A Dichotomy

for Non-Repeating Queries with Negation

in Probabilistic Databases



### Robert Fink and Dan Olteanu



PODS June 24, 2014

# Outline



#### The Dichotomy

#### The Interesting but Hard Queries

The Easy Queries

Leftovers

# Problem Setting

### Relational algebra query language fragment $1 \mbox{RA}^-$

- Included: Equi-joins, selections, projections, difference
- Excluded: Repeating relation symbols (self-joins), unions

#### Tuple-independent probabilistic model

- Each tuple associated with a fresh Boolean random variable x.
- P(x) is the probability that the tuple exists in the database.
- Simplest probabilistic model in the literature.
   Beyond this model, query tractability is quickly lost.
- Used by real-world large-scale probabilistic repositories, e.g., Google Knowledge Vault.

**Query Evaluation Problem:** For a fixed  $1RA^-$  query Q: Given a tuple-independent probabilistic database D and a tuple  $t \in Q(D)$ , compute its marginal probability.

# The Main Result

Data complexity of any  $1RA^-$  query Q on tuple-independent databases:

- Polynomial time if Q is **hierarchical** and
- #P-hard otherwise.

## The Main Result

Data complexity of any  $1RA^-$  query Q on tuple-independent databases:

- Polynomial time if Q is hierarchical and
- #P-hard otherwise.

This result strictly extends a 2004 result by Dalvi and Suciu:

- We added the relational algebra difference operator
  - and moved from conjunctive queries without self-joins to 1RA.
- Same syntactic characterization of tractable queries.
  - The hierarchical property can be recognized in LOGSPACE.
- The reason for tractability is however *different*.

# Hierarchical 1RA<sup>-</sup> Queries

Let [C] be the equivalence class of attribute C in query Q as defined by the transitivity of equi-join conditions and difference operators.

• E.g., C and D are in the same class due to join  $X(C) \bowtie_{C=D} Y(D)$  or difference  $X(C) -_{C \leftrightarrow D} Y(D)$  under attribute mapping  $C \leftrightarrow D$ .

## Hierarchical 1RA<sup>-</sup> Queries

Let [C] be the equivalence class of attribute C in query Q as defined by the transitivity of equi-join conditions and difference operators.

• E.g., C and D are in the same class due to join  $X(C) \bowtie_{C=D} Y(D)$  or difference  $X(C) -_{C \leftrightarrow D} Y(D)$  under attribute mapping  $C \leftrightarrow D$ .

(Boolean\*)  $1RA^{-}$  query Q is hierarchical if

For every pair of distinct attribute equivalence classes [A] and [B], there is no triple of relation symbols R, S, and T in Q such that

•  $R^{[A][\neg B]}$  has attributes in [A] and not in [B],

•  $T^{[\neg A][B]}$  has attributes in [B] and not in [A].

\* For non-Boolean queries, we need not check for equivalence classes with attributes in the query result.

# Examples

Examples of hierarchical queries:

$$\pi_{\emptyset} [(R(A) \bowtie S(A, B)) - T(A, B)]$$

$$\pi_{\emptyset} [(R(A) \times T(B)) - (U(A) \times V(B))]$$

$$\pi_{\emptyset} [(M(A) \times N(B)) - [(R(A) \times T(B)) - (U(A) \times V(B))]]$$

$$\pi_{\emptyset} [(M(A) \times N(B)) - \pi_{A} [(R(A) \times T(B)) - (U(A) \times V(B))]]$$

## Examples

Examples of hierarchical queries:

$$\pi_{\emptyset} [(R(A) \bowtie S(A, B)) - T(A, B)]$$

$$\pi_{\emptyset} [(R(A) \times T(B)) - (U(A) \times V(B))]$$

$$\pi_{\emptyset} [(M(A) \times N(B)) - [(R(A) \times T(B)) - (U(A) \times V(B))]]$$

$$\pi_{\emptyset} [(M(A) \times N(B)) - \pi_{A} [(R(A) \times T(B)) - (U(A) \times V(B))]]$$

Examples of non-hierarchical queries:

$$\pi_{\emptyset} \left[ R(A) \bowtie S(A, B) \bowtie T(B) \right]$$

$$\pi_{\emptyset} \left[ \pi_{B} \left( R(A) \bowtie S(A, B) \right) - T(B) \right]$$

$$\pi_{\emptyset} \left[ T(B) - \pi_{B} \left( R(A) \bowtie S(A, B) \right) \right]$$

$$\pi_{\emptyset} \left[ X(A) \bowtie \left[ R(A) - \pi_{A} \left( T(B) \bowtie S(A, B) \right) \right] \right]$$

# Outline



The Dichotomy

#### The Interesting but Hard Queries

The Easy Queries

Leftovers

## Hardness Proof Idea

Reduction from #P-hard model counting problem for positive 2DNF:

- Given a non-hierarchical 1RA query Q and
- A positive bipartite DNF formula Ψ,
- Construct a tuple-independent database D with
  - size polynomial in the number of variables and clauses in  $\Psi$ , and
  - tuples annotated with variables in  $\Psi$  such that  $\Psi$  annotates Q(D).
- Then  $\#\Psi = 2^n \cdot P_{Q(D)}$ , where
  - $P_{Q(D)}$  is the probability of Q(D),
  - ▶ 1/2 is the probability of each variable in  $\Psi$ , and
  - *n* is the number of variables in  $\Psi$ .

### Example of Hardness Reduction

Input formula and query:

$$\Psi = x_1 y_1 \lor x_1 y_2,$$
  
$$Q = \pi_{\emptyset} \Big[ R(A) - \pi_A \big( T(B) \bowtie S(A, B) \big) \Big]$$

Construct database such that  $\Psi$  annotates Q's (nullary) result:

- Column Φ holds annotations over variables in Ψ.
  - Special annotations:  $\top$  (true),  $\perp$  (false)
- Variables used as constants for the attribute B in T and S.
- **S** $(a, b, \phi)$ : Clause a has variable b exactly when  $\phi$  is true.
- $R(a, \top)$  and  $T(b, \neg b)$ : *a* is a clause and *b* is a variable in  $\Psi$ .

R	Т	S	$T \bowtie S$	$\pi_A(T \bowtie S)$	$R - \pi_A(T \bowtie S)$
AΦ	ΒΦ	ΑΒΦ	ΑΒ Φ	ΑΦ	ΑΦ
1 ⊤	$x_1 \neg \mathbf{x_1}$	$1 x_1 \top$	$1 x_1 \neg \mathbf{x_1}$	$1 \ \neg x_1 \lor \neg y_1$	1 x <sub>1</sub> y <sub>1</sub>
2 ⊤	<i>y</i> <sub>1</sub> ¬ <b>y</b> <sub>1</sub>	$1  y_1 \top$	$1 y_1 \neg \mathbf{y_1}$	$2 \ \neg \textbf{x_1} \lor \neg \textbf{y_2}$	2 x <sub>1</sub> y <sub>2</sub>
	<i>y</i> <sub>2</sub> ¬ <b>y</b> <sub>2</sub>	$1 y_2 \perp$	$1 y_2 \perp$		
		$2 x_1 \top$	2 <i>x</i> <sub>1</sub> ¬ <b>x</b> <sub>1</sub>		
		$2 y_1 \perp$	$2 y_1 \perp$		
		$2 y_2 \top$	2 <i>y</i> <sub>2</sub> ¬ <b>y</b> <sub>2</sub>		

### Example of Hardness Reduction

Input formula and query:

$$\Psi = x_1 y_1 \lor x_1 y_2,$$
  
$$Q = \pi_{\emptyset} \Big[ R(A) - \pi_A \big( T(B) \bowtie S(A, B) \big) \Big]$$

Construct database such that  $\Psi$  annotates Q's (nullary) result:

- Column Φ holds annotations over variables in Ψ.
  - Special annotations:  $\top$  (true),  $\perp$  (false)
- Variables used as constants for the attribute B in T and S.
- **S** $(a, b, \phi)$ : Clause a has variable b exactly when  $\phi$  is true.
- $R(a, \top)$  and  $T(b, \neg b)$ : *a* is a clause and *b* is a variable in  $\Psi$ .

R	<u> </u>	S	$T \bowtie S$	$\pi_A(T \bowtie S)$	R-c	$\pi_A(T \bowtie S)$
ΑΦ	ΒΦ	ΑΒΦ	ΑΒΦ	ΑΦ	Α	Φ
1 ⊤	$x_1 \neg x_1$	$1 x_1 \top$	1 <i>x</i> <sub>1</sub> ¬ <b>x</b> <sub>1</sub>	$1 \ \neg x_1 \lor \neg y_1$	1	$x_1y_1$
2 ⊤	<i>y</i> <sub>1</sub> ¬ <b>y</b> <sub>1</sub>	$1 y_1  op$	$1 y_1 \neg \mathbf{y_1}$	$2 \ \neg \textbf{x_1} \lor \neg \textbf{y_2}$	2	x <sub>1</sub> y <sub>2</sub>
	<i>y</i> <sub>2</sub> ¬ <b>y</b> <sub>2</sub>	$1 y_2 \perp$	$1 y_2 \perp$			
		$2 x_1 \top$	$2 x_1 \neg \mathbf{x_1}$			
		$2 y_1 \perp$	$2 y_1 \perp$			
		2 <i>y</i> <sub>2</sub> ⊤	2 <i>y</i> <sub>2</sub> ¬ <b>y</b> <sub>2</sub>			

Query Q is already hard when T is the only uncertain input relation!

## Hard Query Patterns

There are 48 (!) minimal non-hierarchical query patterns.

- Binary trees with leaves A, AB, and B and inner nodes  $\bowtie$  or -.
  - Some are symmetric and need not be consider separately: A and B can be exchanged, joins are commutative and associative.
  - > Still, many cases left to consider due to the difference operator.



- There is a database construction scheme for each pattern.
- Each non-hierarchical query Q matches a pattern P<sub>x.y</sub>.

## Hard Query Patterns

There are 48 (!) minimal non-hierarchical query patterns.

- Binary trees with leaves A, AB, and B and inner nodes  $\bowtie$  or -.
  - Some are symmetric and need not be consider separately: A and B can be exchanged, joins are commutative and associative.
  - > Still, many cases left to consider due to the difference operator.



- There is a database construction scheme for each pattern.
- Each non-hierarchical query Q matches a pattern  $P_{x.y.}$

 $P_{1.1}$  is the only hard pattern to consider w/o the difference operator!

## Non-hierarchical Queries Match Minimal Hard Patterns

Each non-hierarchical query Q matches a pattern  $P_{x.y}$ :

- There is a total mapping from  $P_{x,y}$  to Q's parse tree that
  - ▶ is identity on inner nodes  $\bowtie$  and -,
  - preserves ancestor-descendant relationships,
  - maps leaves A, AB, B to relations  $R^{[A][\neg B]}, S^{[A][B]}, T^{[\neg A][B]}$ .



• The match preserves the annotation of the query pattern: Q and  $P_{x,y}$  have the same annotation for any input database.

# Outline



The Dichotomy

#### The Interesting but Hard Queries

The Easy Queries

Leftovers

## Evaluation of Hierarchical 1RA<sup>-</sup> Queries

Approach based on knowledge compilation

- For any database D, the probability  $P_{Q(D)}$  of a 1RA<sup>-</sup> query Q is the probability  $P_{\Psi}$  of the query annotation  $\Psi$ .
- Compile  $\Psi$  into poly-size OBDD( $\Psi$ ).
- Compute probability of  $OBDD(\Psi)$  in time linear in its size.

## Evaluation of Hierarchical 1RA<sup>-</sup> Queries

Approach based on knowledge compilation

- For any database D, the probability P<sub>Q(D)</sub> of a 1RA<sup>-</sup> query Q is the probability P<sub>Ψ</sub> of the query annotation Ψ.
- Compile  $\Psi$  into poly-size OBDD( $\Psi$ ).
- Compute probability of  $OBDD(\Psi)$  in time linear in its size.

Distinction from existing tractability results [O. & Huang 2008]:

- 1RA<sup>-</sup> queries w/o difference: Annotations are read-once.
  - Read-once annotations admit linear-size OBBDs.
- 1RA<sup>-</sup> queries: Annotations are <u>not</u> read-once.
  - They admit OBBDs of size linear in the database size <u>but</u> exponential in the query size.

# The Inner Workings

From hierarchical 1RA<sup>-</sup> to RC-hierarchical  $\exists$ -consistent RC<sup> $\exists$ </sup>:

- Translate query Q into an equivalent disjunction of disjunction-free existential relational calculus queries  $Q_1 \lor \cdots \lor Q_k$ .
  - k can be very large for queries with projection under difference!

#### RC-hierarchical:

- For each  $\exists_X(Q')$ , every relation symbol in Q' has variable X.
  - Each of the disjuncts gives rise to a poly-size OBDD.

#### ∃-consistent:

The nesting order of the quantifiers is the same in  $Q_1, \cdots, Q_k$ .

- All OBDDs have compatible variable orders and their disjunction is a poly-size OBDD.
- The OBDD width grows exponentially with k, its height stays linear in the size of the database.
  - Width = maximum number of edges crossing the section between any two consecutive levels.

## Query Evaluation Example

Consider the following query and tuple-independent database:

$$Q = \pi_{\emptyset} \Big[ \big( R(A) \times T(B) \big) - \big( U(A) \times V(B) \big) \Big]$$

R	T	U	V	$R \bowtie T$	$R \bowtie T - U \bowtie V$
AΦ	ВΦ	ΑΦ	ВΦ	ΑΒ Φ	ΑΒ Φ
1 r <sub>1</sub> 2 r <sub>2</sub>	1 t <sub>1</sub> 2 t <sub>2</sub>	1 u <sub>1</sub> 2 u <sub>2</sub>	1 v <sub>1</sub> 2 v <sub>2</sub>	$   \begin{array}{r}     1 & 1 & r_1 t_1 \\     1 & 2 & r_1 t_2 \\     2 & 1 & r_2 t_1 \\     2 & 2 & r_2 t_2   \end{array} $	$ \begin{array}{r} 1 \ 1 \ r_{1}t_{1}\neg(u_{1}v_{1}) \\ 1 \ 2 \ r_{1}t_{2}\neg(u_{1}v_{2}) \\ 2 \ 1 \ r_{2}t_{1}\neg(u_{2}v_{1}) \\ 2 \ 2 \ r_{2}t_{2}\neg(u_{2}v_{2}) \end{array} $

## Query Evaluation Example

Consider the following query and tuple-independent database:

$$Q = \pi_{\emptyset} \Big[ \big( R(A) \times T(B) \big) - \big( U(A) \times V(B) \big) \Big]$$

R	<u> </u>	<u> </u>	V	$R \bowtie T$	$R \bowtie T - U \bowtie V$
ΑΦ	ВΦ	ΑΦ	ВΦ	ΑΒΦ	ΑΒ Φ
1 r <sub>1</sub> 2 r <sub>2</sub>	1 t <sub>1</sub> 2 t <sub>2</sub>	1 u <sub>1</sub> 2 u <sub>2</sub>	1 v <sub>1</sub> 2 v <sub>2</sub>	$ \begin{array}{c} 1 \ 1 \ r_{1}t_{1} \\ 1 \ 2 \ r_{1}t_{2} \\ 2 \ 1 \ r_{2}t_{1} \\ 2 \ 2 \ r_{2}t_{1} \end{array} $	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$

The annotation of Q is:

$$\Psi = r_1 \big[ t_1 (\neg u_1 \lor \neg v_1) \lor t_2 (\neg u_1 \lor \neg v_2) \big] \lor r_2 \big[ t_1 (\neg u_2 \lor \neg v_1) \lor t_2 (\neg u_2 \lor \neg v_2) \big].$$

- Variables entangle in  $\Psi$  beyond read-once factorization.
- This is the pivotal intricacy introduced by the difference operator.

Query Evaluation Example (2)

Translate 
$$Q = \pi_{\emptyset} \Big[ (R(A) \times T(B)) - (U(A) \times V(B)) \Big]$$
 into  $\mathsf{RC}^{\exists}$ :

$$Q_{RC} = \underbrace{\exists_A (R(A) \land \neg U(A)) \land \exists_B T(B)}_{Q_1} \lor \underbrace{\exists_A R(A) \land \exists_B (T(B) \land \neg V(B))}_{Q_2}.$$

Both  $Q_1$  and  $Q_2$  are RC-hierarchical.

•  $Q_1 \lor Q_2$  is  $\exists$ -consistent: Same order  $\exists_A \exists_B$  for  $Q_1$  and  $Q_2$ .

Query annotation:

$$\Psi = \underbrace{(r_1 \neg u_1 \lor r_2 \neg u_2) \land (t_1 \lor t_2)}_{\Psi_1} \lor \underbrace{(r_1 \lor r_2) \land (t_1 \neg v_1 \lor t_2 \neg v_2)}_{\Psi_2}.$$

Both  $\Psi_1$  and  $\Psi_2$  admit linear-size OBDDs.

The two OBDDs have compatible orders and their disjunction is an OBDD whose width is the product of the widths of the two OBDDs.

## Query Evaluation Example (3)

Compile query annotation into OBDD:



# Outline



The Dichotomy

#### The Interesting but Hard Queries

The Easy Queries

### Leftovers

# Dichotomies Beyond 1RA-

Some known dichotomies

- Conjunctive queries w/o self-joins
- Unions of conjunctive queries
- Quantified relational algebra queries

[Dalvi & Suciu 2004] [Dalvi & Suciu 2010] [F. & O. & Rath 2011]

#### Full relational algebra

It is undecidable whether the union of two equivalent relational algebra queries, one hard and one tractable, is tractable.

Non-repeating relational algebra =  $1RA^-$  + union.

- Hierarchical property not enough.
- $\pi_{\emptyset}[(R(A) \bowtie S_1(A, B) \cup T(B) \bowtie S_2(A, B)) S(A, B)]$  is hard, though it is equivalent to a union of two hierarchical 1RA<sup>-</sup> queries.

Non-repeating relational calculus

- $S(x,y) \land \neg R(x)$  is tractable,  $S(x,y) \land (R(x) \lor T(y))$  is hard.
  - Both are non-repeatable, yet not expressible in 1RA<sup>-</sup>.
- Possible (though expensive) approach:
  - ► Translate to RC<sup>∃</sup> and check RC-hierarchical and ∃-consistency.

# Thank you!