## **Probabilistic Databases and Reasoning**

Thomas Lukasiewicz and Dan Olteanu

University of Oxford

#### Probabilistic Databases and Reasoning

This 3-hour tutorial has two main parts:

- 1. Dan Olteanu: *Probabilistic Databases* Now: 8.30am - 10am.
- Thomas Lukasiewicz: Probabilistic Reasoning Next: 10am - 10.30am, then a break, then 11am - 12pm.

Further 1-hour lectures on advanced topics in probabilistic databases:

- 1. DL invited talk today at 12.10pm Dan Suciu: *Lifted Inference in Probabilistic Database*
- KR invited lecture tomorrow at 9.30am
   Dan Suciu: Query compilation: the View from the Database Side

KR features several more papers on probabilistic data and knowledge bases!

#### Probabilistic Databases

For the purpose of the first half of this tutorial:

Probabilistic data =

Relational data

 $^+$ 

Probabilities that measure the degree of uncertainty in the data.

Long-term key challenges:

Models for probabilistic data to capture data and its uncertainty.

Query evaluation = Probabilistic inference
 Query answers are annotated with output probabilities.

# Outline



Dan Suciu Dan Olteanu Christopher Ré Christoph Koch

STATUENS LECTURES ON DATA MANAGEMENT

#### Why Probabilistic Databases?

Probabilistic Data Models

The Query Evaluation Problem

Dichotomies for Query Evaluation

- The Hard Queries
- The Tractable Queries

Ranking Queries

Next Steps

References

#### Research Development Map

We can unify logic and probability by defining distributions over possible worlds that are first-order model structures (objects and relations). Gaifman'64

Early work (80s and 90s):

Basic data models and query processing

Wong'82, Shoshani'82, Cavallo & Pittarelli'87, Barbará'92, Lakshmanan'97,'01, Fuhr& Röllke'97, Zimányi'97, ..

Recent wave (2004 - now):

Computational complexity of query evaluation

Probabilistic database systems

Stanford (Trio), UW (MystiQ), Cornell & Oxford (MayBMS/SPROUT), IBM Almaden & Rice (MCDB), LogicBlox & Technion & Oxford (PPDL), Florida, Maryland, Purdue, Waterloo, Wisconsin, ..

### Why This Interest in Probabilistic Databases?

Probabilistic relational data is commonplace. It accommodates several possible interpretations of the data weighted by probabilities.

 Information extraction: Probabilistic data inferred from unstructured data (e.g., web) text using statistical models
 Google Knowledge Vault, DeepDive, NELL

 Manually entered data Represent several possible readings with MayBMS [Antova'07] Infer missing data with meta-rule semi-lattices [Stoyanovich'11] Manage OCR data with Staccato/Google OCRopus [Kumar'11]

#### Data cleaning

Represent several possible data repairs

## Data integration Google Squared and SPROUT<sup>2</sup>

Risk management (Decision support queries, hypothetical queries); ...

[Beskales'09]

[Fink'11]

#### Information Extraction

Possible segmentations of unstructured text

[Sarawagi'06]

<u>ID</u>	HouseNo	Area	City	PinCode	Р
1	52	Goregaon West	Mumbai	400 062	0.1
1	52-A	Goregaon	West Mumbai	400 062	0.2
1	52-A	Goregaon West	Mumbai	400 062	0.4
1	52	Goregaon	West Mumbai	400 062	0.2

52-A Goregaon West Mumbai 400 076

- Probabilities obtained using probabilistic extraction models (e.g., CRF) The probabilities correlate with the precision of the extraction.
- The output is a ranked list of possible extractions
- Several segmentations are required to cover most of the probability mass and improve recall

Avoid empty answer to queries such as Find areas in 'West Mumbai'

#### Continuously-Improving Information Extraction

#### Recently-Learned Facts twitter

[Mitchell'15]

Refresh

instance	iteration	date learned	confidence
biscutate_swift is an animal	211	18-feb-2011	100.0 🏖 ኛ
<u>pedigree_animals</u> is a <u>mammal</u>	210	17-feb-2011	99.5 🏠 🖑
poppy seed holiday bread is a baked good	212	20-feb-2011	100.0 🍃 🖑
<u>manuel criado de val</u> is a <u>South American person</u>	210	17-feb-2011	99.5 🕼 ኛ
<u>dillon_county_airport</u> is an <u>airport</u>	210	17-feb-2011	93.8 🕼 ኛ
the sports team toronto blue jays was the winner of n1993 world series	212	20-feb-2011	96.9 🗳 🖏
mozart is a person who died at the age of 35	210	17-feb-2011	96.9 🏠 ኛ
<u>peoria</u> and <u>arizona</u> are <u>proxies</u> for eachother	210	17-feb-2011	99.9 🍃 ኛ
wutv tv is a <u>TV affiliate of</u> the network <u>fox</u>	210	17-feb-2011	96.9 🕼 🖑
white stripes collaborates with jack white	210	17-feb-2011	93.8 🖉 🕏

# Manu<sup>e</sup>al?y-enter d census data

MayBMS manages  $10^{10^6}$  possible readings of census data



We want to enter the information from forms like these into a database.

- What is the marital status of the first resp. the second person?
- What are the social security numbers? 185? 186? 785?

[Antova'07]

# Manu<sup>e</sup> I?y-enter d census data

Social Securit	y Number: Name:	281 m2_	; th	
Mari	tal Status:	(1) single	(2) married (4) widowed (2)	
Social Securit Mari	Social Security Number: Name: BLOWA Marital Status: (1) single (2) married (3) divorced (4) widowed (5)			
(TID)	SSN	N	М	
$t_1$	NULL	. Smith	NULL	
$t_2$	NULL	Brown	NULL	

Much of the available information cannot be represented and is lost, e.g.

- Smith's SSN is either 185 or 785; Brown's SSN is either 185 or 186.
- Data cleaning: No two distinct persons can have the same SSN.

### OCR on manually-entered data

#### Staccato

[Kumar'12]



- Stochastic automaton constructed from text using Google OCRopus.
- String *F0 rd* has the highest probability (0.21).
- String *Ford* has lower probability (0.12).

Staccato accommodates several possible readings of the text to increase recall.

#### Web Data Integration with Google Squared

- Tables instead of page links as answers to Google queries [Fink'11]
- Integration of data sources with contradicting information or different schemas, degrees of trust, and degrees of completion
- Confidence values mapped to [0,1]

G	008	e squared	comedy movies		Square it Add	
com	comedy movies					
	Item Nar	ne 💌	Language	Director	Release Date	
×	The Mas	k	English	Chuck Russell	29 July 1994	
×	Scary M	English     Ianguage for the     www.infibeam.ce	e mask om - all 9 sources »	Chuck Russell directed by for The www.infibeam.com	• Mask - all 9 sources »	
		Other possible values		Other possible values	-	
×	Superba	English Langu language for Ma www.freebase.c	a <b>ge</b> Low confidence ask com	John R. Dilworth director for The Ma www.freebase.com	Low confidence ask	
×	Music	<ul> <li>english, french languages for t www.dvdreview.</li> </ul>	1 Low confidence he mask .com	Fiorella Infascelli directed by for The www.freebase.com	Low confidence Mask - all 2 sources »	
×	Knocked	Italian Langua language for Th www.freebase.cd	ge Low confidence ne Mask com	Charles Russell directed by for The www.freebase.com	Low confidence Mask - all 2 sources »	
		Search for more val	ues »	Search for more values	<u>. »</u>	

# Outline



#### Why Probabilistic Databases?

#### Probabilistic Data Models

The Query Evaluation Problem

Dichotomies for Query Evaluation

- The Hard Queries
- The Tractable Queries

Ranking Queries

Next Steps

References

## Revisiting the Census Data Example

Social Secu	rity Number: Name:	28 C m 2	5 .vth
м	arital Status:	(1) single (3) divorced (	<ul> <li>【 (2) married</li> <li>【 (4) widowed</li> </ul>
Social Sect	urity Number: Name:	18 Bro	5 Wh
м	arital Status:	(1) single (3) divorced (1)	(2) married (2) (4) widowed (2)
RID	SSN	Ν	М
$t_1$	NULL	Smith	NULL
$t_2$	NULL	Brown	NULL

NULL values are too uninformative.

We could instead incorporate all available possibilities:

- Smith's SSN is either 185 or 785; Brown's SSN is either 185 or 186.
- Smith's M is either 1 or 2; Brown's M is either 1, 2, 3, or 4.

#### Revisiting the Census Data Example

There are  $2 \times 2 \times 2 \times 4 = 32$  possible readings of our two census entries.



. . .

#### Incomplete Databases

An Incomplete Database is a finite set of database instances  $\mathbf{W} = (W_1, \dots, W_n)$ .

	$W_1$	
SSN	Ν	М
185	Smith	1
185	Brown	1

	$W_2$	
SSN	Ν	М
185	Smith	1
185	Brown	2

Each  $W_i$  is a *possible world*.

	$W_3$	
SSN	Ν	М
185	Smith	1
185	Brown	3

	$W_4$	
SSN	Ν	М
185	Smith	1
185	Brown	4

	$W_5$	
SSN	Ν	М
185	Smith	1
186	Brown	1

	$W_6$	
SSN	Ν	М
185	Smith	1
186	Brown	2

. . .

• •

### Incomplete Databases

An Incomplete Database is a finite set of database instances  $\mathbf{W} = (W_1, \dots, W_n)$ .



 $\rightarrow$  Key challenge: How to succinctly represent incomplete databases?

#### Probabilistic Databases

A Probabilistic Database is  $(\mathbf{W}, P)$ , where **W** is an incomplete database and  $P: \mathbf{W} \to [0,1]$  is a probability distribution:  $\sum_{W_i \in \mathbf{W}} P(W_i) = 1$ .

М

$W_1: P(W_1) = 0.1$			
SSN	Ν	М	
185	Smith	1	
185	Brown	1	

$W_3: P(W_3) = 0.1$		
SSN	Ν	М
185	Smith	1
185	Brown	3

W4 :	$P(W_4) =$	0.1
SSN	Ν	М
185	Smith	1
185	Brown	4

 $W_2: P(W_2) = 0.1$ Ν SSN185Smith185Rrown2

SSN

For <b>W</b> = { $W_1,$	., W <sub>6</sub> },
$\sum_{W_i \in \mathbf{W}} P(W_i) =$	1.

$W_5: P(W_5) = 0.3$		
SSN	Ν	М
185	Smith	1
186	Brown	1

$W_6: P(W_6) = 0.3$		
SSN	Ν	М
185	Smith	1
186	Brown	2

## Succinct Representations of Incomplete/Probabilistic Data

Social Security Number: Name:	<u> 185</u> Smith
Marital Status:	(1) single
Social Security Number:	185 BEDLIG
Name:	

Succinct or-set representation:

[Imielinski'91]

SSN	Ν	М
{ 185,785 }	Smith	{ 1,2 }
$\{ 185, 186 \}$	Brown	{ 1,2,3,4 }

It exploits independence of possible values for different fields:

- Choice for Smith's SSN independent of choice of for Brown's SSN.
- Likewise, the probability distributions associated with these choices are independent (not shown).

### BID: Alternative Representation of Our Or-Set

<u>RID</u>	SSN	Ρ
$t_1$	185	0.7
$t_1$	785	0.3
t <sub>2</sub>	185	0.8
$t_2$	186	0.2

<u>RID</u>	N	Ρ
$t_1$	Smith	1
$t_2$	Brown	1

RID	Μ	Р
$t_1$	1	0.9
$t_1$	2	0.1
t <sub>2</sub>	1	0.25
$t_2$	2	0.25
$t_2$	3	0.25
$t_2$	4	0.25

## BID: Alternative Representation of Our Or-Set

RID	SSN	Р
$t_1$	185	0.7
$t_1$	785	0.3
$t_2$	185	0.8
$t_2$	186	0.2

<u>RID</u>	Ν	Ρ
$t_1$	Smith	1
$t_2$	Brown	1

<u>RID</u>	М	Р
$t_1$	1	0.9
$t_1$	2	0.1
$t_2$	1	0.25
$t_2$	2	0.25
$t_2$	3	0.25
$t_2$	4	0.25

Interpretation:

The tuples within each block with the same key RID are *disjoint* Each world contains one tuple per block, so the tuples within a block are mutually exclusive.

## BID: Alternative Representation of Our Or-Set

<u>RID</u>	SSN	Ρ
$t_1$	185	0.7
$t_1$	785	0.3
t <sub>2</sub>	185	0.8
$t_2$	186	0.2

<u>RID</u>	Ν	Ρ
$t_1$	Smith	1
$t_2$	Brown	1

<u>RID</u>	М	Р
$t_1$	1	0.9
$t_1$	2	0.1
$t_2$	1	0.25
$t_2$	2	0.25
$t_2$	3	0.25
$t_2$	4	0.25

Interpretation:

- The tuples within each block with the same key RID are *disjoint* Each world contains one tuple per block, so the tuples within a block are mutually exclusive.
- Blocks are *independent* of each other.

The choices of tuples within different blocks are independent. The aggregated probability of the worlds taking the first tuple of the first block in each relation is  $0.7 \times 1 \times 0.9 = 0.63$ .

These *block-independent disjoint* (BID) relations are sometimes called x-relations or x-tables. Google squares are prime examples.

### More on BID Databases

BIDs also allow blocks with probabilities less than 1:

<u>RID</u>	SSN	Р
$t_1$	185	0.6
$t_1$	785	0.3
t <sub>2</sub>	185	0.8
t <sub>2</sub>	186	0.2

<u>RID</u>	Ν	Ρ
$t_1$	Smith	0.9
t <sub>2</sub>	Brown	1

<u>RID</u>	М	Р
$t_1$	1	0.8
$t_1$	2	0.1
$t_2$	1	0.25
$t_2$	2	0.25
$t_2$	3	0.25
$t_2$	4	0.25

Interpretation:

There are worlds where the first block of each of the three relations is empty, e.g., the following world:

RID	SSN	Ρ	<u>RID</u>	N	Ρ	<u>RID</u>	М	Р
t <sub>2</sub>	186	0.2	t <sub>2</sub>	Brown	1	t <sub>2</sub>	4	0.25

The probability of this world is

 $0.2 \times 1 \times 0.25 \times (\boldsymbol{1-0.6-0.3}) \times (\boldsymbol{1-0.9}) \times (\boldsymbol{1-0.8-0.1}) = 5 \times 10^{-5}.$ 

Clarification notes to come with the previous slide and to answer questions posed during the tutorial:

\* The two BIDs from the previous two slides are not equivalent since they do not represent the same probabilistic database! Furthermore, by allowing groups with empty instances, some tuples are only partially defined in the column-oriented representation.

\* See [Antova'08] for column-oriented representation of relations with attribute-level uncertainty.

#### TI: Tuple-Independent Databases

*TI databases* are BID databases where each block has exactly one tuple.

TI databases are the simplest and most common probabilistic data model.

RID	SSN	Ρ		RID	N	Ρ	<u>RID</u>	М	P
$t_1$	185	0.7	1	$t_1$	Smith	1	$t_1$	1	0.9
t <sub>2</sub>	185	0.8		t <sub>2</sub>	Brown	1	t <sub>2</sub>	2	0.2

Interpretation:

- Each tuple t is in a random world with its probability p(t).
- A relation with n tuples, whose probabilities are less than 1, has 2<sup>n</sup> possible worlds, since each tuple may be in or out.
- Our TI example has 2<sup>4</sup> worlds: Any subset of the first and third relation and the entire second relation.

5

#### Are BID Databases Enough?

BIDs (and TIs) are good at capturing independence and local choice. What about correlations across blocks?

Enforce the key dependency on SSN in each world. That is: Discard the worlds where both t<sub>1</sub> and t<sub>2</sub> have SSN = 185.

<u>RID</u>	SSN	Ρ
$t_1$	185	0.6
$t_1$	785	0.3
$t_2$	185	0.8
t <sub>2</sub>	186	0.2

#### Are BID Databases Enough?

BIDs (and TIs) are good at capturing independence and local choice. What about correlations across blocks?

Enforce the key dependency on SSN in each world. That is: Discard the worlds where both t<sub>1</sub> and t<sub>2</sub> have SSN = 185.

<u>RID</u>	SSN	Р		<u>RID</u>	SSN	Φ
$t_1$	185	0.6		$t_1$	185	X = 1
$t_1$	785	0.3	$\Rightarrow$	$t_1$	785	X = 2
$t_2$	185	0.8		$t_2$	185	$Y = 1 \land X \neq 1$
$t_2$	186	0.2		$t_2$	186	Y = 2

This constraint is supported by a probabilistic version of *conditional databases*. [Imielinski'84]

Idea: Use random variables to encode correlations between tuples.

- Exclude the world where t<sub>1</sub> and t<sub>2</sub> have the same SSN 185 by using contradicting assignments for variable X.
- Transfer probabilities of tuples to probability distributions of variables.

#### PC: Probabilistic Conditional Databases

A *PC database* is  $(\mathbf{D}, \mathbf{X}, \Phi)$ , where **D** is a relational database, **X** is a set of independent random variables, and  $\Phi$  is a function mapping each tuple in **D** to a propositional formula over **X**.

RID	SSN	Φ	VAR	Dom	Р
$t_1$	185	X = 1	X	1	0.6
$t_1$	785	X = 2	X	2	0.3
$t_2$	185	$Y = 1 \land X \neq 1$	Y	1	0.8
t <sub>2</sub>	186	<i>Y</i> = 2	Y	2	0.2

Interpretation:

- The *world table* (right) lists the probability distribution for each independent random variable in **X**.
- Each total valuation of variables in X defines a world whose probability is the product of probabilities of the variable assignments.
- Each tuple t is conditional on the satisfiability of the formula Φ(t) and is contained in those worlds defined by valuations that satisfy Φ(t).

Clarification notes to come with the previous slide and to answer questions posed during the tutorial:

\* The PC table from the previous slide is not equivalent to the BID table from two slides ago: While the PC table captures the key dependency on SSN, the BID table does not.

\* However, the PC table is not the BID table where the key dependency is enforced: This is because we did not adjust the probabilities of the remaining worlds that satisfy the key dependency.

\* The mechanism for this adjustment is called conditioning, see [Koch'08].

## TIs and BIDs are Special Cases of PCs

Recall our previous TI database example:

RID	SSN	Ρ
$t_1$	185	0.7
t <sub>2</sub>	185	0.8

<u>RID</u>	Ν	Ρ
$t_1$	Smith	1
$t_2$	Brown	1

RID	М	Р
$t_1$	1	0.9
t <sub>2</sub>	2	0.25

Here is a PC encoding of the above TI database:

<u>RID</u>	SSN	Φ	Ρ	RID	N	Φ	Ρ	RID	М	Φ	Ρ
$t_1$	185	<b>s</b> 1	0.7	$t_1$	Smith	<b>n</b> 1	1	$t_1$	1	<b>m</b> 1	0.9
$t_2$	185	<b>s</b> 2	0.8	t <sub>2</sub>	Brown	<b>n</b> <sub>2</sub>	1	t <sub>2</sub>	2	<b>m</b> <sub>2</sub>	0.25

Idea:

- Consider a set of Boolean random variables
- Associate each tuple in the TI database with exactly one of them
- For instance,  $s_1$  annotates  $(t_1, 185)$  and  $P(s_1) = 0.7$
- World table with variable assignments may be stored explicitly

#### Takeaways

Various representations for probabilistic databases of increasing expressiveness.

Most complex: probabilistic conditioned databases. [Imielinski'84]

- ► Trio's ULDBs [Benjelloun'06] and MayBMS' U-relations [Antova'07].
- Completeness: They can represent any probabilistic database.
- Mid-level: block-independent disjoint databases.
   [Barbará'92]
  - MystiQ, Trio, MayBMS, SPROUT<sup>2</sup>.
  - Prime examples of BIDs: Google squares.
  - ► Not complete, but achieve completeness via conjunctive queries over BIDs.

[Poole'93]

Simplest: tuple-independent databases.

- ► *The* norm in real-world repositories like Google's, DeepDive, and NELL.
- Most theoretical work on complexity of query evaluation done for them.
- Not complete even via unions of conjunctive queries.
- However, inference in Markov Logic Networks is captured by relational queries on TI databases! See Dan Suciu's invited DL'16 talk. Also work by Guy van den Broeck. [Jha'12]

# Outline



Dan Suciu Dan Olteanu Christopher Ré Christoph Koch

STATUENS LECTURES ON DATA MANAGEMENT

Why Probabilistic Databases?

Probabilistic Data Models

#### The Query Evaluation Problem

Dichotomies for Query Evaluation

• The Hard Queries

• The Tractable Queries

Ranking Queries

Next Steps

References

#### Possible Worlds Semantics

The underlying semantics of query evaluation in probabilistic databases:

Possible worlds semantics: Given a database  $\mathbf{W} = \{W_1, \dots, W_n\}$  and a query Q, the query answer is  $Q(\mathbf{W}) = \{Q(W_1), \dots, Q(W_n)\}.$ 

#### Possible Worlds Semantics

The underlying semantics of query evaluation in probabilistic databases:

Possible worlds semantics: Given a database  $\mathbf{W} = \{W_1, \dots, W_n\}$  and a query Q, the query answer is  $Q(\mathbf{W}) = \{Q(W_1), \dots, Q(W_n)\}.$ 

Investigations so far followed three main directions:

- 1. Possible and certain query answers for incomplete databases.
- 2. Probabilities of query answers for probabilistic databases.
- 3. Succinct representation of  $Q(\mathbf{W})$  for query languages and data models.

Approaches 1 & 2 close the possible worlds semantics: They compute one relation with answer tuples and possibly their probabilities.

#### Queries on Incomplete Databases

Given query Q and incomplete database **W**:

- An answer t is certain, if  $\forall : W_i \in \mathbf{W}, t \in Q(W_i)$
- An answer t is possible if  $\exists W_i \in \mathbf{W}, t \in Q(W_i)$

	$W_1$	
SSN	Ν	М
185	Smith	1
185	Brown	1

	$W_2$	
SSN	Ν	М
185	Smith	1
185	Brown	2

	$W_3$	
SSN	Ν	М
185	Smith	1
185	Brown	3

	$W_4$	
SSN	Ν	М
185	Smith	1
185	Brown	4

	$W_5$	
SSN	Ν	М
185	Smith	1
186	Brown	1

	$W_6$	
SSN	Ν	М
185	Smith	1
186	Brown	2

#### Queries on Incomplete Databases

Given query Q and incomplete database W:

- An answer t is certain, if  $\forall : W_i \in \mathbf{W}, t \in Q(W_i)$
- An answer t is possible if  $\exists W_i \in \mathbf{W}, t \in Q(W_i)$

	$W_1$	
SSN	Ν	М
185	Smith	1
185	Brown	1

	$W_3$	
SSN	Ν	М
185	Smith	1
185	Brown	3

	$W_5$	
SSN	Ν	М
185	Smith	1
186	Brown	1

	$W_2$	
SSN	Ν	М
185	Smith	1
185	Brown	2

	$W_4$	
SSN	Ν	М
185	Smith	1
185	Brown	4

	$W_6$	
SSN	N	М
185	Smith	1
186	Brown	2

Let  $\mathbf{W} = \{W_1, \ldots, W_6\}.$ 

Query  $\exists_N \exists_M Census(S, N, M)$  has certain answer (185) and possible answers (185) and (186).

Query ∃<sub>S</sub>∃<sub>M</sub>Census(S, N, M) has the same possible and certain answers (Smith) and (Brown).
#### Queries on Incomplete Databases

Several studies on this started back in the 90s for various models, in particular conditional databases. [Abiteboul'91, O.'08a]

Hard tasks already for positive relational algebra:

- Tuple possibility is NP-complete
- Tuple certainty is coNP-complete

We next focus on probabilistic databases.

## Queries on Probabilistic Databases

Given query Q and probabilistic database  $(\mathbf{W}, P)$ : The Marginal Probability of an answer t is:  $P(t) = \sum \{P(W_i) \mid W_i \in \mathbf{W}, t \in Q(W_i)\}.$ 

$W_1: P(W_1) = 0.1$		
SSN	Ν	М
185	Smith	1
185	Brown	1

$W_3: P(W_3) = 0.1$		
SSN	Ν	М
185	Smith	1
185	Brown	3

$W_2: P(W_2) = 0.1$				
SSN	Ν	М		
185	Smith	1		
185	Brown	2		
- 50		105 010001 2		

$W_4: P(W_4) = 0.1$		
SSN	Ν	М
185	Smith	1
185	Brown	4

$W_5: P(W_5) = 0.3$		
SSN	Ν	М
185	Smith	1
186	Brown	1

$W_6: P(W_6) = 0.3$		
SSN	Ν	Μ
185	Smith	1
186	Brown	2

## Queries on Probabilistic Databases

Given query Q and probabilistic database (**W**, P): The Marginal Probability of an answer t is:  $P(t) = \sum \{P(W_i) \mid W_i \in \mathbf{W}, t \in Q(W_i)\}.$ 

$W_1: P(W_1) = 0.1$		
SSN	Ν	М
185	Smith	1
185	Brown	1

$W_3: P(W_3) = 0.1$		
SSN	Ν	М
185	Smith	1
185	Brown	3

$W_5: P(W_5) = 0.3$		
SSN	Ν	М
185	Smith	1
186	Brown	1

$W_2: P(W_2) = 0.1$		
SSN	Ν	М
185	Smith	1
185	Brown	2

$W_4: P(W_4) = 0.1$		
SSN	Ν	Μ
185	Smith	1
185	Brown	4

$W_6: P(W_6) = 0.3$		
SSN	Ν	М
185	Smith	1
186	Brown	2

- Let  $\mathbf{W} = \{ W_1, \dots, W_6 \}.$
- $\exists_N \exists_M Census(S, N, M):$ P(185) = 1 and P(186) = 0.6. $\exists_S \exists_M Census(S, N, M):$
- P(Smith) = P(Brown) = 1.

These are trivial queries! Computing the marginal probability is hard in general!

# Queries on Probabilistic Databases

Given query Q and probabilistic database (**W**, P): The Marginal Probability of an answer t is:  $P(t) = \sum \{P(W_i) \mid W_i \in \mathbf{W}, t \in Q(W_i)\}.$ 



 $\rightarrow$  Key challenge: Which queries admit efficient (polynomial time) computation of marginal probabilities for their answers?

# Representability of Query Answers

For a given query language Q and data model W: For any query  $Q \in Q$  and database  $\mathbf{W} \in W$ , is there  $\overline{Q} \in Q$  such that  $\overline{Q}(\mathbf{W}) = \{Q(W_i) \mid W_i \in \mathbf{W}\}$  and can be represented in W?



- This holds for relational algebra and PC databases: [Imielinski'84]  $\overline{Q}(T)$  is an extension of Q to also compute the *query lineage*.
- This does not hold for BIDs and TIs, but query lineage still useful for computing marginal probabilities of query answers on BIDs and TIs.
- This idea is also used by Trio and MayBMS. [Das Sarma'06, Antova'08]

# Query Lineage by Example

				(	)rdorc			Lineit	tem	
(	Custome	r	akay	ckov	data	φ.	okey	disc	ckey	Φ
ckey	name	Φ	OKEY	Скеу		Ψ	1	0.1	1	<i>z</i> <sub>1</sub>
1	Joe	X1	T	1	1995-01-10	<i>y</i> 1	1	0.2	1	22
2	Dan	Xo	2	1	1996-01-09	<i>y</i> 2	3	0.4	2	72
-	Dun	172	3	2	1994-11-11	<i>y</i> 3	2	0.1	2	-3
							3	0.1	2	24

Query asking for the dates of discounted orders shipped to customer 'Joe':  $\exists_C \exists_O \exists_D Customer(C, Joe), Orders(O, C, D), Lineitem(O, S, C), S > 0$ 

Query answer and lineage					
odate	Φ				
1995-01-10	$x_1y_1z_1 + x_1y_1z_2$				

 $\overline{Q}$  does Q and propagates the input conditions  $\Phi$  to the answers:

- join of tuples leads to conjunction of their conditions
- union/disjunction of tuples leads to disjunction of their conditions.

Query lineage traces the computation of an answer back to its input.

# Marginal Probabilities via Query Lineage

The marginal probability of a query answer is the probability of its lineage.

How to compute the lineage probability?

$x_1$	$y_1$	$z_1$	<i>z</i> <sub>2</sub>	$x_1y_1z_1 + x_1y_1z_2$	Probability
0	*	*	*	0	0
1	0	*	*	0	0
1	1	0	0	0	0
1	1	0	1	1	$P(x_1) \cdot P(y_1) \cdot (1 - P(z_1)) \cdot P(z_2)$
1	1	1	0	1	$P(x_1) \cdot P(y_1) \cdot P(z_1) \cdot (1 - P(z_2))$
1	1	1	1	1	$P(x_1) \cdot P(y_1) \cdot P(z_1) \cdot P(z_2)$

 $P(x_1y_1z_1 + x_1y_1z_2) = P(x_1) \cdot P(y_1) \cdot [1 - (1 - P(z_1)(1 - P(z_2))].$ 

Going over its truth table is exponential in the number of variables.

Two ideas:

[O.'08b]

Read-once lineage factorization

$$x_1y_1z_1 + x_1y_1z_2 = x_1y_1(z_1 + z_2)$$

Lineage compilation into polysize decision diagrams.

#### Where Are We Now?

We know how to compute the query answers using a simple query extension that also computes the query lineage.

We do not know yet how to compute the marginal probabilities of query answers efficiently.

Next part of the tutorial:

 Analyze the complexity of computing marginal probabilities as a function of database size and query structure.

# Outline



Why Probabilistic Databases? Probabilistic Data Models The Query Evaluation Problem Dichotomies for Query Evaluation • The Hard Queries

• The Tractable Queries

Ranking Queries

Next Steps

References

# Short Recap on Complexity Class #P (Sharp P)

#P =Class of functions f(x) for which there exists a PTIME non-deterministic Turing machine M such that f(x) = number of accepting computations of Mon input x. [Valiant'79]

#### Class of **counting problems** associated with decision problems in NP:

- **SAT** (given formula  $\phi$ , is  $\phi$  satisfiable?) is NP-complete
- #SAT (given formula φ, count # of satisfying assignments) is #P-complete

A PTIME machine with a #P oracle can solve any problem in polynomial hierarchy with one #P query. [Toda'91]

#SAT is #P-complete already for bipartite positive DNFs! [Provan'83]

• .. yet SAT is trivially PTIME for DNFs.

### Dichotomies for Queries on Probabilistic Databases

The following property has been observed for several classes  ${\cal Q}$  of relational queries on TI databases:

The data complexity of every query in Q is either **polynomial time** or **#P-hard**.

## Dichotomies for Queries on Probabilistic Databases

The following property has been observed for several classes  ${\cal Q}$  of relational queries on TI databases:

The data complexity of every query in Q is either **polynomial time** or **#P-hard**.

Examples of such classes Q of relational queries:

- NCQ: non-repeating conjunctive queries [Dalvi'07]
- NCQs under functional dependencies [O.'09]
- Quantified queries (division, set comparisons)
   [Fink'11]
- UCQ: unions of conjunctive queries [Dalvi'12]
- RNCQ: ranking NCQ [O.'12]
- 1RA<sup>-</sup>: NCQ's relational algebra counterpart extended with negation [Fink'16]

# Syntactic Characterizations of Tractable Queries

The tractable queries in (R)NCQ and  $1RA^-$  admit an efficient syntactic characterization via the *hierarchical* property.

A (Boolean) NCQ or  $1RA^-$  query Q is hierarchical if:

For every pair of distinct variables A and B in Q, there is no triple of relation symbols R, S, and T in Q such that:

- $R^{A \neg B}$  has query variable A and not B,
- $S^{AB}$  has both query variables A and B, and
- $T^{\neg AB}$  has query variable *B* and not in *A*.

Non-hierarchical queries:

- $\blacksquare \exists_{A} \exists_{B} [R(A) \land S(A, B) \land T(B)]$
- $\exists_{B} \left[ \exists_{A} (R(A) \land S(A, B)) \land \neg T(B) \right]$

 $\blacksquare \exists_B \Big[ T(B) \land \neg \exists_A \big( R(A) \land S(A, B) \big) \Big]$ 

Non-hierarchical queries:

- $\blacksquare \exists_{A} \exists_{B} [R(A) \land S(A, B) \land T(B)]$
- $\exists_B \left[ \exists_A (R(A) \land S(A, B)) \land \neg T(B) \right]$
- $\blacksquare \exists_B \Big[ T(B) \land \neg \exists_A \big( R(A) \land S(A, B) \big) \Big]$



Hierarchical queries:

- $\blacksquare \exists_A \exists_B [(R(A) \land S(A, B)) \land \neg T(A, B)]$
- $\blacksquare \exists_A \exists_B [(R(A) \land T(B)) \land \neg (U(A) \land V(B))]$

$$= \exists_{A} \exists_{B} \Big[ (M(A) \land N(B)) \land \neg \big[ (R(A) \land T(B)) \land \neg (U(A) \land V(B)) \big] \Big]$$

Hierarchical queries:

- $\blacksquare \exists_{A} \exists_{B} [ (R(A) \land S(A, B)) \land \neg T(A, B) ]$
- $\blacksquare \exists_A \exists_B [(R(A) \land T(B)) \land \neg (U(A) \land V(B))]$

$$= \exists_{A} \exists_{B} \left[ (M(A) \land N(B)) \land \neg [(R(A) \land T(B)) \land \neg (U(A) \land V(B))] \right]$$



# Outline



Why Probabilistic Databases? Probabilistic Data Models The Query Evaluation Problem Dichotomies for Query Evaluation

#### • The Hard Queries

• The Tractable Queries

Ranking Queries

Next Steps

References

#### Hardness Proof Idea

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

- Given a non-hierarchical  $1RA^-$  query Q and
- Any positive bipartite DNF formula Ψ over disjoint sets **X** and **Y** of random variables.
- $\#\Psi$  can be computed using linearly (in most cases constantly) many calls to an oracle for P(Q), where Q is evaluated on tuple-independent databases with sizes polynomial in the size of  $\Psi$ .

#### Simplest Example of Hardness Reduction

[Grädel'98, Dalvi'07]

Input formula and query:

• 
$$\Psi = x_1 y_1 \lor x_1 y_2 \lor x_2 y_1$$
 over sets  $\mathbf{X} = \{x_1, x_2\}, \mathbf{Y} = \{y_1, y_2\}$   
•  $Q = \exists_A \exists_B \left[ R(A) \land S(A, B) \land T(B) \right]$ 

Construct a TI database **D** such that  $\Psi$  annotates  $Q(\mathbf{D})$ :

- Column  $\Phi$  holds random variables in  $\Psi$ .
  - ► Notation: ⊤ (true)
- Variables also used as constants for A and B.
- $S(x_i, y_j, \top)$ :  $x_i y_j$  is a clause in  $\Psi$ .
- **R** $(x_i, \mathbf{x_i})$  and  $T(y_j, \mathbf{y_j})$ :  $x_i$  is a variable in **X** and  $y_j$  is a variable in **Y**.

R	Т	S	$R \land S \land T$	Q
ΑΦ	ΒΦ	ΑΒΦ	ΑΒΦ	φ
<i>x</i> <sub>1</sub> <b>x</b> <sub>1</sub>	<i>y</i> <sub>1</sub> <b>y</b> <sub>1</sub>	$x_1 y_1 \top$	x <sub>1</sub> y <sub>1</sub> x <sub>1</sub> y <sub>1</sub>	() Ψ
<i>x</i> <sub>2</sub> <b>x</b> <sub>2</sub>	<i>Y</i> 2 <b>Y</b> 2	$x_1 y_2 \top$	<i>x</i> <sub>1</sub> <i>y</i> <sub>2</sub> <b>x</b> <sub>1</sub> <i>y</i> <sub>2</sub>	
		$x_2 y_1 \top$	<i>x</i> <sub>2</sub> <i>y</i> <sub>1</sub> <b>x</b> <sub>2</sub> <i>y</i> <sub>1</sub>	

#### Simplest Example of Hardness Reduction

[Grädel'98, Dalvi'07]

Input formula and query:

• 
$$\Psi = x_1 y_1 \lor x_1 y_2 \lor x_2 y_1$$
 over sets  $\mathbf{X} = \{x_1, x_2\}, \mathbf{Y} = \{y_1, y_2\}$   
•  $Q = \exists_A \exists_B \left[ R(A) \land S(A, B) \land T(B) \right]$ 

Construct a TI database **D** such that  $\Psi$  annotates  $Q(\mathbf{D})$ :

- Column  $\Phi$  holds random variables in  $\Psi$ .
  - ► Notation: ⊤ (true)
- Variables also used as constants for A and B.
- $S(x_i, y_j, \top)$ :  $x_i y_j$  is a clause in  $\Psi$ .
- **R** $(x_i, \mathbf{x_i})$  and  $T(y_j, \mathbf{y_j})$ :  $x_i$  is a variable in **X** and  $y_j$  is a variable in **Y**.

R	Т	S	$R \land S \land T$	Q
ΑΦ	ΒΦ	Α Β Φ	ΑΒΦ	φ
<i>x</i> <sub>1</sub> <b>x</b> <sub>1</sub>	<i>y</i> <sub>1</sub> <b>y</b> <sub>1</sub>	$x_1 y_1 \top$	x <sub>1</sub> y <sub>1</sub> x <sub>1</sub> y <sub>1</sub>	() Ψ
<i>x</i> <sub>2</sub> <b>x</b> <sub>2</sub>	<i>Y</i> 2 <b>Y</b> 2	$x_1 y_2 \top$	<i>x</i> <sub>1</sub> <i>y</i> <sub>2</sub> <b>x</b> <sub>1</sub> <i>y</i> <sub>2</sub>	
		$x_2 y_1 \top$	<i>x</i> <sub>2</sub> <i>y</i> <sub>1</sub> <b>x</b> <sub>2</sub> <i>y</i> <sub>1</sub>	

Query Q is the only minimal hard pattern in case of queries without negation!

#### A Surprising Example of Hardness Reduction

Input formula and query:

• 
$$\Psi = x_1 y_1 \lor x_1 y_2$$
 over sets  $\mathbf{X} = \{x_1\}, \mathbf{Y} = \{y_1, y_2\}$   
•  $Q = \exists_A \Big[ R(\mathbf{A}) \land \neg \exists_B (T(B) \land S(\mathbf{A}, B)) \Big]$ 

Construct a TI database **D** such that  $\Psi$  annotates  $Q(\mathbf{D})$ :

■  $S(i, b, \top)$ : Clause *i* in  $\Psi$  has variable *b*.

■  $R(i, \top)$  and  $T(b, \neg b)$ : *i* is a clause and *b* is a variable in  $\Psi$ .

R	T	5	$T \wedge S$	$\exists_{B}(T \land S)$	$R \land$	$\neg \exists_{B}(T \land S)$
ΑΦ	ΒΦ	ΑΒΦ	ΑΒΦ	ΑΦ	A	Φ
1 ⊤	$x_1 \neg x_1$	$1 x_1 \top$	$1 x_1 \neg \mathbf{x_1}$	$\boxed{1 \ \neg \textbf{x_1} \lor \neg \textbf{y_1}}$	1	x1y1
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	$x_1y_2$
	<i>y</i> <sub>2</sub> ¬ <b>y</b> <sub>2</sub>	$2 x_1 \top$	$2 x_1 \neg \mathbf{x_1}$			
		2 <i>y</i> <sub>2</sub> ⊤	2 <i>y</i> <sub>2</sub> ¬ <b>y</b> <sub>2</sub>			

[Fink'16]

### A Surprising Example of Hardness Reduction

Input formula and query:

$$\Psi = x_1 y_1 \lor x_1 y_2 \text{ over sets } \mathbf{X} = \{x_1\}, \mathbf{Y} = \{y_1, y_2\}$$
$$Q = \exists_A \Big[ R(A) \land \neg \exists_B \big( T(B) \land S(A, B) \big) \Big]$$

Construct a TI database **D** such that  $\Psi$  annotates  $Q(\mathbf{D})$ :

■  $S(i, b, \top)$ : Clause *i* in  $\Psi$  has variable *b*.

■  $R(i, \top)$  and  $T(b, \neg b)$ : *i* is a clause and *b* is a variable in  $\Psi$ .

R	T	5	$T \wedge S$	$\exists_{B}(T \land S)$	$R \land$	$\neg \exists_{B}(T \land S)$
ΑΦ	ΒΦ	<u>Α</u> Β Φ	ΑΒΦ	ΑΦ	Α	Φ
1 ⊤	$x_1 \neg \mathbf{x_1}$	$1 x_1 \top$	$1 x_1 \neg \mathbf{x_1}$	$\boxed{1 \ \neg \textbf{x_1} \lor \neg \textbf{y_1}}$	1	x1y1
2 ⊤	$y_1 \neg \mathbf{y_1}$	$1 y_1 \top$	$1 y_1 \neg \mathbf{y_1}$	$2 \ \neg \textbf{x_1} \lor \neg \textbf{y_2}$	2	$x_1y_2$
	<i>y</i> <sub>2</sub> ¬ <b>y</b> <sub>2</sub>	$2 x_1 \top$	2 <i>x</i> <sub>1</sub> ¬ <b>x</b> <sub>1</sub>			
		$2 y_2 \top$	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!

[Fink'16]

# Outline



Why Probabilistic Databases?
Probabilistic Data Models
The Query Evaluation Problem
Dichotomies for Query Evaluation
The Hard Queries

#### • The Tractable Queries

Ranking Queries

Next Steps

References

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

Approach based on knowledge compilation

- For any TI 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 lineage Ψ.
- 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 TI 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 lineage Ψ.
- Compile  $\Psi$  into poly-size OBDD( $\Psi$ ).
- Compute probability of  $OBDD(\Psi)$  in time linear in its size.

Lineage of tractable 1RA<sup>-</sup> queries:

- Read-once for queries without negation (so NCQ) [O.'08b]
   It admits linear-size OBBDs.
- **Not** read-once for queries with negation [Fink'16]
  - It admits 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  $Q_1 \lor \cdots \lor Q_k$ .

#### 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  $Q_1, \dots, Q_k$ .

- 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.

Similar ideas used for the evaluation of inversion-free UCQs. [Jha'13]

# Query Evaluation Example (1/3)

Consider the following query and TI database:

$$Q = \exists_A \exists_B \Big[ \big( R(A) \land T(B) \big) \land \neg \big( U(A) \land V(B) \big) \Big]$$

R	T	U	V	$R \wedge T$	$R \wedge T \wedge \neg (U \wedge V)$
AΦ	ВΦ	ΑΦ	ВΦ	ΑΒ Φ	ΑΒ Φ
1 r <sub>1</sub>	1 t <sub>1</sub>	1 u <sub>1</sub>	1 v <sub>1</sub>	1 1 r <sub>1</sub> t <sub>1</sub>	$1 \ 1 \ r_1 t_1 \neg (u_1 v_1)$
2 r <sub>2</sub>	2 t <sub>2</sub>	2 u <sub>2</sub>	2 v <sub>2</sub>	1 2 r <sub>1</sub> t <sub>2</sub>	$1 \ 2 \ r_1 t_2 \neg (u_1 v_2)$
				2 1 r <sub>2</sub> t <sub>1</sub>	$2 \ 1 \ r_2 t_1 \neg (u_2 v_1)$
				2 2 r <sub>2</sub> t <sub>2</sub>	2 2 $r_2 t_2 \neg (u_2 v_2)$

# Query Evaluation Example (1/3)

Consider the following query and TI database:

$$Q = \exists_A \exists_B \Big[ \big( R(A) \land T(B) \big) \land \neg \big( U(A) \land V(B) \big) \Big]$$

R	T	U	V	$R \wedge T$	$R \wedge T \wedge \neg (U \wedge V)$
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>	1 1 $r_1 t_1$ 1 2 $r_1 t_2$ 2 1 $r_2 t_1$ 2 2 $r_2 t_2$	$ \begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$

The lineage 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 negation.

Query Evaluation Example (2/3)

 $\mathsf{Translate} \ \ Q = \exists_A \exists_B \Big[ \big( R(A) \land T(B) \big) \land \neg \big( U(A) \land V(B) \big) \Big] \ \mathsf{into} \ \mathsf{RC}^\exists :$ 

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

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 (3/3)

Compile query annotation into OBDD:



# Outline



• The Hard Queries

• The Tractable Queries

#### **Ranking Queries**

#### Ranking Answers in Probabilistic Databases

Given a NCQ query Q, a TI database **D**, and any two answers  $t_1, t_2 \in Q(\mathbf{D})$ , does  $P(t_1) \leq P(t_2)$  hold?

Motivation:

- Probabilities are mere degrees of uncertainty in the data and are not otherwise meaningful to the user.
- Users mostly care about the ranking of answers in decreasing order of their probabilities or about a few most likely answers.

# Ranking versus Query Evaluation

#### Two complementary observations

- 1. Probability computation for distinct answers may share a common factor
  - That can be computed only once
    - Save computation time for both query evaluation and ranking!
  - Or that can be uniformly ignored for all answers.
    - For ranking purposes, we may ignore computationally hard tasks!

Ranking is computationally easier than query evaluation.

# Ranking versus Query Evaluation

#### Two complementary observations

- 1. Probability computation for distinct answers may share a common factor
  - That can be computed only once
    - Save computation time for both query evaluation and ranking!
  - Or that can be uniformly ignored for all answers.
    - For ranking purposes, we may ignore computationally hard tasks!

Ranking is computationally easier than query evaluation.

- 2. To compute the **exact ranking** of query answers, **approximate probabilities** of the individual answers may suffice.
  - Compute lower and upper bounds on these probabilities.
  - Incrementally refine the bounds to the extent needed to rank the answers.

# Share Query Plans and Anytime Approximation

#### Approach with two main ingredients

[0.'12]

- $1. \ \mbox{Share} \ \mbox{query} \ \mbox{plans} \ \mbox{to} \ \mbox{detect} \ \mbox{factors} \ \mbox{common to} \ \mbox{query} \ \mbox{answers}$ 
  - Static analysis on the query structure to identify subqueries whose computation can be shared across distinct query answers.
  - Equivalently, they identify factors shared by lineage of query answers.
- 2. Ranking based on anytime deterministic approximate inference
  - Incremental compilation of lineage with shared factors into BDDs
  - Each compilation step refines lower and upper bounds on lineage probabilities
### Share Query Plans and Anytime Approximation

#### Approach with two main ingredients

- [0.'12]
- 1. Share query plans to detect factors common to query answers
  - Static analysis on the query structure to identify subqueries whose computation can be shared across distinct query answers.
  - Equivalently, they identify factors shared by lineage of query answers.
- 2. Ranking based on anytime deterministic approximate inference
  - Incremental compilation of lineage with shared factors into BDDs
  - Each compilation step refines lower and upper bounds on lineage probabilities

Alternative approach using FPRAS-based Monte Carlo

[Ré'07]

- Ranking with probabilistic guarantee only
- Not truly incremental
- Black box approach, structure and common factors of query lineage not exploited.

### Example

List topics posted by users who have mentioned their followers:  $Q(X) = \exists_Y \exists_Z \exists_U \text{Trends}(X, Y), \text{Follows}(Y, Z), \text{Mentions}(U, Y, Z), \text{Tweets}(U, Y).$ (User Y contributed to trendy topic X, user Y follows user Z, user Y mentions user Z in tweet U, tweet U of user Y.)

A share plan for Q is as follows



and corresponds to the following rewriting:

$$Q(X) = \text{Trends}(X, Y), Q'(Y)$$
  
 $Q'(Y) = \text{Follows}(Y, Z), \text{Mentions}(U, Y, Z), \text{Tweets}(U, Y)$ 

Several answers (X-values) can be paired with the same value y of the variable Y and thus share the lineage Q'(y).

For any value y, the query Q'(y) is non-hierarchical and thus #P-hard!

## Outline



• The Hard Queries • The Tractable Queries Next Steps

References

- PPDL, semantics given by a notion of probabilistic chase [Bárány'16]
- Incorporating ontologies
- Vast literature (including MLNs) but missing the declarativity aspect!

- PPDL, semantics given by a notion of probabilistic chase [Bárány'16]
- Incorporating ontologies
- Vast literature (including MLNs) but missing the declarativity aspect!
- Further push the barrier on complexity
  - See Dan Suciu's advanced lecture on lifted inference!
  - Understand tractability for probabilistic programs at large

- PPDL, semantics given by a notion of probabilistic chase [Bárány'16]
- Incorporating ontologies
- Vast literature (including MLNs) but missing the declarativity aspect!
- Further push the barrier on complexity
  - See Dan Suciu's advanced lecture on lifted inference!
  - Understand tractability for probabilistic programs at large
- Tractability not sufficient in practice, develop approximations with predictable performance

- PPDL, semantics given by a notion of probabilistic chase [Bárány'16]
- Incorporating ontologies
- Vast literature (including MLNs) but missing the declarativity aspect!
- Further push the barrier on complexity
  - See Dan Suciu's advanced lecture on lifted inference!
  - Understand tractability for probabilistic programs at large
- Tractability not sufficient in practice, develop approximations with predictable performance
- This tutorial assumed CWA. What about OWA, e.g., Google squares? See Guy van den Broeck et al's talk on open-world probabilistic databases!

- PPDL, semantics given by a notion of probabilistic chase [Bárány'16]
- Incorporating ontologies
- Vast literature (including MLNs) but missing the declarativity aspect!
- Further push the barrier on complexity
  - See Dan Suciu's advanced lecture on lifted inference!
  - Understand tractability for probabilistic programs at large
- Tractability not sufficient in practice, develop approximations with predictable performance
- This tutorial assumed CWA. What about OWA, e.g., Google squares? See Guy van den Broeck et al's talk on open-world probabilistic databases!
- Build open-source systems, provide benchmarks

# Outline



Dan Suciu Dan Olteanu Christopher Ré Christoph Koch

STATUESIS LECTURES ON DATA MANAGEMENT

Why Probabilistic Databases?

Probabilistic Data Models

The Query Evaluation Problem

Dichotomies for Query Evaluation

- The Hard Queries
- The Tractable Queries

Ranking Queries

Next Steps

#### References

### References: Main Reference



Dan Suciu, Dan Olteanu, Christopher Ré, Christoph Koch: Probabilistic Databases. Synthesis Lectures on Data Management, Morgan & Claypool Publishers 2011 http://www.morganclaypool.com/doi/abs/10.2200/ S00362ED1V01Y201105DTM016

### References: Some Probabilistic Database Systems (1/3)

Stanford: Trio

Parag Agrawal, Omar Benjelloun, Anish Das Sarma, Chris Hayworth, Shubha U. Nabar, Tomoe Sugihara, Jennifer Widom. *Trio: A System for Data, Uncertainty, and Lineage.* VLDB 2006: 1151-1154

UW: MystiQ

Jihad Boulos, Nilesh N. Dalvi, Bhushan Mandhani, Shobhit Mathur, Christopher Ré, Dan Suciu. *MYSTIQ: a system for finding more answers by using probabilities*. SIGMOD 2005: 891-893

 Cornell & Oxford: MayBMS & SPROUT
 Jiewen Huang, Lyublena Antova, Christoph Koch, Dan Olteanu. MayBMS: a probabilistic database management system. SIGMOD 2009: 1071-1074

Dan Olteanu, Jiewen Huang, Christoph Koch. *SPROUT: Lazy vs. Eager Query Plans for Tuple-Independent Probabilistic Databases.* ICDE 2009: 640-651

## References: Some Probabilistic Database Systems (2/3)

IBM Almaden & Rice: MCDB

Ravi Jampani, Fei Xu, Mingxi Wu, Luis Leopoldo Perez, Chris Jermaine, Peter J. Haas. *The Monte Carlo Database system: Stochastic analysis close to the data*. ACM Trans. Database Syst. 36(3): 18 (2011)

 LogicBlox & Technion & Oxford: PPDL
 [Bárány'16] Vince Bárány, Balder ten Cate, Benny Kimelfeld, Dan
 Olteanu, Zografoula Vagena. *Declarative Probabilistic Programming with Datalog.* ICDT 2016: 7:1-7:19

Maryland: PrDB

Prithviraj Sen, Amol Deshpande, Lise Getoor. *PrDB: managing and exploiting rich correlations in probabilistic databases*. VLDB J. 18(5): 1065-1090 (2009)

## References: Some Probabilistic Database Systems (3/3)

#### UC Berkeley: Bayestore

Daisy Zhe Wang, Eirinaios Michelakis, Minos N. Garofalakis, Joseph M. Hellerstein. *BayesStore: managing large, uncertain data repositories with probabilistic graphical models.* PVLDB 1(1): 340-351 (2008)

#### Purdue: Orion

Sarvjeet Singh, Chris Mayfield, Sagar Mittal, Sunil Prabhakar, Susanne E. Hambrusch, Rahul Shah. *Orion 2.0: native support for uncertain data*. SIGMOD 2008: 1239-1242

### Waterloo: Probabilistic ranking

Ihab F. Ilyas, Mohamed A. Soliman. *Probabilistic Ranking Techniques in Relational Databases*. Synthesis Lectures on Data Management, Morgan & Claypool Publishers 2011

many, many others..

### References: Applications (1/2)

Google Knowledge Vault

Xin Dong, Evgeniy Gabrilovich, Geremy Heitz, Wilko Horn, Ni Lao, Kevin Murphy, Thomas Strohmann, Shaohua Sun, Wei Zhang: Knowledge vault: a web-scale approach to probabilistic knowledge fusion. KDD 2014: 601-610

#### DeepDive

Feng Niu, Ce Zhang, Christopher Re, Jude W. Shavlik. *DeepDive: Web-scale Knowledge-base Construction using Statistical Learning and Inference*. VLDS 2012: 25-28

#### NELL

[Mitchell'15] Tom M. Mitchell et al. *Never-Ending Learning*. AAAI 2015: 2302-2310

#### Census data

**[Antova'07]** Lyublena Antova, Christoph Koch, Dan Olteanu. 10<sup>10<sup>6</sup></sup> Worlds and Beyond: Efficient Representation and Processing of Incomplete Information. ICDE 2007: 606-615

### References: Applications (2/2)

 Creating probabilistic databases
 [Stoyanovich'11] Julia Stoyanovich, Susan B. Davidson, Tova Milo, Val Tannen. Deriving probabilistic databases with inference ensembles. ICDE 2011: 303-314
 [Sarawagi'06] Rahul Gupta, Sunita Sarawagi. Creating Probabilistic Databases from Information Extraction Models. VLDB 2006: 965-976

 OCR with Staccato
 [Kumar'11] Arun Kumar, Christopher Ré. Probabilistic Management of OCR Data using an RDBMS. PVLDB 5(4): 322-333 (2011)

#### Possible data repairs

[Beskales'09] George Beskales, Mohamed A. Soliman, Ihab F. Ilyas, Shai Ben-David. *Modeling and Querying Possible Repairs in Duplicate Detection.* PVLDB 2(1): 598-609 (2009)

#### Querying open-world Google squares

**[Fink'11]** Robert Fink, Andrew Hogue, Dan Olteanu, Swaroop Rath. *SPROUT<sup>2</sup>: a squared query engine for uncertain web data.* SIGMOD 2011: 1299-1302

### References: Data Models (1/2)

Or-sets

[Imielinski'91] Tomasz Imielinski, Shamim A. Naqvi, Kumar V. Vadaparty. Incomplete Objects - A Data Model for Design and Planning Applications. SIGMOD 1991: 288-297

- World-set decompositions
  [O.'08a] Dan Olteanu, Christoph Koch, and Lyublena Antova. World-set decompositions: Expressiveness and efficient algorithms. Theor. Comput. Sci., 403:265-284, 2008
- Block-independent disjoint, x-relations, x-tables
  [Barbará'92] D. Barbará, H. Garcia-Molina, and D. Porter. *The management of probabilistic data*. IEEE Trans. on Knowl. and Data Eng., 4:487-502, 1992

[Poole'93] David Poole. Probabilistic horn abduction and bayesian networks. Artif. Intell., 64(1):81-129, 1993

### References: Data Models (1/2)

Conditional databases (c-tables)

[Imielinski'84] Tomasz Imielinski, Witold Lipski Jr. Incomplete Information in Relational Databases. J. ACM 31(4): 761-791 (1984)

#### Trio's ULDBs

**[Benjelloun'06]**: Omar Benjelloun, Anish Das Sarma, Alon Halevy, and Jennifer Widom. *ULDBs: databases with uncertainty and lineage*. In Proc. 32nd Int. Conf. on Very large Data Bases, pages 953-964, 2006

#### MayBMS' U-relations

**[Antova'08]** Lyublena Antova, Thomas Jansen, Christoph Koch, Dan Olteanu. *Fast and Simple Relational Processing of Uncertain Data*. ICDE 2008: 983-992

 Markov Logic Networks via queries over TI databases
 [Jha'13] Abhay Kumar Jha, Dan Suciu. Probabilistic Databases with MarkoViews. PVLDB 5(11): 1160-1171 (2012)

### References: Query Evaluation (1/2)

**Recent overview:** Nilesh Dalvi, Dan Olteanu. *Query Processing on Probabilistic Data*. Encyclopedia of Database Systems, 2nd ed. 2016

- Possible/certain answers in conditional databases (c-tables)
  [Abiteboul'91] Serge Abiteboul, Paris Kanellakis, and Gösta Grahne. On the representation and querying of sets of possible worlds. Theor. Comput. Sci., 78:159-187, 1991.
  [O.'08a] Dan Olteanu, Christoph Koch, and Lyublena Antova. World-set decompositions: Expressiveness and efficient algorithms. Theor. Comput. Sci., 403:265-284, 2008
- Representability of query answers, query lineage
  [Imielinski'84] Tomasz Imielinski, Witold Lipski Jr. Incomplete
  Information in Relational Databases. J. ACM 31(4): 761-791 (1984)
  [Das Sarma'06] Anish Das Sarma, Omar Benjelloun, Alon Y. Halevy,
  Jennifer Widom. Working Models for Uncertain Data. ICDE 2006: 7
  [Antova'08] Lyublena Antova, Thomas Jansen, Christoph Koch, Dan
  Olteanu. Fast and Simple Relational Processing of Uncertain Data. ICDE 2008: 983-992

### References: Query Evaluation (2/2)

Factorization of query lineage

**[0.08b]** Dan Olteanu and Jiewen Huang. Using OBDDs for efficient query evaluation on probabilistic databases. In Proc. 2nd Int. Conf. on Scalable Uncertainty Management, pages 326-340, 2008.

Dan Olteanu, Jakub Zavodny. *Factorised representations of query results:* size bounds and readability. ICDT 2012: 285-298

**[Jha'13]** Abhay Kumar Jha, Dan Suciu. *Knowledge Compilation Meets Database Theory: Compiling Queries to Decision Diagrams.* Theory Comput. Syst. 52(3): 403-440 (2013)

Ranking query answers

**[Ré'07]** Christopher Ré, Nilesh N. Dalvi, and Dan Suciu. *Efficient top-k query evaluation on prob- abilistic data*. ICDE 2007: 886-895

**[O.'12]** Dan Olteanu, Hongkai Wen. *Ranking Query Answers in Probabilistic Databases: Complexity and Efficient Algorithms.* ICDE 2012: 282-293

### References: Complexity of Probabilistic Query Evaluation

Complexity class #P

**[Valiant'79]** L. G. Valiant. *The complexity of computing the permanent*. Theor. Comput. Sci., 8(2):189-201, 1979

**[Toda'91]** Seinosuke Toda. *PP is as Hard as the Polynomial-Time Hierarchy.* SIAM J. Comput. 20(5): 865-877 (1991)

Hardness results

**[Provan'83]** J. Scott Provan and Michael O. Ball. *The complexity of counting cuts and of computing the probability that a graph is connected.* SIAM Journal on Computing, 12(4):777-788, 1983

**[Grädel'98]** Erich Gradel, Yuri Gurevich, and Colin Hirsch. The complexity of query reliability. PODS 1998: 227-234

References: Dichotomies for Probabilistic Query Evaluation

Non-repeating CQ

[Dalvi'07] Nilesh N. Dalvi, Dan Suciu. *Efficient query evaluation on probabilistic databases.* VLDB J. 16(4): 523-544 (2007)

Under functional dependencies: **[O.'09]** Dan Olteanu, Jiewen Huang, Christoph Koch: SPROUT: Lazy vs. Eager Query Plans for Tuple-Independent Probabilistic Databases. ICDE 2009: 640-651

Ranking: **[O.'12]** Dan Olteanu, Hongkai Wen. Ranking Query Answers in Probabilistic Databases: Complexity and Efficient Algorithms. ICDE 2012: 282-293

Union of conjunctive queries (UCQs)
 [Dalvi'12] Nilesh N. Dalvi, Dan Suciu. The dichotomy of probabilistic inference for unions of conjunctive queries. J. ACM 59(6): 30 (2012)

 Non-repeating conjunctive relational algebra (with negation)
 [Fink'16] Robert Fink, Dan Olteanu. Dichotomies for Queries with Negation in Probabilistic Databases. ACM Trans. Database Syst. 41(1): 4 (2016)

### References: Futher References

Conditioning probabilistic databases
 [Koch'08] Christoph Koch, Dan Olteanu. Conditioning probabilistic databases. PVLDB 1(1): 313-325 (2008)

#### Probabilistic XML

Michael Benedikt, Evgeny Kharlamov, Dan Olteanu, Pierre Senellart. Probabilistic XML via Markov Chains. PVLDB 3(1): 770-781 (2010)

Benny Kimelfeld, Pierre Senellart. *Probabilistic XML: Models and Complexity*. Advances in Probabilistic Databases for Uncertain Information Management 2013: 39-66

#### Beyond relational queries

Lei Chen, Xiang Lian. *Query Processing over Uncertain Databases*. Synthesis Lectures on Data Management, Morgan & Claypool Publishers 2012

## Thank you!