# Brief Tutorial on Probabilistic Databases

Dan Suciu

University of Washington

# About This Talk

- Probabilistic databases
  - Tuple-independent
  - Query evaluation
- Statistical relational models
  - Representation, learning, inference in FO
  - Reasoning/learning = lifted inference
- Sources:
  - Book 2011 [S.,Olteanu,Re,Koch]
  - Upcoming F&T survey [van Den Broek,S]

# Background: Relational databases

Smoker

| X | Y |
|---|---|
| Alice | 2009 |
| Alice | 2010 |
| Bob | 2009 |
| Carol | 2010 |

Friend

| X | Z |
|---|---|
| Alice | Bob |
| Alice | Carol |
| Bob | Carol |
| Carol | Bob |

Database    D =

Query: $Q(z) = \exists x \, (Smoker(x,'2009') \wedge Friend(x,z))$

$Q(D) =$

| Z |
|---|
| Bob |
| Carol |

Constraint:

$Q = \forall x \, (Smoker(x,'2010') \Rightarrow Friend(x,'Bob'))$

$Q(D) = $ true

# Probabilistic Database

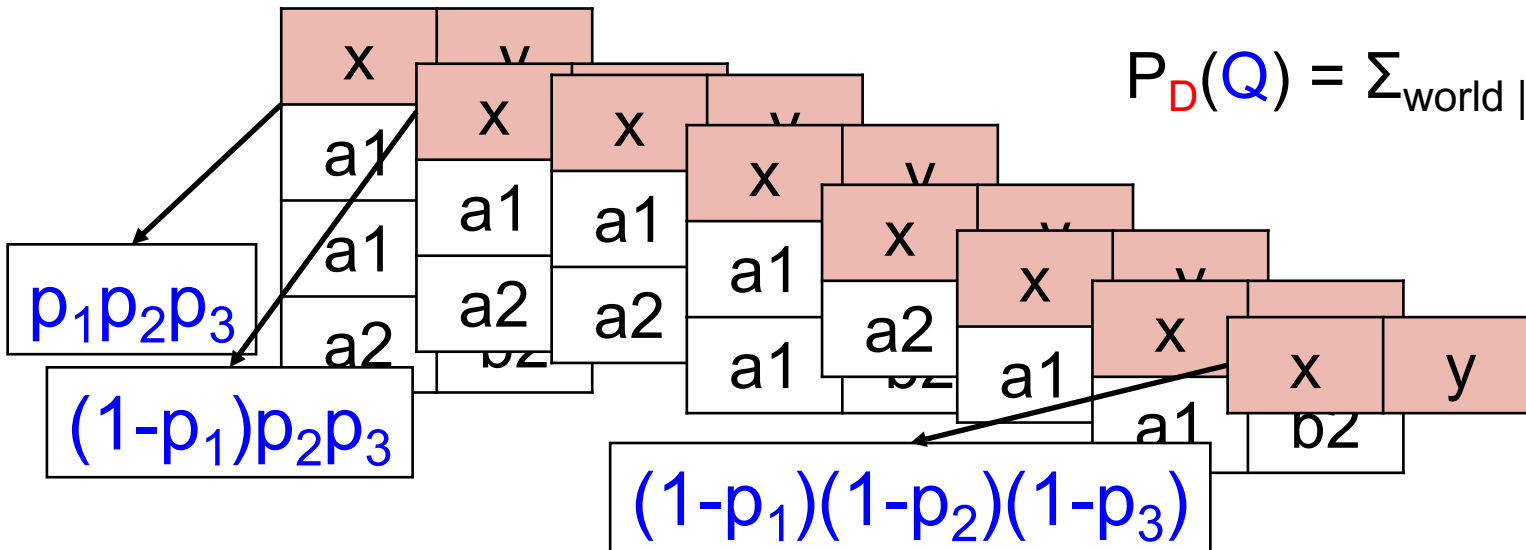Probabilistic database D:

| x | y | P |
|---|---|---|
| a1 | b1 | $p_1$ |
| a1 | b2 | $p_2$ |
| a2 | b2 | $p_3$ |

Possible worlds semantics:

$\Sigma_{world} P_D(world) = 1$

$P_D(Q) = \Sigma_{world \models Q} P_D(world)$

| x | y |
|---|---|
| a1 | |
| a1 | |
| a2 | |

$p_1 p_2 p_3$

$(1-p_1)p_2 p_3$

| x | y |
|---|---|
| a1 | |
| a2 | |

| x | y |
|---|---|
| a1 | |
| a2 | |

| x | y |
|---|---|
| a1 | |
| a1 | |

| x | y |
|---|---|
| a2 | |

| x | y |
|---|---|
| a1 | |

| x | y |
|---|---|
| b2 | |

| x | y |
|---|---|

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

# Outline

- Model Counting

- Small Dichotomy Theorem

- Dichotomy Theorem

- Query Compilation

- Conclusions, Open Problems

# Model Counting

- Given propositional Boolean formula F, compute the number of models #F

**Example**:

F = (X1 $\lor$ X2) $\land$ (X2 $\lor$ X3) $\land$ (X3 $\lor$ X1)

#F = 4

| X1 | X2 | X3 | F |
|----|----|----|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | **1** |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | **1** |
| 1 | 1 | 0 | **1** |
| 1 | 1 | 1 | **1** |

[Valiant'79]  #P-hard, even for 2CNF

# Probability of a Formula

- Each variable $X$ has a probability $p(X)$;
- $P(F)$ = probability that $F$=true,
  when each $X$ is set to true independently

**Example**:

$F = (X1 \lor X2) \land (X2 \lor X3) \land (X3 \lor X1)$

$P(F) = (1-p1)*p2*p3 +$
$\qquad p1*(1-p2)*p3 +$
$\qquad p1*p2*(1-p3) +$
$\qquad p1*p2*p3$

If $p(X) = \frac{1}{2}$ for all $X$, then $P(F) = \#F / 2^n$

| X1 | X2 | X3 | F |
|----|----|----|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | **1** |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | **1** |
| 1 | 1 | 0 | **1** |
| 1 | 1 | 1 | **1** |

# Algorithms for Model Counting

[Gomes, Sabharwal, Selman'2009]
Based on full search DPLL:

- Shannon expansion.
$$\#F = \#F[X=0] + \#F[X=1]$$

- Caching.
Store #F, look it up later

- Components. If Vars(F1) ∩ Vars(F2) = ∅:
$$\#(F1 \wedge F2) = \#F1 * \#F2$$

# Relational Representation (1/2)

- Fix an FO sentence $Q$ and a domain $\Delta$

- Ground atom $\rightarrow$ Boolean variable

**Definition** The lineage $F_{Q,\Delta}$ is:

$F_{Q,\Delta} = Q$          if $Q$ = ground atom

$F_{Q1 \wedge Q2,\Delta} = F_{Q1,\Delta} \wedge F_{Q2,\Delta}$    same for $\vee$, $\rightarrow$, $\neg$

$F_{\forall x.Q,\Delta} = \bigwedge_{a \in \Delta} F_{Q[a/x],\Delta}$

$F_{\exists x.Q,\Delta} = \bigvee_{a \in \Delta} F_{Q[a/x],\Delta}$

$Q = \forall x (Student(x) \Rightarrow Person(x))$

$F_{Q,[n]} = (Student(1) \Rightarrow Person(1)) \wedge \ldots \wedge (Student(n) \Rightarrow Person(n))$

# Relational Representation (2/2)

- For a database $D$, denote

$$F_{Q,D} = F_{Q,domain(D)}$$

  where all tuples not in $D$ are set to *false*

- $F_{Q,\Delta}$ or $F_{Q,D}$ is called the *lineage* or the *provenance* or the *grounding* of $Q$

# Weighted FO Model Counting

- Probabilities of ground atoms in $D$ = probabilities of Boolean variables p($X$)

- Fix $Q$.  Given $D$, compute $P(F_{Q,D})$

| x | y | P |
|---|---|---|
| a1 | b1 | $p_1$ |
| a1 | b2 | $p_2$ |
| a2 | b2 | $p_3$ |

- Simple fact: $P_D(Q) = P(F_{Q,D})$

# This Talk

Fix a query $Q$:

- What is the complexity of $P_D(Q)$ in the size of $D$?

- What is the best runtime of a DPLL-based algorithm on $F_{Q,D}$ in the size of $D$?

# Discussion: Correlations

[Domingos&Richardson'06]  MLN = popular FO framework for Machine Learning tasks

Lise Getoor's talk today

Smoker(x) $\wedge$ Friends(x,y) $\rightarrow$ Smoker(y), weight = 2.3

**Theorem** [Jha,S'11] One can construct effectively D s.t.
$P_{MLN, \Delta}(Q) = P_D(Q \mid \Gamma) = P_D(Q \wedge \Gamma) / P_D(\Gamma)$

# Outline

- Model Counting

- Small Dichotomy Theorem

- Dichotomy Theorem

- Query Compilation

- Conclusions, Open Problems

# Background: Query Plans

Q(z) = R(z,x), S(x,y),T(y,u), u=123

Query plan = expressions
over the input relation

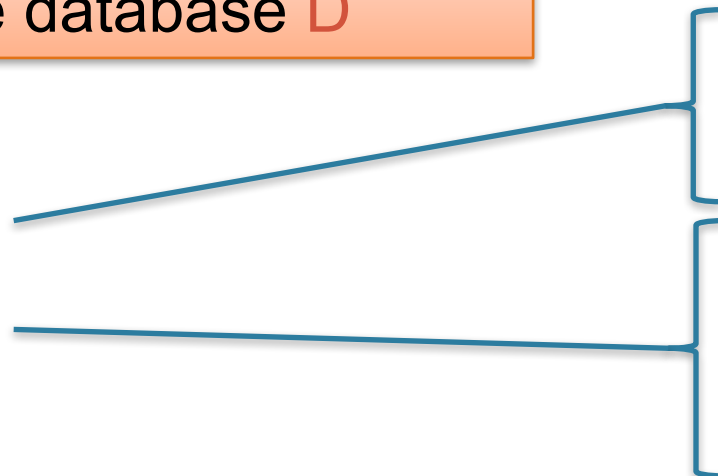Operators = selection, projection,
join, union, difference

$\Pi_z$

$\bowtie_x$

$\sigma_{u=123}$

$\bowtie_y$

R(z,x)    S(x,y)    T(y,u)

# An Example

$Q() = R(x), S(x,y)$ $\qquad = \exists x \exists y(R(x) \wedge S(x,y))$

$P_D(Q) = 1- \{1- p1*[\ 1-(1-q1)*(1-q2)\ ]\} *$

$\qquad\qquad \{1- p2*[\ 1-(1-q3)*(1-q4)*(1-q5)]\}$

One can compute $P_D(Q)$ in PTIME in the size of the database D

R

| x | P |
|---|---|
| a1 | p1 |
| a2 | p2 |
| a3 | p3 |

S

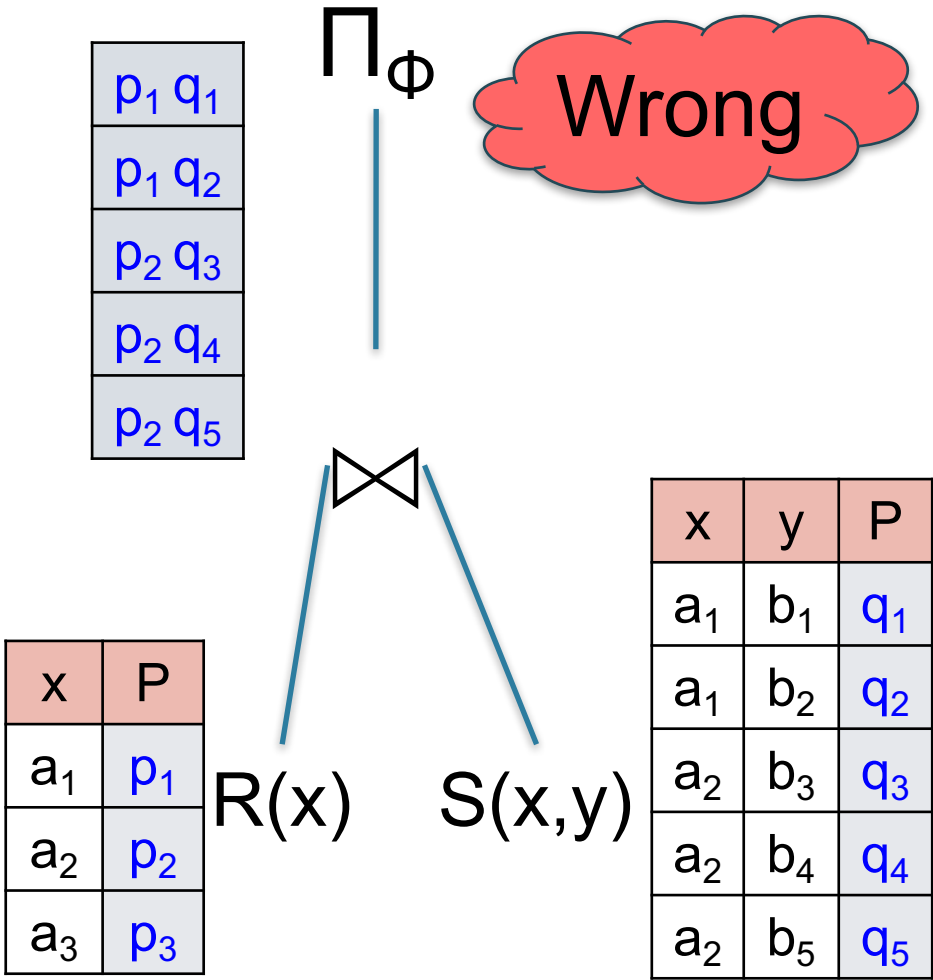| x | y | P |
|---|---|---|
| a1 | b1 | q1 |
| a1 | b2 | q2 |
| a2 | b3 | q3 |
| a2 | b4 | q4 |
| a2 | b5 | q5 |

# Extensional Plans

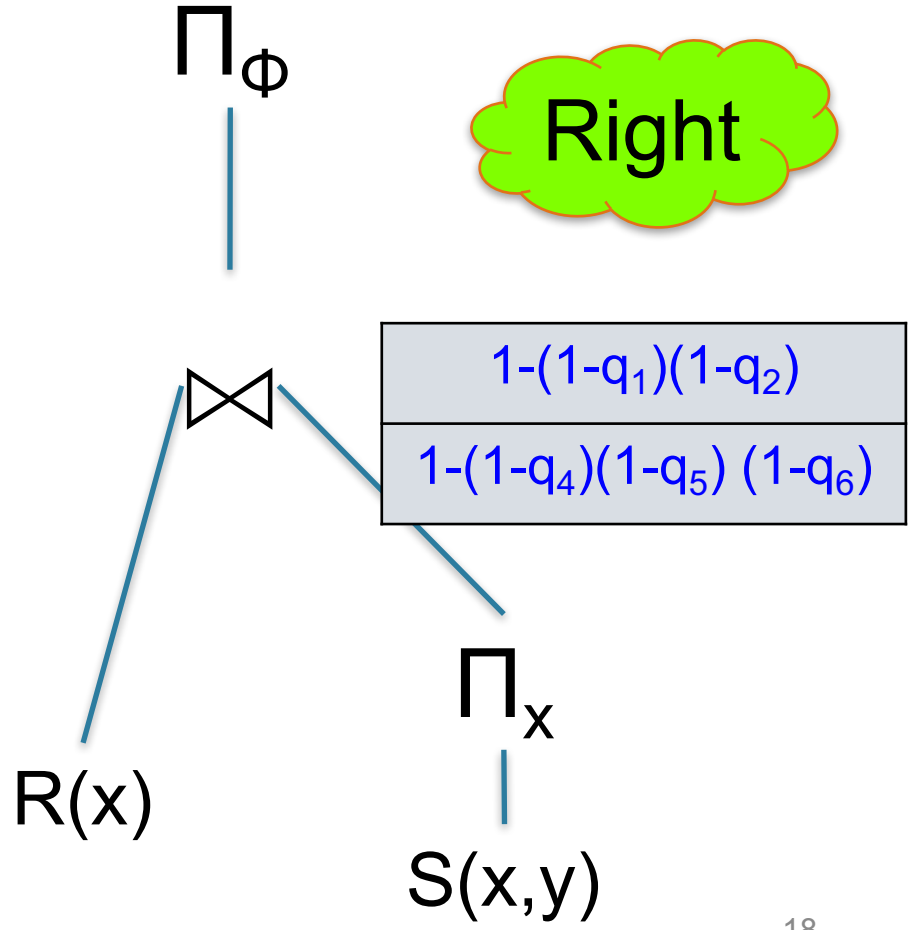- Modify each operator to compute output probabilities, assuming independent events

$$Q() = R(x), S(x,y)$$

$$P(Q) = 1 - [1-p_1*(1-(1-q_1)*(1-q_2))]$$
$$*[1- p_2*(1-(1-q_3)*(1-q_4)*(1-q_5))]$$

$$1-(1-p_1q_1)(1-p_1q_2)(1-p_2q_3)(1-p_2q_4)(1-p_2q_5)$$

$$1-\{1-p_1[1-(1-q_1)(1-q_2)]\}*$$
$$\{1-p_2[1-(1-q_4)(1-q_5)(1-q_6)]\}$$

$\Pi_\Phi$

Wrong

$\Pi_\Phi$

Right

| $p_1 q_1$ |
|---|
| $p_1 q_2$ |
| $p_2 q_3$ |
| $p_2 q_4$ |
| $p_2 q_5$ |

| | $1-(1-q_1)(1-q_2)$ |
|---|---|
| | $1-(1-q_4)(1-q_5)(1-q_6)$ |

$\bowtie$

| x | y | P |
|---|---|---|
| $a_1$ | $b_1$ | $q_1$ |
| $a_1$ | $b_2$ | $q_2$ |
| $a_2$ | $b_3$ | $q_3$ |
| $a_2$ | $b_4$ | $q_4$ |
| $a_2$ | $b_5$ | $q_5$ |

$\bowtie$

R(x)       S(x,y)

R(x)

$\Pi_x$

| x | P |
|---|---|
| $a_1$ | $p_1$ |
| $a_2$ | $p_2$ |
| $a_3$ | $p_3$ |

S(x,y)

# Safe Queries

**Definition** A plan for Q is _safe_ if it computes the probabilities correctly.

Q is _safe_ if it has a safe plan.

- In AI, computing Q using a safe plan is called _lifted inference_

- _Safe query_ = _Liftable query_
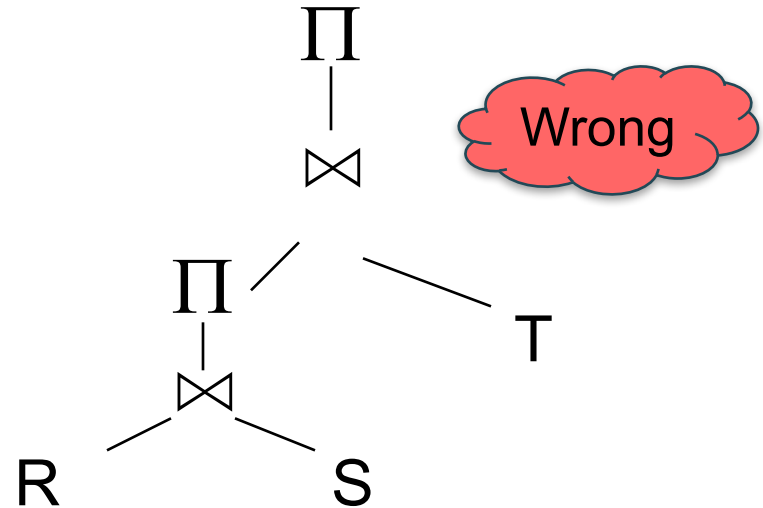
- If Q is safe then $P_D(Q)$ is in PTIME

# Unsafe Queries

R

| X | P |
|----|----|
| x1 | p1 |
| x2 | p2 |

S

| X | Y |
|----|----|
| x1 | y1 |
| x1 | y2 |
| x2 | y2 |

T

| Y | P |
|----|----|
| y1 | q1 |
| y2 | q2 |

$\Pi$
|
$\bowtie$

Wrong

$\Pi$
|
$\bowtie$

T

R     S

$H_0() = R(x), S(x,y), T(y)$

**Theorem.** [Dalvi&S.2004] $P_D(H_0)$ is #P-hard

However:
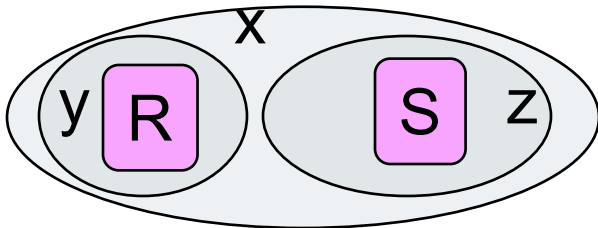1. This plan computes an upper bound [VLDB'15]
2. Use samples on T [VLDB'16]

Wolfgang Gatterbauer's talk today

# Hierarchical Queries

Fix Q;  at(x) = set of atoms (=literals) containing the variable x

**Definition**  Q is hierarchical  if forall variables x, y:
$at(x) \subseteq at(y)$   or   $at(x) \supseteq at(y)$   or   $at(x) \cap at(y) = \varnothing$
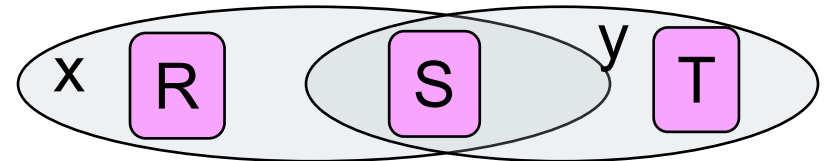
Hierarchical

Q() = R(x,y), S(x,z)

Non-hierarchical

$H_0$() =  R(x), S(x,y), T(y)

# The Small Dichotomy Theorem

Non-repeating Conjunctive Query =
        = Conjunctive Query "without self-joins"
        = "Simple" conjunctive query

[Dalvi&S.04]

**__Theorem__** Let $Q$ be a non-repeating CQ
- If $Q$ is hierarchical, then $P_D(Q)$ is in PTIME.
- If $Q$ is not hierarchical then $P_D(Q)$ is #P-hard.

By duality, the same holds for a non-repeating clause

# Summary so Far

| Complexity of $P_D(Q)$ | Non-repeating CQ<br>Non-repeating clause |
|---|---|
| PTIME | Hierarchical |
| #P - hard | Non-hierarchical |

# Outline

- Model Counting

- Small Dichotomy Theorem

- Dichotomy Theorem

- Query Compilation

- Conclusions, Open Problems

# The Rules for Lifted Inference

Preprocess Q (omitted from this talk; see book),
then apply these rules (some have preconditions)

$$P(\neg Q) = 1 - P(Q)$$ negation

$$P(Q1 \wedge Q2) = P(Q1)P(Q2)$$
$$P(Q1 \vee Q2) = 1 - (1 - P(Q1))(1 - P(Q2))$$

Independent
join / union

$$P(\exists z\ Q) = 1 - \Pi_{a \in \text{Domain}}\ (1 - P(Q[a/z]))$$
$$P(\forall z\ Q) = \Pi_{a \in \text{Domain}}\ P(Q[a/z])$$

Independent project

$$P(Q1 \wedge Q2) = P(Q1) + P(Q2) - P(Q1 \vee Q2)$$
$$P(Q1 \vee Q2) = P(Q1) + P(Q2) - P(Q1 \wedge Q2)$$

Inclusion/
exclusion

# FO$^{un}$ = Unate FO

An FO sentence is *unate* if:

- Negations occur only on atoms

- Every relational symbol R either occurs only positively, or only negatively

FO$^{un}$ = FO restricted to unate sentences

# Dichotomy Theorem

[Dalvi&S'12]

**Theorem** For any $Q$ in $\forall^* FO^{un}$ (or $\exists^* FO^{un}$)
- If rules succeed, then $P_D(Q)$ in PTIME in $|D|$
- If rules fail, then $P_D(Q)$ is #P hard in $|D|$

Note: Unions of Conjunctive queries (UCQ)
is essentially $\exists^* FO^{un}$

# Example: Liftable Query

$Q_J() = S(x_1,y_1), R(y_1), S(x_2,y_2), T(y_2)$

$= [S(x_1,y_1),R(y_1)] \wedge [S(x_2,y_2),T(y_2)]$

$Q_1$        $Q_2$

$P(Q_J) = P(Q_1) + P(Q_2) - P(Q_1 \vee Q_2)$

PTIME (have seen before)

$y = y1 = y2$

$Q_1 \vee Q_2 = \exists y [S(x_1,y),R(y) \vee S(x_2,y)),T(y)]$

$P(Q_1 \vee Q_2) =$
$= 1 - \Pi_{b \in Domain} (1 - P[S(x_1,b),R(b) \vee S(x_2,b)),T(b)])$
$= 1 - \Pi_{b \in Domain} (1 - P[S(x_1,b)] * P[R(b) \vee T(b)]) = \ldots$ etc

Runtime $= O(n^2)$.

# Example: Liftable Query

$$Q_J = \forall x_1 \forall y_1 \forall x_2 \forall y_2 \, (S(x_1,y_1) \lor R(y_1) \lor S(x_2,y_2) \lor T(y_2))$$

$$= [\forall x_1 \forall y_1 S(x_1,y_1) \lor R(y_1)] \lor [\forall x_2 \forall y_2 S(x_2,y_2) \lor T(y_2)]$$

$Q_1$         $Q_2$

$$P(Q_J) = P(Q_1) + P(Q_2) - P(Q_1 \land Q_2)$$

PTIME (have seen before)

$y = y1 = y2$

$$Q_1 \land Q_2 = \forall y \, [(\forall x_1 S(x_1,y) \lor R(y)) \land (\forall x_2 S(x_2,y)) \lor T(y)]$$
$$= \forall y \, [\forall x \, S(x,y) \lor (R(y) \land T(y))]$$

$$P(Q_1 \land Q_2) = \Pi_{b \in \text{Domain}} \, P[\forall x.S(x,b) \lor (R(b) \land T(b))] = \ldots \text{etc}$$

Runtime $= O(n^2)$.

# Unliftable Queries $H_k$

Will drop $\forall$ to reduce clutter

$H_0 = R(x) \vee S(x,y) \vee T(y)$

$H_1 = [R(x_0) \vee S(x_0,y_0)] \wedge [S(x_1,y_1) \vee T(y_1)]$

$H_2 = [R(x_0) \vee S_1(x_0,y_0)] \wedge [S_1(x_1,y_1) \vee S_2(x_1,y_1)] \vee [S_2(x_2,y_2) \vee T(y_2)]$

$H_3 = [R(x_0) \vee S_1(x_0,y_0)] \wedge [S_1(x_1,y_1) \vee S_2(x_1,y_1)] \wedge [S_2(x_2,y_2) \vee S_3(x_2,y_2)] \wedge [S_3(x_3,y_3) \vee T(y_3)]$
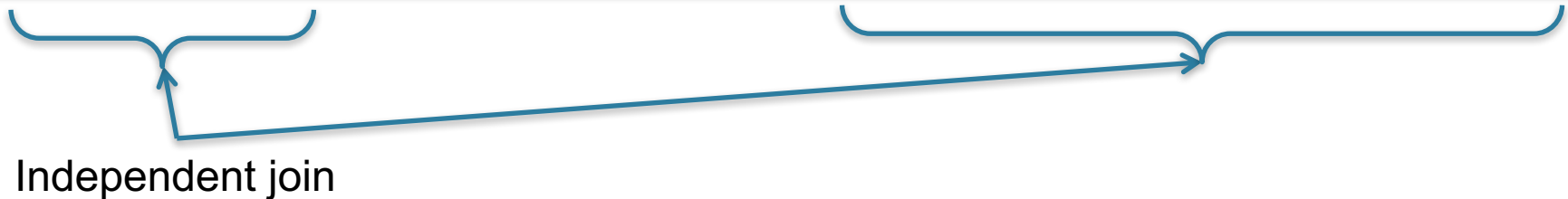
. . .

Every $H_k$, $k \geq 1$ is hierarchical

**Theorem.** [Dalvi&S'12] Every query $H_k$ is #P-hard

# A Closer Look at $H_k$

If we drop any one clause → in PTIME

$H_3 = [R(x_0) \lor S_1(x_0,y_0)] \land [S_1(x_1,y_1) \lor S_2(x_1,y_1)] \land [S_2(x_2,y_2) \lor S_3(x_2,y_2)] \land [S_3(x_3,y_3) \lor T(y_3)]$

Independent join

# Summary so Far

| Complexity of $P_D$(Q) | Non-repeating CQ Non-repeating clauses | $\exists^*FO^{un}$ $\forall^*FO^{un}$ |
|---|---|---|
| PTIME | Hierarchical | Rules succeed |
| #P - hard | Non-hierarchical | Rules fail |

# Outline

- Model Counting

- Small Dichotomy Theorem

- Dichotomy Theorem

- Query Compilation

- Conclusions, Open Problems

# Lifted v.s. Grounded Inference

- To compute $P_D(Q)$:
  compute the lineage $F_{Q,D}$
  use DPLL-based algorithm for $P(F_{Q,D})$

- For which queries $Q$ can this be in PTIME?

- [Huang&Darwiche'2005] The _trace_ of a DPLL-based algorithm is "decision-DNNF"
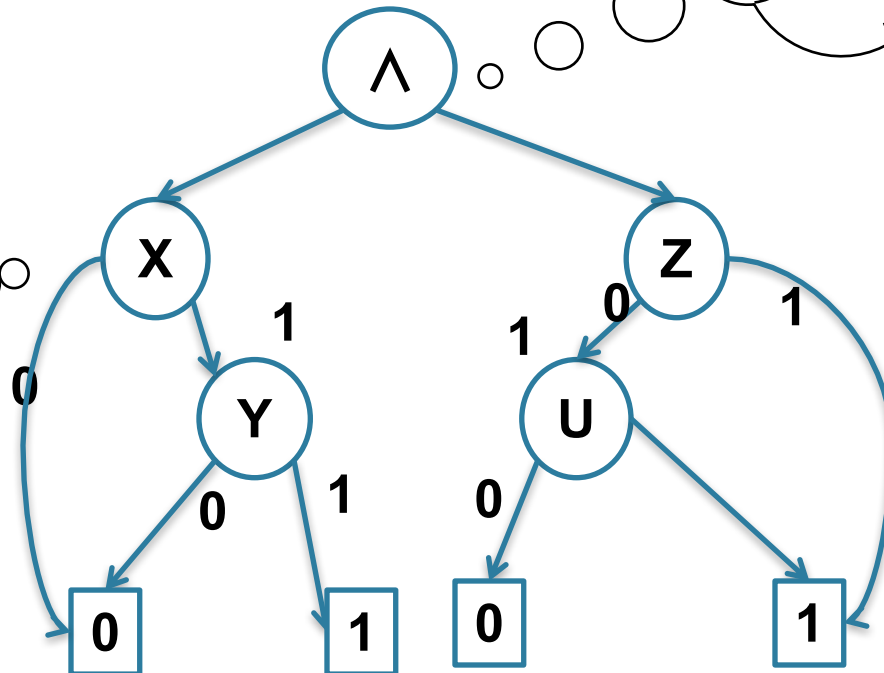
# Decision-DNNF

**Def** [Darwiche] A Decision-DNNF is a rooted DAG where:
- Internal nodes are decision or $\wedge$
- Sink nodes are 0 or 1

Children of $\wedge$ have disjoint sets of variables

Every root-to-sink path reads each variable at most once

# Notations

$$H_{k0} = \forall x \forall y\, R(x) \vee S_1(x,y)$$
$$H_{k1} = \forall x \forall y\, S_1(x,y) \vee S_2(x,y)$$
$$H_{k2} = \forall x \forall y\, S_2(x,y) \vee S_3(x,y)$$
$$\ldots$$
$$\ldots$$
$$H_{kk} = \forall x \forall y\, S_k(x,y) \vee T(y)$$

$$f(Z_0, Z_1, \ldots, Z_k) = \text{a Boolean function}$$

$$Q = f(H_{k0}, H_{k1}, \ldots, H_{kk})$$

Example: $f = Z_0 \wedge Z_1 \wedge \ldots \wedge Z_k$  then  $f(H_{k0}, H_{k1}, \ldots, H_{kk}) = H_k$

# Easy/Hard Queries

[Beame'14]

**Theorem** Let $Q = f(H_{k0}, H_{k1}, \ldots, H_{kk})$ where $f(Z_0, Z_1, \ldots, Z_k)$ is a monotone Boolean function.
- Any Decision-DNNF for $F_{Q,[n]}$ has size $2^{\Omega(\sqrt{n})}$.
- $P_D(Q)$ is in PTIME iff $\mu_Q(0, 1) = 0$

$\mu$ = Möbius function of the implicates lattice of $Q$

Consequence: Any DPLL-based algorithm takes time $2^{\Omega(\sqrt{n})}$, even if the query is in PTIME!

# Cancellations

$$Q_W = (H_{30} \wedge H_{32}) \vee$$
$$(H_{30} \wedge H_{33}) \vee$$
$$(H_{31} \wedge H_{33})$$

$$H_{30} = \forall x \forall y \, R(x) \vee S_1(x,y)$$
$$H_{31} = \forall x \forall y \, S_1(x,y) \vee S_2(x,y)$$
$$H_{32} = \forall x \forall y \, S_2(x,y) \vee S_3(x,y)$$
$$H_{33} = \forall x \forall y \, S_3(x,y) \vee T(y)$$

$$P(Q_W) = P(H_{30} \wedge H_{32}) + P(H_{30} \wedge H_{33}) + P(H_{31} \wedge H_{33}) +$$

$$- P(H_{30} \wedge H_{32} \wedge H_{33}) - P(H_{30} \wedge H_{31} \wedge H_{33})$$

$$- P(H_{30} \wedge H_{31} \wedge H_{32} \wedge H_{33})$$

$$= H_3 \text{ (hard !)}$$

$$+ P(H_{30} \wedge H_{31} \wedge H_{32} \wedge H_{33})$$
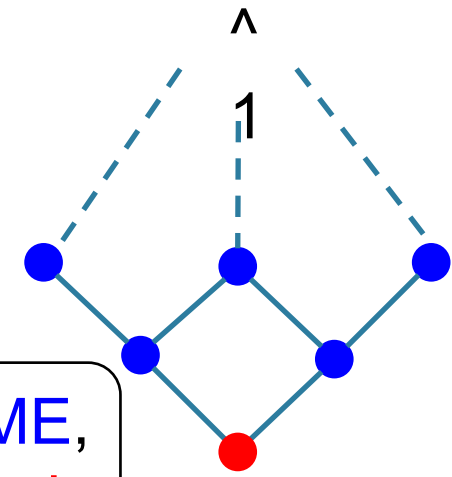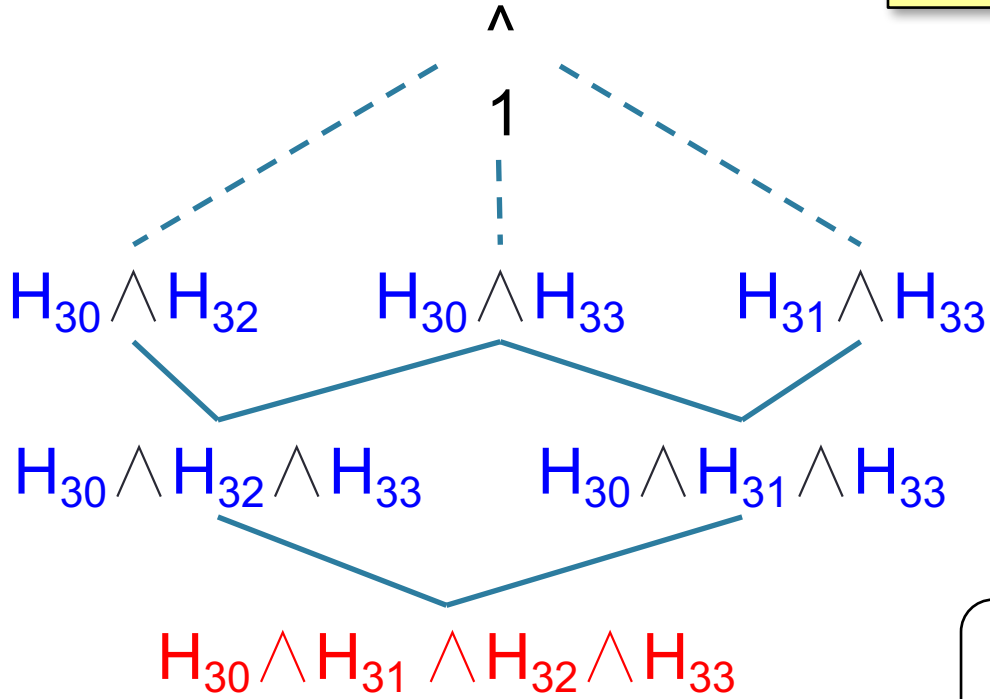
Also $= H_3$

$$P(Q_W) \text{ is in PTIME}$$

# The CNF Lattice

$$Q_W = (H_{30} \wedge H_{32}) \vee$$
$$(H_{30} \wedge H_{33}) \vee$$
$$(H_{31} \wedge H_{33})$$

**Definition**. The DNF lattice of
$Q = Q_1 \vee Q_2 \vee \ldots$ is:
- Elements are prime implicants
- Order is implication

$\wedge$

$1$

$H_{30} \wedge H_{32}$      $H_{30} \wedge H_{33}$      $H_{31} \wedge H_{33}$

$H_{30} \wedge H_{32} \wedge H_{33}$      $H_{30} \wedge H_{31} \wedge H_{33}$

$H_{30} \wedge H_{31} \wedge H_{32} \wedge H_{33}$

$\wedge$

$1$

Nodes • in PTIME,
Nodes • #P hard.

# The Möbius' Function

**Def**. The Möbius function:

$$\mu(\hat{1},\hat{1}) = 1 \qquad \mu(u,\hat{1}) = -\sum_{u < v \leq \hat{1}} \mu(v,\hat{1})$$

**Möbius' Inversion Formula**:

$$P(Q) = -\sum_{Qi < \hat{1}} \mu(Qi,\hat{1}) \, P(Qi)$$

# The Möbius' Function

**Def**. The Möbius function:

$$\mu(\hat{1}, \hat{1}) = 1 \qquad \mu(u, \hat{1}) = -\sum_{u < v \leq \hat{1}} \mu(v, \hat{1})$$

**Möbius' Inversion Formula**:

$$P(Q) = -\sum_{Qi < \hat{1}} \mu(Qi, \hat{1}) \, P(Qi)$$

# The Möbius' Function

**Def**. The Möbius function:

$$\mu(\hat{1}, \hat{1}) = 1 \qquad \mu(u, \hat{1}) = -\sum_{u < v \le \hat{1}} \mu(v, \hat{1})$$

**Möbius' Inversion Formula**:

$$P(Q) = -\sum_{Qi < \hat{1}} \mu(Qi, \hat{1}) \, P(Qi)$$

# The Möbius' Function

**Def**. The Möbius function:

$$\mu(\hat{1}, \hat{1}) = 1 \qquad \mu(u, \hat{1}) = - \sum_{u < v \leq \hat{1}} \mu(v, \hat{1})$$

**Möbius' Inversion Formula**:

$$P(Q) = - \sum_{Qi < \hat{1}} \mu(Qi, \hat{1}) \, P(Qi)$$

# The Möbius' Function

**Def**. The Möbius function:

$$\mu(\hat{1}, \hat{1}) = 1 \qquad \mu(u, \hat{1}) = -\sum_{u < v \leq \hat{1}} \mu(v, \hat{1})$$

**Möbius' Inversion Formula**:

$$P(Q) = -\sum_{Qi < \hat{1}} \mu(Qi, \hat{1}) \, P(Qi)$$

# The Möbius' Function

**Def**. The Möbius function:

$$\mu(\hat{1}, \hat{1}) = 1 \qquad \mu(u, \hat{1}) = - \sum_{u < v \leq \hat{1}} \mu(v, \hat{1})$$

**Möbius' Inversion Formula**:

$$P(Q) = \sum_{Qi < \hat{1}} \mu(Qi, \hat{1}) \, P(Qi)$$

# Summary

| | nr CQ nr clause | $\exists^*FO^{un}$ $\forall^*FO^{un}$ | $f(H_{k0},\ldots,H_{kk})$ |
|---|---|---|---|
| PTIME | Hierarchical | Rules succeed | $\mu(0,1) = 0$ |
| #P - hard | Non-hierarchical | Rules fail | $\mu(0,1) \neq 0$ |
| DPLL | | | $2^{\Omega(\sqrt{n})}$ |

# Möbius Über Alles



Non-hierarchical $H_0$

#P-hard  hierarchical $H_1$

$H_2$

$Q_W$  $Q_9$  $H_3$

PTIME

Poly-size FBDD, dec-DNNF

$Q_V$

Poly-size OBDD, SDD =
= inversion-free

$Q_J$

Read Once  $Q_U$

$\forall^* FO^{un}$ or
$\exists^* FO^{un}$ ($\approx$ UCQ)

Open

# Outline

- Model Counting

- Small Dichotomy Theorem

- Dichotomy Theorem

- Query Compilation

- Conclusions, Open Problems

# Summary

- Query evaluation on probabilistic databases = weighted model counting

- Each query Q defines a different WMC problem

- Dichotomy: depending on Q, WMC is in PTIME or #P-hard

- Using a DPLL-based algorithm on the grounded Q is suboptimal

# Discussion: Extensions

Open problems: extend the dichotomy theorem to:

- Mixed probabilistic/deterministic relations
- Functional dependencies
- Interpreted predicates: <, ≠

Open problem: complexity of MAP

# Discussion: Symmetric Relations

- A relation R is *symmetric* if all ground tuples have the same probability

- [van den Broeck'14] For every Q in $FO^2$, **P**(Q) is in PTIME on symmetric databases.

- [Beame'15] Hardness results.

- In general the complexity is open

# Discussion: Negation

[Fink&Olteanu'14] Restrict FO to non-repeating exrpessions

Dan Olteanu's talk next

- Theorem Hierarchical expressions are in PTIME, non-hierarchical are #P-hard.

[Gribkoff,S.,v.d.Broeck'14]   $\forall^*$FO or $\exists^*$FO

- Need resolution compute some queries with negation

Open problem: completeness/dichotomy?

# Thank You!

# BACKUP

# Weighted Model Counting

- Each variable $X$ has a weight $w(X)$;
- Weight of a model $= \prod_{X=true} w(X)$
- WMC($F$) = sum of weights of models of $F$

**Example**:
$F = (X1 \vee X2) \wedge (X2 \vee X3) \wedge (X3 \vee X1)$

WMC($F$) = w2*w3 +
             w1*w3 +
             w1*w2 +
             w1*w2*w3

| X1 | X2 | X3 | F |
|----|----|----|---|
| 0  | 0  | 0  | 0 |
| 0  | 0  | 1  | 0 |
| 0  | 1  | 0  | 0 |
| 0  | 1  | 1  | 1 |
| 1  | 0  | 0  | 0 |
| 1  | 0  | 1  | 1 |
| 1  | 1  | 0  | 1 |
| 1  | 1  | 1  | 1 |

Set $w(X) = 1$: then WMC($F$) = #F

# Probability of a Formula

- Each variable $X$ has a probability $p(X)$;
- $P(F)$ = probability that $F$=true,
  when each $X$ is set to true independently

**Example**:
$F = (X1 \lor X2) \land (X2 \lor X3) \land (X3 \lor X1)$

$P(F) = (1-p1)*p2*p3 +$
$\quad\quad\quad p1*(1-p2)*p3 +$
$\quad\quad\quad p1*p2*(1-p3) +$
$\quad\quad\quad p1*p2*p3$

| X1 | X2 | X3 | F |
|----|----|----|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | **1** |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | **1** |
| 1 | 1 | 0 | **1** |
| 1 | 1 | 1 | **1** |

Set $w(X) = p(X)/(1-p(X))$
Then $P(F) = WMC(F) / Z$, where $Z = \Pi_X (1+w(X))$

# Discussion: Dichotomy for #SAT

- [Creignou&Hemann'96] consider model counting #F where the formula F is given by *generalized clauses*

- **Dichotomy Theorem**: #F is in PTIME when all clauses are affine, #P-hard otherwise

- Not helpful in our context:
  - If $Q$ is a UCQ, then the clauses of $F_{Q,D}$ are of the form $X \vee Y \vee Z$ … and are not affine;
  - But $P(F_{Q,D})$ is not always #P-hard, because $Q$ restricts the struture of the clauses

# Warm-up: Weights

Replace probabilities with weights:

R:

| x | y | w |
|----|----|-------|
| a1 | b1 | $w_1$ |
| a2 | b1 | $w_2$ |
| a3 | b2 | $w_3$ |

S:

| y | w |
|----|-------|
| b1 | $v_1$ |
| b2 | $v_2$ |
| b3 | $v_3$ |

$P_D(\text{world}) = \text{Weight(world)}/Z$

$Z = \Sigma_{\text{world'}}\ \text{Weight(world')}$

Weight of a possible world:

$$Z = (1+v_1)\,(1+v_2)\,(1+v_3) \\ (1+w_1)\,(1+w_2)\,(1+w_3)$$

R:

S:

Weight(

| x | y |
|----|----|
| a1 | b1 |
| a2 | b1 |

| y |
|----|
| b2 |

) = $w_1 w_2 v_2$

# Markov Logic Networks

Replace probabilities with weights:

Add soft constraints:

R:

| x | y | w |
|---|---|---|
| a1 | b1 | $w_1$ |
| a2 | b1 | $w_2$ |
| a3 | b2 | $w_3$ |

S:

| y | w |
|---|---|
| b1 | $v_1$ |
| b2 | $v_2$ |
| b3 | $v_3$ |

$R(x,y) \Rightarrow S(y)$    $w_4$

# Markov Logic Networks

Replace probabilities with weights:

Add soft constraints:

$$R(x,y) \Rightarrow S(y) \qquad w_4$$

R:

| x | y | w |
|----|----|-------|
| a1 | b1 | $w_1$ |
| a2 | b1 | $w_2$ |
| a3 | b2 | $w_3$ |

S:

| y | w |
|----|-------|
| b1 | $v_1$ |
| b2 | $v_2$ |
| b3 | $v_3$ |

Weight of a possible world:

Weight(

R:

| x | y |
|----|----|
| a1 | b1 |
| a3 | b2 |

S:

| y |
|----|
| b1 |

) = $w_1 w_3 v_1 w_4 w_4 w_4 w_4 w_4$

# Markov Logic Networks

Replace probabilities with weights:

Add soft constraints:

R:

| x | y | w |
|----|----|-------|
| a1 | b1 | $w_1$ |
| a2 | b1 | $w_2$ |
| a3 | b2 | $w_3$ |

S:

| y | w |
|----|-------|
| b1 | $v_1$ |
| b2 | $v_2$ |
| b3 | $v_3$ |

$$R(x,y) \Rightarrow S(y) \qquad w_4$$

$$P_{MLN}(\text{world}) = \text{Weight(world)}/Z$$

$$Z = \Sigma_{\text{world'}} \; \text{Weight(world')}$$

Weight of a possible world:

Weight(

R:

| x | y |
|----|----|
| a1 | b1 |
| a3 | b2 |

S:

| y |
|----|
| b1 |

) = $w_1 w_3 v_1 w_4 w_4 w_4 w_4 w_4$

# Markov Logic Networks

Replace probabilities with weights:

Add soft constraints:

R:

| x | y | w |
|---|---|---|
| a1 | b1 | $w_1$ |
| a2 | b1 | $w_2$ |
| a3 | b2 | $w_3$ |

S:

| y | w |
|---|---|
| b1 | $v_1$ |
| b2 | $v_2$ |
| b3 | $v_3$ |

$$R(x,y) \Rightarrow S(y) \qquad w_4$$

$$P_{MLN}(\text{world}) = \text{Weight(world)}/Z$$

$$Z = \Sigma_{\text{world'}} \ \text{Weight(world')}$$

$$Z = \frac{(1+v_1)(1+v_2)(1+v_3)}{(1+w_1)(1+w_2)(1+w_3)}$$

$Z$ is #P-hard to compute

Weight of a possible world:

Weight(

R:

| x | y |
|---|---|
| a1 | b1 |
| a3 | b2 |

S:

| y |
|---|
| b1 |

) = $w_1 w_3 v_1 w_4 w_4 w_4 w_4 w_4$

# Discussion

Weights v.s. probabilities

- Soft constraints *with probabilities* may be inconsistent

- Soft constraints *with weights (≠ 0, ∞)* always consistent

Weight values have no semantics!

- Learned from training data

Inconsistent:
S(x):      p=0.5
S(x)∧R(x):  p=0.9

Consistent:
S(x):      w=5
S(x)∧R(x): w=9

# MLN's to Tuple-Independent PDB

Replace probabilities with weights:

Soft constraint:

R:

| x | y | w |
|----|----|-------|
| a1 | b1 | $w_1$ |
| a2 | b1 | $w_2$ |
| a3 | b2 | $w_3$ |

S:

| y | w |
|----|-------|
| b1 | $v_1$ |
| b2 | $v_2$ |
| b3 | $v_3$ |

$R(x,y) \Rightarrow S(y)$    $w_4$

Replace with hard constraint:

# MLN's to Tuple-Independent PDB

Replace probabilities with weights:

Soft constraint:

R:

| x | y | w |
|---|---|---|
| a1 | b1 | $w_1$ |
| a2 | b1 | $w_2$ |
| a3 | b2 | $w_3$ |

S:

| y | w |
|---|---|
| b1 | $v_1$ |
| b2 | $v_2$ |
| b3 | $v_3$ |

$$R(x,y) \Rightarrow S(y) \qquad w_4$$

Replace with hard constraint:

$$\Gamma \equiv \forall x \forall y$$
$$A(x,y) \Leftrightarrow (R(x,y) \Rightarrow S(y))$$

New relation A:

| x | y | w |
|---|---|---|
| a1 | b1 | $w_4$ |
| a1 | b2 | $w_4$ |
| a1 | b3 | $w_4$ |
| a1 | b1 | $w_4$ |
| … | | |

# MLN's to Tuple-Independent PDB

Replace probabilities with weights:

Soft constraint:

R:

| x | y | w |
|---|---|---|
| a1 | b1 | $w_1$ |
| a2 | b1 | $w_2$ |
| a3 | b2 | $w_3$ |

S:

| y | w |
|---|---|
| b1 | $v_1$ |
| b2 | $v_2$ |
| b3 | $v_3$ |

$R(x,y) \Rightarrow S(y)$   $w_4$

Replace with hard constraint:

$\Gamma \equiv \forall x \forall y$
$A(x,y) \Leftrightarrow (R(x,y) \Rightarrow S(y))$

New relation A:

| x | y | w |
|---|---|---|
| a1 | b1 | $w_4$ |
| a1 | b2 | $w_4$ |
| a1 | b3 | $w_4$ |
| a1 | b1 | $w_4$ |
| … | | |

Weight(

R:

| x | y |
|---|---|
| a1 | b1 |
| a3 | b2 |

S:

| y |
|---|
| b1 |

A:

| x | y |
|---|---|
| a1 | b1 |
| a1 | b2 |
| a2 | b1 |
| a2 | b2 |
| a3 | b1 |

) = $w_1 w_3 v_1 w_4 w_4 w_4 w_4 w_4$

# MLN's to Tuple-Independent PDB

Replace probabilities with weights:

Soft constraint:

$R(x,y) \Rightarrow S(y)$    $w_4$

R:

| x | y | w |
|----|----|-------|
| a1 | b1 | $w_1$ |
| a2 | b1 | $w_2$ |
| a3 | b2 | $w_3$ |

S:

| y | w |
|----|-------|
| b1 | $v_1$ |
| b2 | $v_2$ |
| b3 | $v_3$ |

Replace with hard constraint:

$\Gamma \equiv \forall x \forall y$
$A(x,y) \Leftrightarrow (R(x,y) \Rightarrow S(y))$

New relation A:

| x | y | w |
|----|----|-------|
| a1 | b1 | $w_4$ |
| a1 | b2 | $w_4$ |
| a1 | b3 | $w_4$ |
| a1 | b1 | $w_4$ |
| … | | |

Weight(

R:

| x | y |
|----|----|
| a1 | b1 |
| a3 | b2 |

S:

| y |
|----|
| b1 |

A:

| x | y |
|----|----|
| a1 | b1 |
| a1 | b2 |
| a2 | b1 |
| a3 | b1 |

) $= w_1 w_3 v_1 w_4 w_4 w_4 w_4 w_4$

**Theorem**: $P_{MLN}(Q) = P_D(Q \mid \Gamma)$

# Improved Translation

Soft constraint:

$$R(x,y) \Rightarrow S(y) \qquad w_4$$

Replace with hard constraint:

$$\Gamma \equiv \forall x \forall y$$
$$A(x,y) \Rightarrow (R(x,y) \Rightarrow S(x))$$

Replace $\Leftrightarrow$ with $\Rightarrow$

A clause remains a clause!

New weight is $w_4$-1

Probability may be < 0 !!! That's OK

**Theorem**:  $P_{MLN}(Q) = P_D(Q \mid \Gamma)$

New relation A:

| x | y | w |
|---|---|---|
| a1 | b1 | $w_4$-1 |
| a1 | b2 | $w_4$-1 |
| a1 | b3 | $w_4$-1 |
| a1 | b1 | $w_4$-1 |
| ... | | |

[VLDB'2016]

# SlimShot = SafePlans + Sample



Maximum Relative Error: Domain Size 100, 10200 Random Variables

Legend (left chart):
- Alchemy (0.09)
- Tuffy - Original (2.29)
- Tuffy - AI2 (0.18)
- CondSample (0.08)
- Cond bound for $\delta = 0.10$ (1.41E+04 samples)
- ImportanceSample (0.03)
- IS bound for $\delta = 0.10$ (2612 samples)

Legend (right chart):
- Alchemy (0.01)
- Tuffy - Original (0.04)
- Tuffy - AI2 (0.03)
- CondSample (0.00)
- Cond bound for $\delta = 0.10$ (3265 samples)
- ImportanceSample (0.00)
- IS bound for $\delta = 0.10$ (284 samples)

Axis labels: Maximum Relative Error (log) vs Number of Samples (log)

SlimShot

SlimShot needs N≈200 for Accuracy=10%

Accuracy = f(Number of Samples)
Lower is better

VLDB 2016
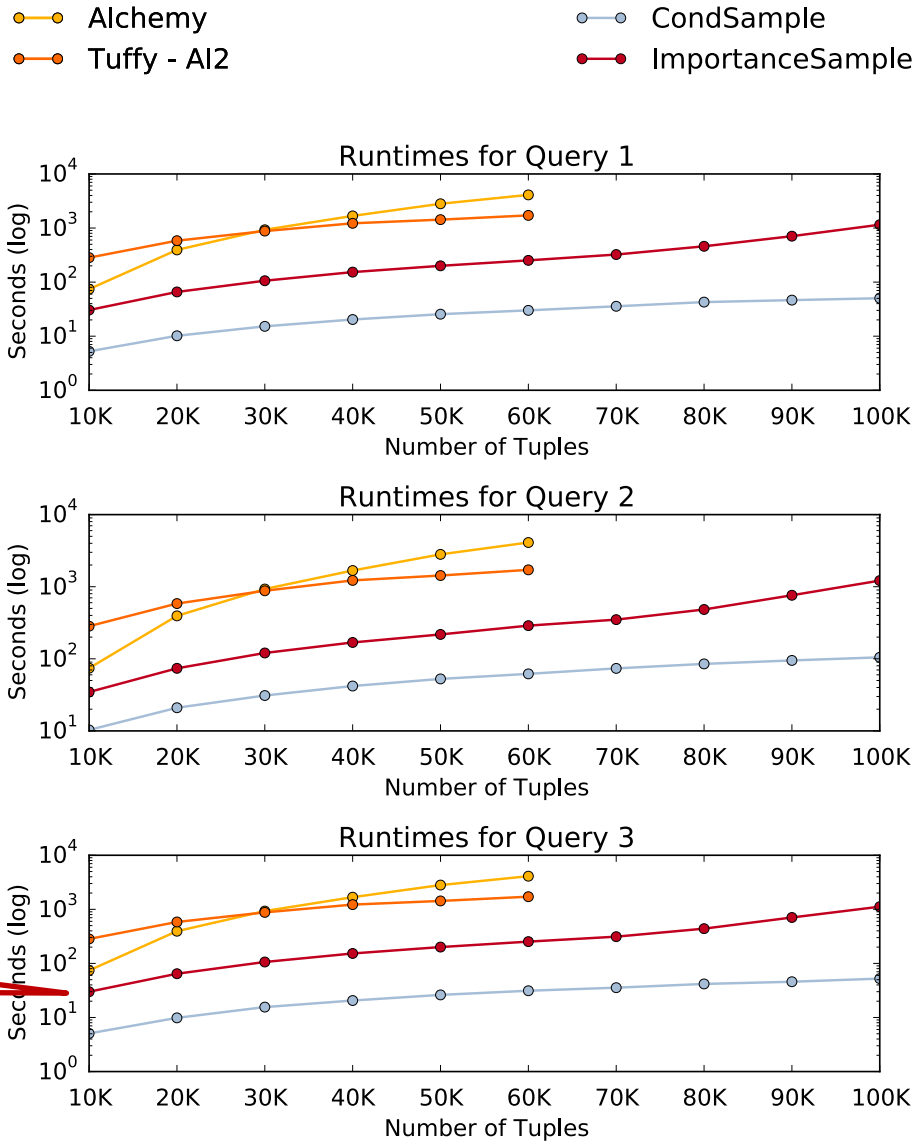
# Runtime

Runtime = f(N),  where N=10000

Runtime = f(precision)

# Scalability

# Duality

- The _dual_ of a query Q is the formula obtained by the following transformations:
  - $\wedge / \vee \; \rightarrow \; \vee / \wedge \qquad \forall / \exists \; \rightarrow \; \exists / \forall$
- Q and its dual have the same complexity

Query:

$H_0() = R(x), S(x,y), T(y)$

$\exists x \exists y (R(x) \wedge S(x,y) \wedge T(y))$

Dual query:

$H_0 = \forall x \forall y (R(x) \vee S(x,y) \vee T(y))$

$\top \sqsubseteq \forall S.R$

# Example: Liftable Clause

$$Q = \forall x \forall y \; S(x,y) \Rightarrow R(y) \qquad = \forall y \; (\exists x \; S(x,y) \Rightarrow R(y))$$

$$P(Q) = \Pi_{b \in Domain} \; P(\; \exists x \; S(x,b) \Rightarrow R(b))$$

Indep. $\forall$

$$P(Q) = \Pi_{b \in Domain} \; [1 - P(\exists x \; S(x,b)) \times (1 - P(R(b)))]$$

Indep. or:
$P(X \Rightarrow Y) =$
$= P(\neg X \vee Y)$
$= P(X)\,(1 - P(Y))$

$$P(Q) = \Pi_{b \in Domain} \; [1 - (1 - \Pi_{a \in Domain} \; (1 - P(S(a,b)))) \times (1 - P(R(b)))]$$

Indep. $\exists$

Lookup the probabilities in D

Runtime = $O(n^2)$.