	2
	y	1
	y	2
	z	1
	z	2

U_1	D	T_R	Id
		a	1
	$x \mapsto 1$	b	2
	$x \mapsto 2$	b	3
	$x \mapsto 1$	c	3
	$x \mapsto 2$	c	2
		d	4

U_2	D	T_R	Type
		a	Tank
		b	Transport
		c	Tank
	$y \mapsto 1$	d	Tank
	$y \mapsto 2$	d	Transport

U_3	D	T_R	Faction
		a	Friend
		b	Friend
		c	Enemy
	$z \mapsto 1$	d	Friend
	$z \mapsto 2$	d	Enemy

Discussions

- Show the c-table $R(\text{TS}, \text{Id}, \text{Type Faction})$
- Give the view expression that recovers R
- Is this a uniform lineage ?
- Do we need W ?

[Antova'2008]

Show the c-Table

D1	D2	D3	<u>TS</u>	Id	Type	Faction
			a	1	Tank	Friend
x → 1			b	2	Transport	Friend
x → 2			b	3	Transport	Friend
x → 1			c	3	Tank	Enemy
x → 2			c	2	Tank	Enemy
	y → 1	z → 1	d	4	Tank	Friend
	y → 1	z → 2	d	4	Tank	Enemy
	y → 2	z → 1	d	4	Transport	Friend
	y → 2	z → 2	d	4	Transport	Enemy

[Antova'2008]

Why We Need W

Need to store the probability distribution of the the variables separately

<i>W</i>	Var	Rng	Pr
	x	1	0.1
	x	2	0.9
	y	1	0.3
	y	2	0.7
	z	1	0.6
	z	2	0.4

“Reduced” U-database

Definition A U-database is reduced if every tuple in every U-relation can appear in a possible world

U_1	D	T	A
	$c_1 \mapsto 1$	t_1	a_1
	$c_2 \mapsto 1$	t_2	a_2

U_2	D	T	B
	$c_1 \mapsto 1$	t_1	b_1
	$c_1 \mapsto 2$	t_1	b_2

This is not reduced
WHY ?

[Antova'2008]

More Examples

$$S = \pi_{\text{Id}}(\sigma_{\text{Type}='Tank' \wedge \text{Faction}='Enemy'}(R))$$

U_4	D_1	D_2	T_S	Id
	$x \mapsto 1$		c	3
	$x \mapsto 2$		c	2
	$y \mapsto 1$	$z \mapsto 2$	d	4

$$(S \ s_1) \bowtie_{s_1.\text{Id} \neq s_2.\text{Id}} (S \ s_2)$$

U_5	D_1	D_2	D_3	T_{s_1}	T_{s_2}	Id ₁	Id ₂
	$x \mapsto 1$	$y \mapsto 1$	$z \mapsto 2$	c	d	3	4
	$x \mapsto 2$	$y \mapsto 1$	$z \mapsto 2$	c	d	2	4
	$y \mapsto 1$	$z \mapsto 2$	$x \mapsto 1$	d	c	4	3
	$y \mapsto 1$	$z \mapsto 2$	$x \mapsto 2$	d	c	4	2

Summary of U-Relations

U-relations = Vertical partition + uniform lineage

- Vertical partition: not necessarily disjoint attrs
 - Consistency check needed if multiple occurrences of A
- Reduced U-relations = semijoin reduction
- “Parsimonious query translation”:
 - Select-project-join queries over possible worlds →
select-project-join queries over U-relations

The Design Problem

- How do we “normalize” a probabilistic database ?
- Combine:
 - MVD = multi-valued dependencies (read book..)
 - Probabilistic networks (next)
- NOTE: This is largely unexplored area

Conditional Independence

- Recall: $X \perp\!\!\!\perp Y \mid Z$ means:

$$P(XY \mid Z) = P(X \mid Z) * P(Y \mid Z)$$

- Equivalently:

$$P(XYZ) = P(XZ) * P(YZ) / P(Z)$$

Review: Rule of Thumb #3

Conditional independence =
MVD + some probability identities

Conditional Independence=MVD

- If $X \perp\!\!\!\perp Y \mid Z$ then $X \twoheadrightarrow Y$
- Stronger: $P(X, Y, Z) = P_1(X, Z) \bowtie P_2(Y, Z)$
 - This identity holds between *probabilistic tables*

Probabilistic Networks

- Many variables V_1, V_2, \dots, V_n
- Given joint distribution on their values
- Goal of a PN: to capture all conditional independences
- An edge (V_i, V_j) means, intuitively, that V_i, V_j are dependent
- Lack of an edge means they are independent

Quiz Time

- Three random variables V_1, V_2, V_3
- Suppose their joint distribution is given by a function $p(v_1, v_2, v_3) = f(v_1, v_2) * g(v_2, v_3)$, where f, g are arbitrary positive functions

QUESTION What independence relation holds between V_1, V_2, V_3 ?

Quiz Time

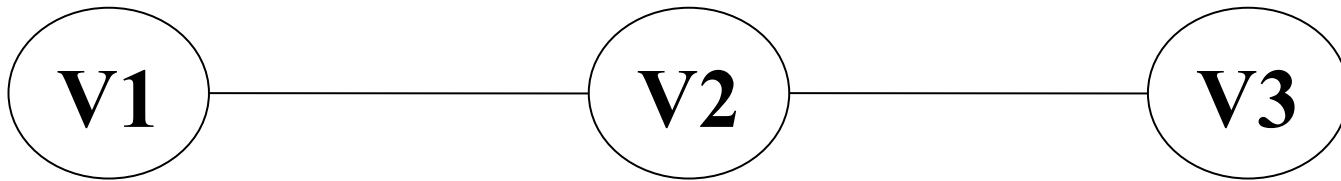
- Three random variables V_1, V_2, V_3
- Suppose their joint distribution is given by a function $p(v_1, v_2, v_3) = f(v_1, v_2) * g(v_2, v_3)$, where f, g are arbitrary positive functions

QUESTION What independence relation holds between V_1, V_2, V_3 ?

ANSWER $V_1 \perp\!\!\!\perp V_3 \mid V_2$

(proof in class)

A Probabilistic Network



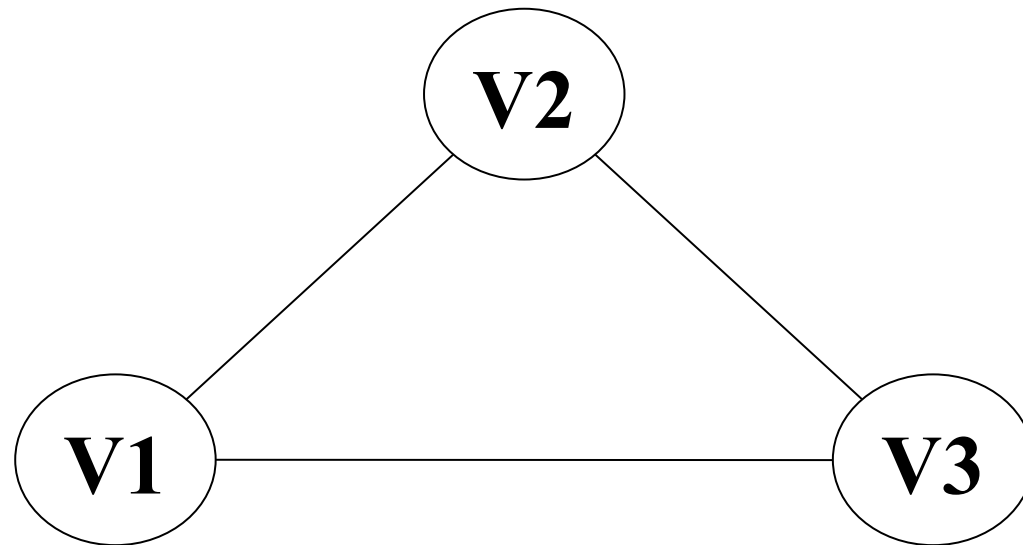
Quiz Time 2

- Now suppose

$$p(v_1, v_2, v_3) = f(v_1, v_2) * g(v_2, v_3) * h(v_1, v_3)$$

QUESTION What independence relation holds between V_1 , V_2 , V_3 ?

Another Probabilistic Network



Answer: NONE !

Hierarchical Decomposition

- $P(V_1, \dots, V_n)$ is a *Hierarchical Distribution* if there exists functions f_1, \dots, f_k , s.t.:
 - $P = f_1 * f_2 * \dots * f_k$
 - $f_i(V_{ci})$ depends on a subset $V_{ci} \subseteq \{V_1, \dots, V_n\}$
- Also called: *Probabilistic Graphical Model*
- f_1, \dots, f_k (or the set V_{ci}) are called *factors*

Probabilistic Network

- Define a graph $G = (V, E)$ s.t.
 - $V = \{V_1, \dots, V_n\}$
 - $E = \{(V_i, V_j) \mid \text{exists } \mathcal{V}_c \text{ s.t. } V_i \in \mathcal{V}_c, V_j \in \mathcal{V}_c\}$
- A *clique* is $C \subseteq V$ s.t. any two nodes in C are connected
- A *separator* is $S = C \cap C'$, where C, C' =cliques

[Cowell'99]

Probabilistic Network

- Let C_1, C_2, \dots be all the cliques, and S_1, S_2, \dots all separators
- Then:

$$P(V) = P(\mathcal{V}_{C_1}) * P(\mathcal{V}_{C_2}) * \dots / (P(\mathcal{V}_{S_1}) * P(\mathcal{V}_{S_2}) * \dots)$$

Research question: how does decompose P into D/I tables + views ?

Query Evaluation on D/I Databases

Problem Statement

- Given:
 - A disjoint/independent probdb PDB
 - A Boolean conjunctive query Q
- Compute the probability $Q(\text{PDB})$

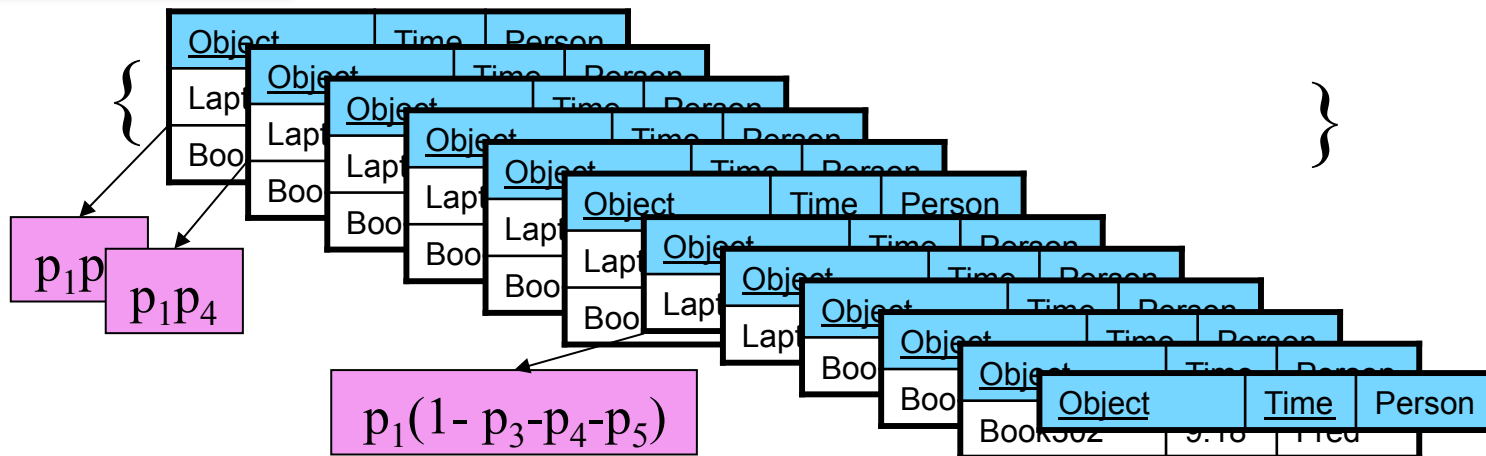
Why does it suffice to restrict to D/I databases ?

Review: D/I Databases

S =

<u>Object</u>	<u>Time</u>	Person	P
LaptopX77	9:07	John	p ₁
LaptopX77	9:07	Jim	p ₂
Book302	9:18	Mary	p ₃
Book302	9:18	John	p ₄
Book302	9:18	Fred	p ₅

Rep(S) =



Review Query Semantics

Semantics 1: Possible Sets of Answers

A probability distributions on sets of tuples

$$\Pr(Q = A) = \sum_{I \in \text{Inst. } Q(I) = A} \Pr(I)$$

Semantics 2: Possible Tuples

A probability function on tuples

$$\Pr(t \in Q) = \sum_{I \in \text{Inst. } t \in Q(I)} \Pr(I)$$

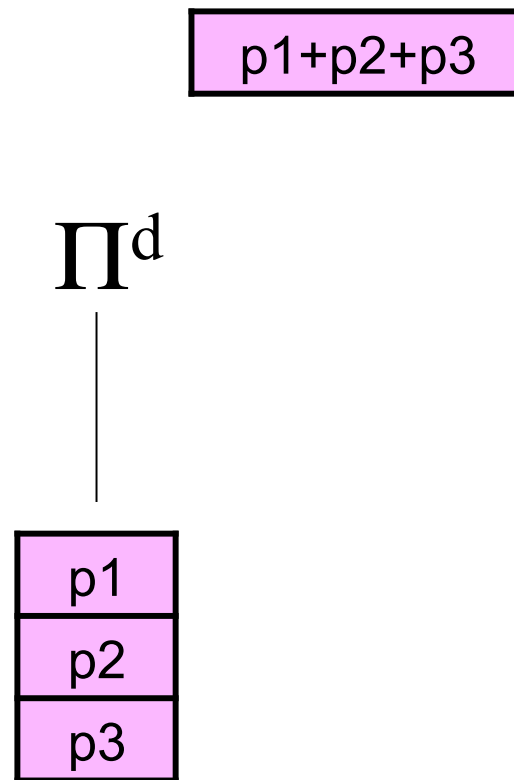
Extensional Operators

Object	Person	Location	P
Laptop77	John	L45	p1
	Jim	L45	p2
	Jim	L66	p3
Book302	Mary	L66	p4
	Mary	L45	p5
	Jim	L66	p6
	John	L45	p7
	Fred	L45	p8

$q(z) :- \text{HasObject}^P(\mathbf{Book302}, y, z)$

Location	P
L66	p4+p6
L45	p5+p7+p8

Disjoint Project



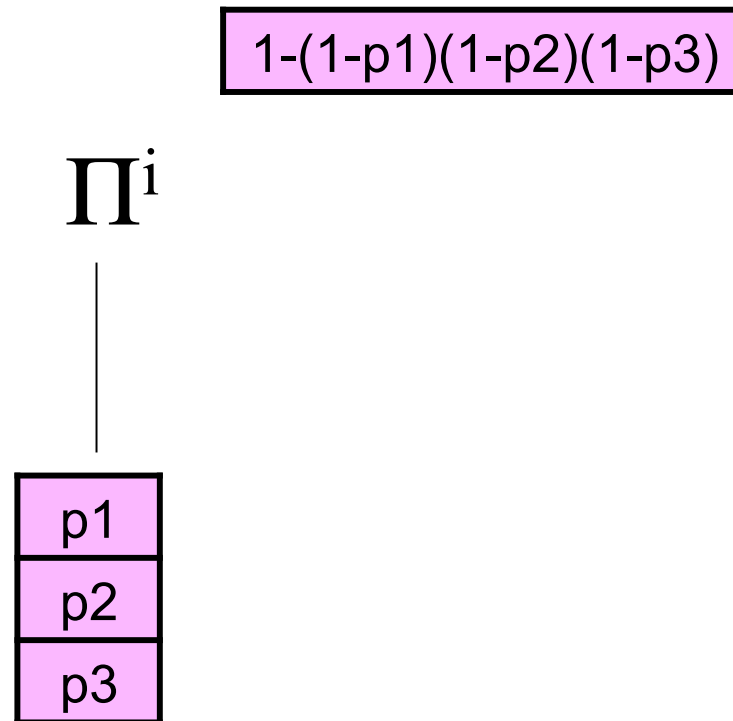
Extensional Operators

Object	Person	Location	P
Laptop77	John	L45	p1
	Jim	L45	p2
	Jim	L66	p3
Book302	Mary	L66	p4
	Mary	L45	p5
	Jim	L66	p6
	John	L45	p7
	Fred	L45	p8

Person	Location	P
Jim	L66	$1-(1-p3)(1-p6)$
John	L45	$1-(1-p1)(1-p7)$
...		

$q(y,z) :- \text{HasObject}^p(\underline{x},y,z)$

Independent Project



$$q(y) :- \text{Movie}^p(\underline{x}, y), \text{Review}^p(\underline{x}, z), z > 3$$

Example

Movie

id	year	P
m42	1995	p1
m99	2002	p2
m76	2002	p3

Review

mid	rating	P
m42	7	q1
m42	4	q2
m42	9	q3
m99	7	q4
m99	5	q5
m76	6	q6

Answer

year	P
1995	$p1 \times (1 - (1 - q1) \times (1 - q2) \times (1 - q3))$
2002	$1 - (1 - p2 \times (1 - (1 - q4) \times (1 - q5))) \times (1 - p3 \times q6)$

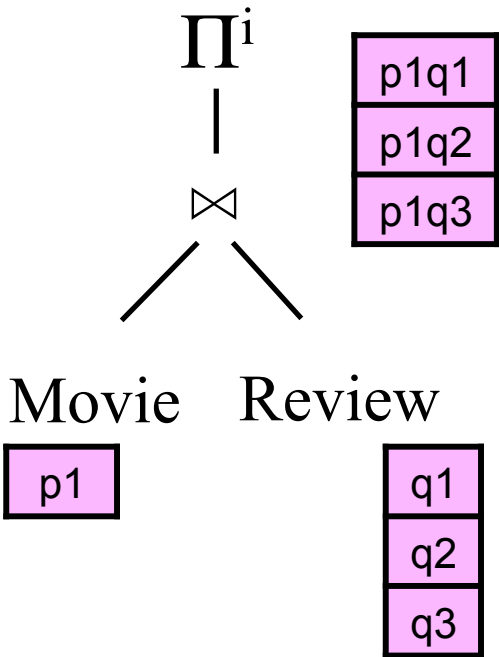
$q(y) :- \text{Movie}^p(\underline{x}, y), \text{Review}^p(\underline{x}, \underline{z})$

$q(1995)$

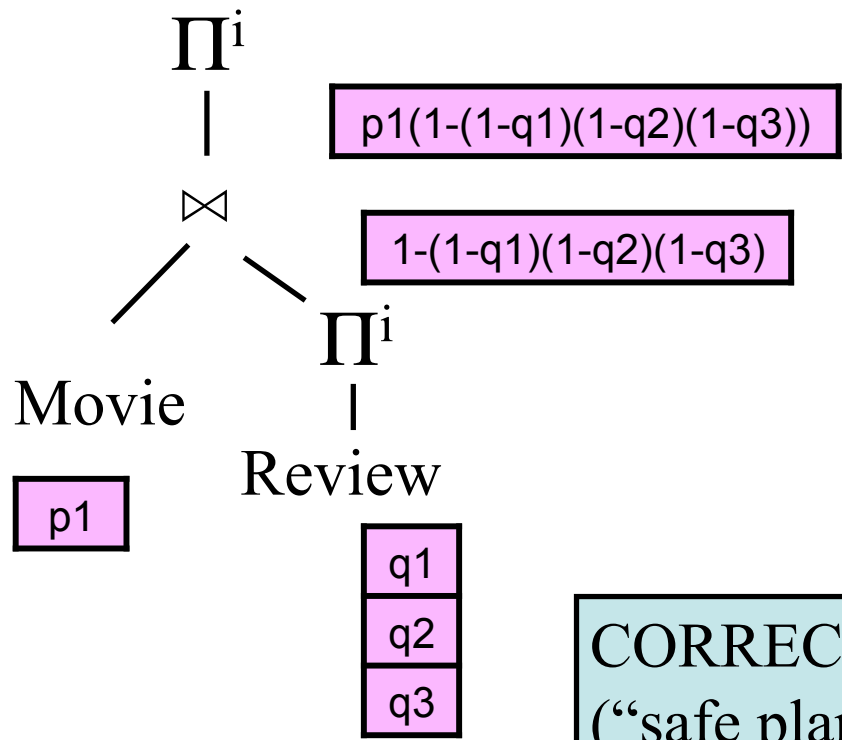
Answer depends on query plan !

$1 - (1 - p_1 q_1)(1 - p_1 q_2)(1 - p_1 q_3)$

$1 - (1 - p_1(1 - (1 - q_1)(1 - q_2)(1 - q_3)))(1 - \dots) \dots$



INCORRECT



CORRECT
("safe plan")

Safe Plans are Efficient

- Very efficient: run almost as fast as regular queries
- Require only simple modifications of the relational operators
- Or can be translated back into SQL and sent to any RDBMS

Can we always generate a safe plan ?

A Hard Query

R^p

<u>A</u>	<u>B</u>	P
a	x1	p1
a	x2	p2

S

<u>B</u>	<u>C</u>
x1	y1
x1	y2
x2	y1

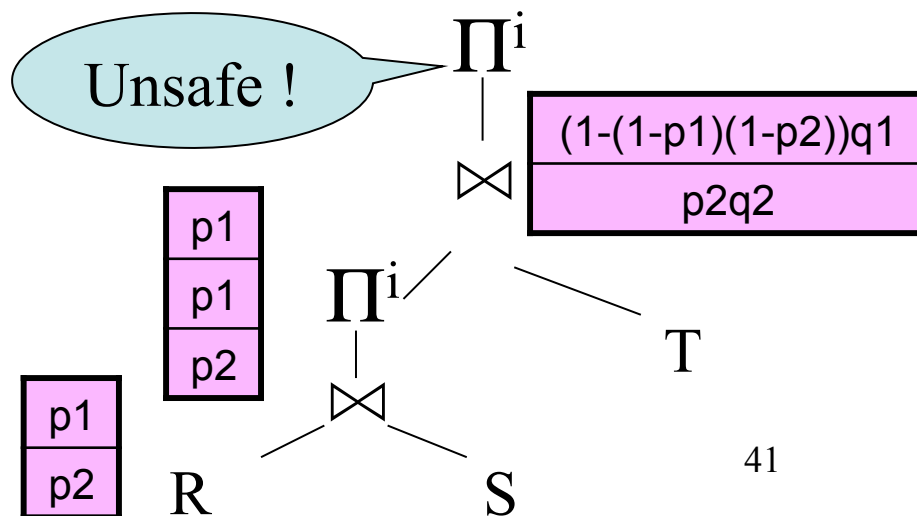
T^p

<u>C</u>	<u>D</u>	P
y1	c	q1
y2	c	q2

$h(u,v) :- R^p(\underline{u},\underline{x}), S(\underline{x},\underline{y}), T^p(\underline{y},\underline{v})$

$h(a,c)$

There is no safe plan !



Independent Queries

Let q_1, q_2 be two boolean queries

Definition q_1, q_2 are “independent” if $P(q_1, q_2) = P(q_1) P(q_2)$

Also: $P(q_1 \vee q_2) = 1 - (1 - P(q_1))(1 - P(q_2))$

Quiz: which are independent ?

q1	q2	Indep.?
Movie ^p (<u>m41</u> , <u>y</u>)	Review ^p (<u>m41</u> , <u>z</u>)	
Movie ^p (<u>m42</u> , <u>y</u>), Review ^p (<u>m42</u> , <u>z</u>)	Movie ^p (<u>m77</u> , <u>y</u>), Review ^p (<u>m77</u> , <u>z</u>)	
Movie ^p (<u>m42</u> , <u>y</u>), Review ^p (<u>m42</u> , <u>z</u>)	Movie ^p (<u>m42</u> , <u>1995</u>)	
Movie ^p (<u>m42</u> , <u>y</u>), Review ^p (<u>m42</u> , <u>7</u>)	Movie ^p (<u>m42</u> , <u>y</u>), Review ^p (<u>m42</u> , <u>4</u>)	
R ^p (<u>x</u> , <u>y</u> , <u>z</u> , <u>z</u> , <u>u</u>), R ^p (<u>x</u> , <u>x</u> , <u>x</u> , <u>y</u> , <u>y</u>)	R ^p (<u>a</u> , <u>a</u> , <u>b</u> , <u>b</u> , <u>c</u>)	

Answers

q1	q2	Indep.?
Movie ^p (<u>m41</u> , <u>y</u>)	Review ^p (<u>m41</u> , <u>z</u>)	YES
Movie ^p (<u>m42</u> , <u>y</u>),Review ^p (<u>m42</u> , <u>z</u>)	Movie ^p (<u>m77</u> , <u>y</u>),Review ^p (<u>m77</u> , <u>z</u>)	YES
Movie ^p (<u>m42</u> , <u>y</u>),Review ^p (<u>m42</u> , <u>z</u>)	Movie ^p (<u>m42</u> , <u>1995</u>)	NO
Movie ^p (<u>m42</u> , <u>y</u>),Review ^p (<u>m42</u> , <u>7</u>)	Movie ^p (<u>m42</u> , <u>y</u>),Review ^p (<u>m42</u> , <u>4</u>)	NO
R ^p (<u>x</u> , <u>y</u> , <u>z</u> , <u>z</u> , <u>u</u>), R ^p (<u>x</u> , <u>x</u> , <u>x</u> , <u>y</u> , <u>y</u>)	R ^p (<u>a</u> , <u>a</u> , <u>b</u> , <u>b</u> , <u>c</u>)	YES

Prop If no two subgoals unify then q1,q2 are independent

Detour: Independent Queries

- Let q = a Boolean conjunctive query
- A critical tuple is a tuple t s.t.
 - There exists an instance I s.t.
 - $q(I)$ is true
 - $q(I - \{t\})$ is false
 - (in other words, $I \models q$ and $I - \{t\} \not\models q$)
- Denote $\text{crit}(q)$ = the set of critical tuples

Detour: Independent Queries

Note: *necessary* but not *sufficient* condition

Theorem Queries q_1, q_2 are independent iff $\text{crit}(q_1) \cap \text{crit}(q_2) = \emptyset$

q_1 :- Movie^p(m42,y),Review^p(m42,z)

Crit(q_1) = ?

q_2 :- Movie^p(m77,y),Review^p(m77,z)

Crit(q_2) = ?

Detour: Independent Queries

Note: *necessary* but not *sufficient* condition

$q1 :- R(x, y, z, z, u), R(x, x, x, y, y).$

$q2 :- R(a, a, b, b, c),$

Theorem Independence is Π^p_2 complete [Miklau&S'04]
Reducible to query containment [Machanavajjhala&Gehrke'06]

Disjoint Queries

Let q_1, q_2 be two boolean queries

Definition q_1, q_2 are “disjoint” if $P(q_1, q_2) = 0$

Iff q_1, q_2 depend on two disjoint tuples t_1, t_2

Quiz: which are disjoint ?

q1	q2	?
HasObject ^p (<u>'book'</u> , <u>'9'</u> , 'Mary', x)	HasObject ^p (<u>'book'</u> , <u>'9'</u> , 'Jim', x)	
HasObject ^p (<u>'book'</u> , <u>t</u> , 'Mary', x)	HasObject ^p (<u>'book'</u> , <u>t</u> , 'Jim', x)	
HasObject ^p (<u>'book'</u> , <u>'9'</u> , u, x)	HasObject ^p (<u>'book'</u> , <u>'9'</u> , v, x)	

Answers

q1	q2	?
HasObject ^p (<u>'book'</u> , <u>'9'</u> , 'Mary', x)	HasObject ^p (<u>'book'</u> , <u>'9'</u> , 'Jim', x)	Y
HasObject ^p (<u>'book'</u> , <u>t</u> , 'Mary', x)	HasObject ^p (<u>'book'</u> , <u>t</u> , 'Jim', x)	N
HasObject ^p (<u>'book'</u> , <u>'9'</u> , u, x)	HasObject ^p (<u>'book'</u> , <u>'9'</u> , v, x)	N

Proposition q1, q2 are “disjoint” if they contain subgoals g1, g2:

Have the same values for the key attributes

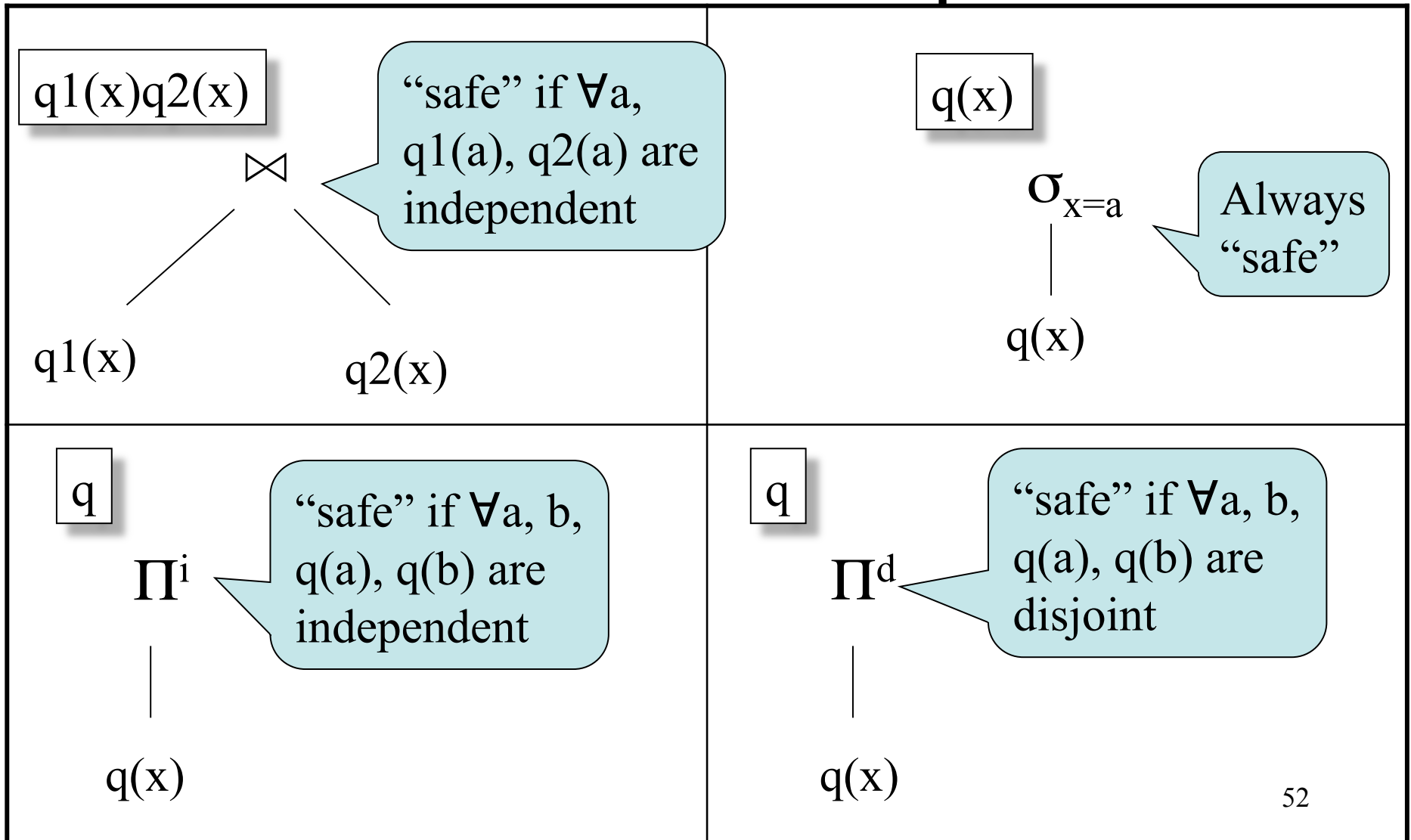
these values are constants

have at least one different constant in the non-key attributes

Definition of Safe Operators

- Safe join = left and right are indep.
- Safe independent project = duplicate tuples are independent
- Safe disjoint project = duplicate tuples are disjoint
- Safe select = any select is safe 😊

Definition of Safe Operators



$q(y^c) :- \text{Movie}^p(\underline{x}, y^c), \text{Review}^p(\underline{x}, \underline{z})$

y^c “is a constant”

Example 1

$q1 :- \text{Movie}(\underline{x}, y^c), \text{Review}(\underline{x}, \underline{z})$

Π^i

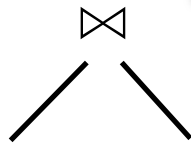
Unsafe

Because these are dependent:

$q1(m42, 7) = \text{Movie}(m42, y^c), \text{Review}(m42, 7)$

$q1(m42, 4) = \text{Movie}(m42, y^c), \text{Review}(m42, 7)$

$q1(x, z) :- \text{Movie}(\underline{x}, y^c), \text{Review}(\underline{x}, \underline{z})$



Movie Review

$q(y^c) :- \text{Movie}^p(\underline{x}, y^c), \text{Review}^p(\underline{x}, \underline{z})$

y^c “is a constant”

Example 2

$q1 :- \text{Movie}(\underline{x}, y^c), \text{Review}(\underline{x}, \underline{z})$

Π^i

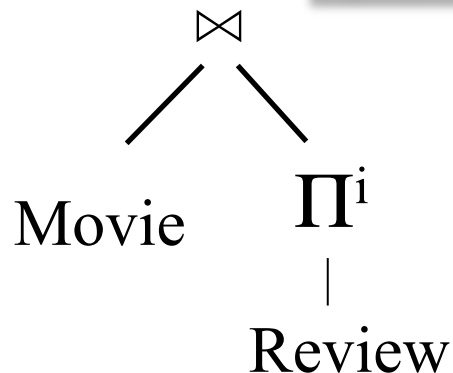
Safe !

Now these are independent !

$q1(m42) = \text{Movie}(m42, y^c), \text{Review}(m42, z)$

$q1(m77) = \text{Movie}(m77, y^c), \text{Review}(m77, z)$

$q1(x) :- \text{Movie}(\underline{x}, y^c), \text{Review}(\underline{x}, \underline{z})$



Complexity Class #P

Definition #P is the class of functions $f(x)$ for which there exists a PTIME non-deterministic Turing machine M s.t.
 $f(x)$ = number of accepting computations of M on input x

Examples:

SAT = “given formula Φ , is Φ satisfiable ?”
= NP-complete

#SAT = “given formula Φ , count # of satisfying assignments”
= #P-complete

[Valiant'79]

[Provan&Ball'83]

Examples

Class	Example	SAT	#SAT
3CNF	$(X \vee Y \vee Z) \wedge (\neg X \vee U \vee W) \dots$	NP	#P
2CNF	$(X \vee Y) \wedge (\neg X \vee U) \dots$	PTIME	#P
Positive, partitioned 2CNF	$(X_1 \vee Y_1) \wedge (X_1 \vee Y_4) \wedge$ $(X_2 \vee Y_1) \wedge (X_3 \vee Y_1) \dots$	PTIME	#P
Positive, partitioned 2DNF	$(X_1 \wedge Y_1) \vee (X_1 \wedge Y_4) \vee$ $(X_2 \wedge Y_1) \vee (X_3 \wedge Y_1) \dots$	PTIME	#P

Here NP, #P means “NP-complete, #P-complete”

See also [Graedel et al. 98]

#P-Hard Queries

hd1 :- $R^p(\underline{x}), S(\underline{x}, \underline{y}), T^p(\underline{y})$

Theorem The query hd1 is #P-hard

Proof: Reduction from partitioned, positive 2DNF

E.g. $\emptyset = x1 y1 \vee x2 y1 \vee x1 y2 \vee x3 y2$ reduces to

R^p

<u>A</u>	P
x1	0.5
x2	0.5
x3	0.5

S

<u>A</u>	<u>B</u>
x1	y1
x2	y1
x1	y2
x3	y2

T^p

<u>B</u>	P
y1	0.5
y2	0.5

$$\#\emptyset = P(\text{hd1}) * 2^n$$

Where We Are

- We have seen a query that can be evaluated with a safe plan
 - Very efficient
- We have seen a query whose data complexity is #P hard
- What is the general picture ?

PTIME Queries

$R(\underline{\mathbf{x}}, \underline{\mathbf{y}}), S(\underline{\mathbf{x}}, \underline{\mathbf{z}})$

$R(\underline{\mathbf{x}}, y), S(\underline{\mathbf{y}}), T(\underline{\mathbf{a}}, y)$

$R(\underline{\mathbf{x}}), S(\underline{\mathbf{x}}, \underline{\mathbf{y}}), T(\underline{\mathbf{y}}), U(\underline{\mathbf{u}}, y), W(\underline{\mathbf{a}}, u)$

• • •

#P-Hard Queries

hd1 = $R(\underline{\mathbf{x}}), S(\underline{\mathbf{x}}, \underline{\mathbf{y}}), T(\underline{\mathbf{y}})$

hd2 = $R(\underline{\mathbf{x}}, y), S(\underline{\mathbf{y}})$

hd3 = $R(\underline{\mathbf{x}}, y), S(x, \underline{\mathbf{y}})$

• • •

Will discuss next how to decide their complexity and how evaluate PTIME queries

Dichotomy for a Language L

- Fix a query language L
- The dichotomy property is:
 - Every query in L is either in PTIME or #P-hard
- Note that this does not follow from general principles
 - It may be false for some languages L

Dichotomy Property

LANG: CQ = conjunctive queries
CQ¹ = conjunctive queries without self-joins

Theorem The dichotomy property holds for:

CQ¹ and independent dbs.

CQ¹ and disjoint/independent dbs.

CQ and independent dbs.

We'll start these today and continue next lecture

Hierarchical Queries

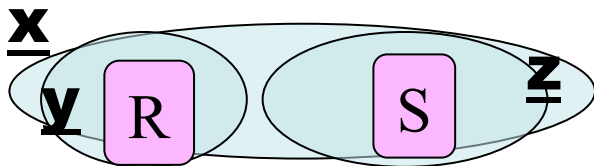
$sg(x)$ = set of subgoals containing the variable x in a key position

Definition A query q is *hierarchical* if for all x, y :

$$sg(x) \subseteq sg(y) \quad \text{or} \quad sg(y) \subseteq sg(x) \quad \text{or} \quad sg(x) \cap sg(y) = \emptyset$$

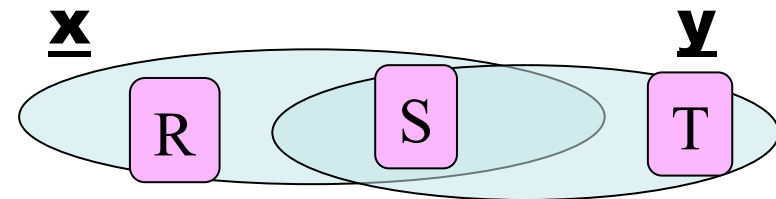
Hierarchical

$$q = R(\underline{x}, \underline{y}), S(\underline{x}, \underline{z})$$



Non-hierarchical

$$h1 = R(\underline{x}), S(\underline{x}, \underline{y}), T(\underline{y})$$



Case 1: CQ¹ + Independent

Note that in this case:

- CQ¹ (conjunctive queries, no self-joins):
 - $R(\underline{\mathbf{x}}, \underline{\mathbf{y}}), S(\underline{\mathbf{y}}, \underline{\mathbf{z}})$ OK
 - $R(\underline{\mathbf{x}}, \underline{\mathbf{y}}), R(\underline{\mathbf{y}}, \underline{\mathbf{z}})$ Not OK
- Independent tuples only:
 - $R(\underline{\mathbf{x}}, \underline{\mathbf{y}})$ OK
 - $S(\underline{\mathbf{y}}, \underline{\mathbf{z}})$ Not OK

[Dalvi&S'2004]

CQ¹ + Independent

Theorem For all $q \in \text{CQ}^1$:

q is hierarchical, has a safe plan, and is in PTIME,

OR

q is not hierarchical and is #P-hard

The PTIME Queries

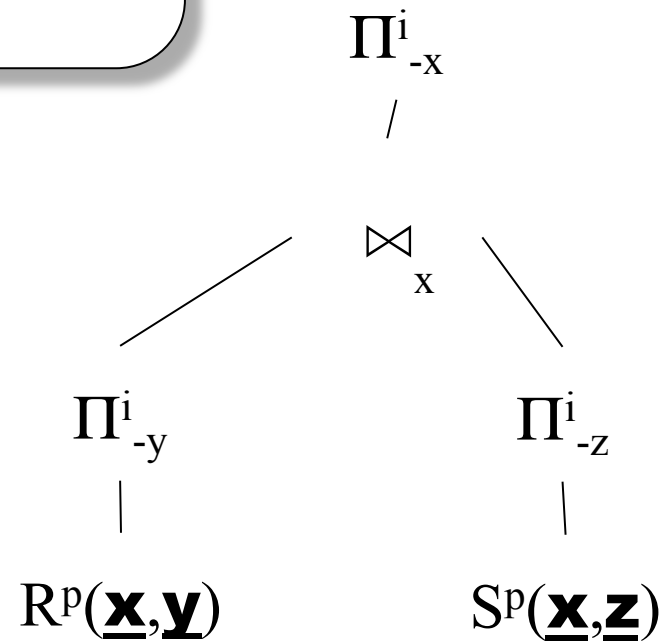
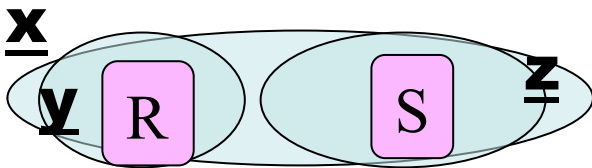
Algorithm: convert a Hierarchy to a Safe Plan

Root variable $u \rightarrow \Pi_{-u}^i$

Connected components \rightarrow Join

Single subgoal \rightarrow Leaf node

$q = R(\underline{x}, \underline{y}), S(\underline{x}, \underline{z})$



Independent project

$$P(q) =$$

$$1 - (1 - p_1(1 - (1 - q_1)(1 - q_2))) * (1 - p_2(1 - (1 - q_3)(1 - q_4)(1 - q_5)))$$

$$q =$$

$$R(\underline{x}, \underline{y}), S(\underline{x}, \underline{z})$$

 Π_{-x}
 \bowtie_x
 Π_{-y}
 Π_{-z}
 $R^p(\underline{x}, \underline{y})$
 $S^p(\underline{x}, \underline{z})$

A	P
a ₁	$p_1(1 - (1 - q_1)(1 - q_2))$
a ₂	$p_2(1 - (1 - q_3)(1 - q_4)(1 - q_5))$

A	P
a ₁	$1 - (1 - q_1)(1 - q_2)$
a ₂	$1 - (1 - q_3)(1 - q_4)(1 - q_5)$

<u>A</u>	<u>C</u>	P
a ₁	c ₁	q ₁
a ₁	c ₂	q ₂
a ₂	c ₃	q ₃
a ₂	c ₄	q ₄
a ₂	c ₅	q ₅

<u>A</u>	<u>B</u>	P
a ₁	b ₁	p ₁
a ₂	b ₂	p ₂

[D&S'2004]

The #P-Hard Queries

Are precisely the non-hierarchical queries. Example:

hd1 :- R(\mathbf{x}), S(\mathbf{x} , \mathbf{y}), T(\mathbf{y})

More general:

q :- ..., R(\mathbf{x} , ...), S(\mathbf{x} , \mathbf{y} , ...), T(\mathbf{y} , ...) , ...

Theorem Testing if q is PTIME or #P-hard is in AC^0

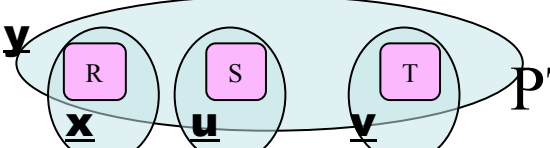
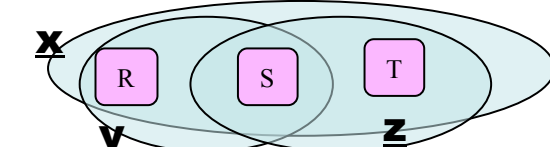
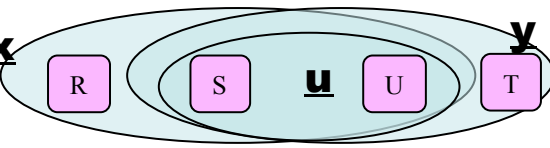
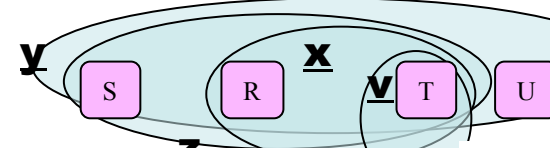
Quiz: What is their complexity ?

q	P TIME or #P ?
$R(\underline{x}, \underline{y}), S(\underline{y}, \underline{a}, \underline{u}), T(\underline{y}, \underline{y}, \underline{v})$	
$R(\underline{x}, \underline{y}), S(\underline{x}, \underline{y}, \underline{z}), T(\underline{x}, \underline{z})$	
$R(\underline{x}, \underline{a}), S(\underline{y}, \underline{u}, \underline{x}), T(\underline{u}, \underline{y}), U(\underline{x}, \underline{y})$	
$R(\underline{x}, \underline{y}, \underline{z}), S(\underline{z}, \underline{u}, \underline{y}), T(\underline{y}, \underline{v}, \underline{z}, \underline{x}), U(\underline{y})$	

Hint...

q	PTIME or #P ?
$R(\underline{x}, \underline{y}), S(\underline{y}, \underline{a}, \underline{u}), T(\underline{y}, \underline{y}, \underline{v})$	
$R(\underline{x}, \underline{y}), S(\underline{x}, \underline{y}, \underline{z}), T(\underline{x}, \underline{z})$	
$R(\underline{x}, \underline{a}), S(\underline{y}, \underline{u}, \underline{x}), T(\underline{u}, \underline{y}), U(\underline{x}, \underline{y})$	
$R(\underline{x}, \underline{y}, \underline{z}), S(\underline{z}, \underline{u}, \underline{y}), T(\underline{y}, \underline{v}, \underline{z}, \underline{x}), U(\underline{y})$	

...Answer

q	PTIME or #P ?
$R(\underline{x}, \underline{y}), S(\underline{y}, \underline{a}, \underline{u}), T(\underline{y}, \underline{y}, \underline{v})$	
$R(\underline{x}, \underline{y}), S(\underline{x}, \underline{y}, \underline{z}), T(\underline{x}, \underline{z})$	
$R(\underline{x}, \underline{a}), S(\underline{y}, \underline{u}, \underline{x}), T(\underline{u}, \underline{y}), U(\underline{x}, \underline{y})$	
$R(\underline{x}, \underline{y}, \underline{z}), S(\underline{z}, \underline{u}, \underline{y}), T(\underline{y}, \underline{v}, \underline{z}, \underline{x}), U(\underline{y})$	

Summary

- We have discussed only the simplest case: CQ w/o self-joins, on independent dbs
- Next time:
 - Add FDs at the representation level
 - Extend to independent/disjoint dbs
 - Extend to arbitrary CQs (with self-joins)