# **Probabilistic Databases**

Dan Olteanu (Oxford) London, November 2017

The National Archives



For the purpose of this introductory talk:

Probabilistic data =

Relational data

+

• Probabilities that measure the degree of uncertainty in the data.

Long-term research challenges:

- Models for probabilistic data to capture data and its uncertainty.
- Query evaluation = Probabilistic inference
   Query answers are annotated with output probabilities.

Early work (80s and 90s):

• Basic data models and query processing

Wong'82, Shoshani'82, Cavallo & Pittarelli'87, Barbara'92, Lakshmanan'97,'01, Fuhr& Roellke'97, Zimanyi'97.

Recent wave (2004 - now):

- Computational complexity of query evaluation
- Probabilistic database systems, just a few examples:
  - UW (MystiQ)
  - Stanford (Trio)
  - Cornell & Oxford (MayBMS/SPROUT)
  - IBM Almaden & Rice (MCDB)
  - Maryland, Waterloo, UBC, Florida, Purdue, Wisconsin
  - LogicBlox & Technion & Oxford (PPDL)

#### Why Probabilistic Databases?

Probabilistic Data Models

The Query Evaluation Problem

Query Evaluation: Complexity and Algorithms

Live Demo with MayBMS

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

<ul> <li>Manually entered data</li> </ul>	
Represent several possible readings with MayB	MS [Antova'07]
Infer missing data with meta-rule semi-lattices	[Stoyanovich'11]
Manage OCR data with Staccato/Google OCR	Ropus [Kumar'12]
• Data cleaning Represent several possible data repairs	[Beskales'09]
• Data integration Google Squared and SPROUT <sup>2</sup>	[Fink'11]

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

Possible segmentations of unstructured text

<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				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'

## Never-Ending Language Learner (NELL) database

#### Recently-Learned Facts twitter

Refresh

instance	iteration	date learned	confidence
<u>biscutate_swift</u> is an <u>animal</u>	211	18-feb-2011	100.0 🕼 🖏
<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 🖉 🖏
<u>dillon_county_airport</u> is an <u>airport</u>	210	17-feb-2011	93.8 🖓 🖏
the sports team toronto blue jays was the winner of n1993 world series	212	20-feb-2011	96.9 🖓 🖏
<u>mozart</u> is a person who <u>died at the age of</u> 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 🏠 🖏
<u>wutv_tv</u> 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 🖉 🖏

[Mitchell'15]

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

а

# MayBMS manages 10<sup>106</sup> possible readings of census data

[Antova'07]

Social Security Number:	185
Name:	Smith
Marital Status:	(1) single
	- ~ ]
Social Security Number:	186
Name:	Brown

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?

Social Security	y Number:	<u>280</u>	2		
Mari	tal Status:	(1) single (3) divorced	(2) married <b>*</b> (4) widowed <b>□</b>		
Social Securit	Social Security Number: 85				
	Name: REDWA				
Mari	Marital Status: (1) single (2) married (3) divorced (4) widowed				
$(\Pi D)$	22IN	IN	IVI		
$t_1$	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 lists of page links as answers to Google queries
- 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	iedy movi	ies			
	Item Nar	ne 💌	Language	Director 💌 🗙	Release Date
X	The Mas	k	English	Chuck Russell	29 July 1994
×	Scary M	English     Ianguage for th     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 som	<ul> <li>John R. Dilworth director for The Ma www.freebase.com</li> </ul>	Low confidence Isk
×	Music	english, french languages for ti www.dvdreview.	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.c	<b>ge</b> 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	»

Why Probabilistic Databases?

#### Probabilistic Data Models

The Query Evaluation Problem

Query Evaluation: Complexity and Algorithms

Live Demo with MayBMS

#### **Revisiting the Census Data Example**



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.



13/48

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	Μ
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 <sub>4</sub>	
SSN	Ν	Μ
185	Smith	1
185	Brown	4

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

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

. . .

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



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

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

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

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

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

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

For 
$$\mathbf{W} = \{W_1, \dots, W_6\}$$
, $\sum_{W_i \in \mathbf{W}} P(W_i) = 1.$ 

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

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

### Succinct Representations of Incomplete/Probabilistic Data

Social Security Number: Name:	<u> 185</u> Smith
Marital Status:	(1) single ⊠ (2) married (3) divorced □ (4) widowed □
Social Security Number:	$\frac{185}{Brown$

Succinct or-set representation:

[Imielinski'91]

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

It exploits independence of different fields:

- Choice for Smith's SSN independent of choice of for Brown's SSN.
- The probability distributions associated with these choices are independent.

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

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

<u>RID</u>	N	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 <sub>2</sub>	1	0.25
t <sub>2</sub>	2	0.25
t <sub>2</sub>	3	0.25
t <sub>2</sub>	4	0.25

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

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

<u>RID</u>	N	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 <sub>2</sub>	1	0.25
t <sub>2</sub>	2	0.25
t <sub>2</sub>	3	0.25
t <sub>2</sub>	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	P
$t_1$	1	0.9
$t_1$	2	0.1
t <sub>2</sub>	1	0.25
$t_2$	2	0.25
t <sub>2</sub>	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.

BIDs also allow blocks with probabilities less than 1:

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

RID	Ν	Р
$t_1$	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 <sub>2</sub>	1	0.25
$t_2$	2	0.25
t <sub>2</sub>	3	0.25
t <sub>2</sub>	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>
$t_2$	186	0.2	t <sub>2</sub>

RID	N	Ρ
t <sub>2</sub>	Brown	1

RID	М	Р
t <sub>2</sub>	4	0.25

The probability of this world is

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

*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	<u>RID</u>	Ν	Ρ	<u>RID</u>	Μ	Р
$t_1$	185	0.7	$t_1$	Smith	1	$t_1$	1	0.9
$t_2$	185	0.8	t <sub>2</sub>	Brown	1	$t_2$	2	0.25

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.

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  $t_1$  and  $t_2$  have SSN = 185.

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

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		RID	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  $t_1$  and  $t_2$  have the same SSN 185 by using contradicting assignments for variable *X*.
- Transfer probabilities of tuples to probability distributions of variables.

A *PC database* is  $(D, X, \Phi)$ , 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	Ρ
$t_1$	185	X = 1	X	1	0.6
$t_1$	785	X = 2	X	2	0.3
t <sub>2</sub>	185	$Y = 1 \land X \neq 1$	Y	1	0.8
t <sub>2</sub>	186	Y = 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).

Recall our previous TI database example:

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

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

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

Here is a PC encoding of the above TI database:

RID	SSN	Φ	Ρ	RID	Ν	Φ	Ρ	RID	М	Φ	Р
$t_1$	185	<b>s</b> 1	0.7	$t_1$	Smith	<b>n</b> 1	1	$t_1$	1	<b>m</b> <sub>1</sub>	0.9
$t_2$	185	<b>s</b> <sub>2</sub>	0.8	t <sub>2</sub>	Brown	<b>n</b> <sub>2</sub>	1	t <sub>2</sub>	2	<i>m</i> <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's U-relations [Antova'07].
  - Completeness: They can represent any probabilistic database.
- Mid-level: block-independent disjoint databases.
  - 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.

[Zimanyi'97]

[Barbará'92]

- The norm in real-world repositories like Google's, DeepDive, and NELL.
- Not complete even via unions of conjunctive queries.
- Most theoretical work on complexity of query evaluation done for them.

Why Probabilistic Databases?

Probabilistic Data Models

The Query Evaluation Problem

Query Evaluation: Complexity and Algorithms

Live Demo with MayBMS

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  $Q(\mathbf{W}) = \{Q(W_1), \dots, Q(W_n)\}.$ 

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(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	Μ
185	Smith	1
185	Brown	4

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

	$W_6$	
SSN	Ν	Μ
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	Μ
185	Smith	1
185	Brown	3

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

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

	147	
	VV4	
SSN	Ν	Μ
185	Smith	1
185	Brown	4

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

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

- Query  $Q_1(s) = Census(s, n, m)$ has certain answer (185) and possible answers (185) and (186).
- Query
   Q<sub>2</sub>(n) = Census(s, n, m)
   has the same possible and
   certain answers (Smith)
   and (Brown).

Several studies on this date back to 90s for various models, in particular conditional databases. [Abiteboul'91, Grahne'91, O.'08]

Hard tasks already for positive relational algebra:

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

We next focus 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	Μ
185	Smith	1
185	Brown	1

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

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

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

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

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

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	Μ		
185	Smith	1		
185	Brown	1		

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

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

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

- Let  $\mathbf{W} = \{ W_1, \ldots, W_6 \}.$
- $Q_1(s) = Census(s, n, m)$ : P(185) = 1 and P(186) = 0.6.
- $Q_2(n) = Census(s, n, m)$ : P(Smith) = P(Brown) = 1.

 $Q_1$  and  $Q_2$  are trivial queries! Computing the marginal probability is hard in general! 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	V	$V_2$ :	$P(W_2) =$	0.1	
	SSN	Ν	М	S	5N	Ν	Μ	Let $\mathbf{W} = \{W_1, \ldots, W_6\}$
	185	Smith	1	1	85	Smith	1	
	185	Brown	1	1	85	Brown	2	• $Q_1(s) = Census(s, n, n)$
								D(195) = 1 and
	W3 :	$P(W_3) =$	0.1	V	V4 :	$P(W_4) =$	0.1	F(105) = 1 and
	SSN	Ν	Μ	S	5N	Ν	Μ	P(186) = 0.6.
	185	Smith	1	1	85	Smith	1	
	185	Brown	3	1	85	Brown	4	• $Q_2(n) = \text{Census}(s, n,$
ĺ								P(Smith) = P(Brown)
	$W_5$ :	$P(W_5) =$	0.3	V	V <sub>6</sub> :	$P(W_6) =$	0.3	$O_1$ and $O_2$ are trivial $\sigma$
	SSN	Ν	Μ	S	SN	Ν	Μ	
	185	Smith	1	1	85	Smith	1	Computing the margina
	186	Brown	1	1	86	Brown	2	probability is hard in ge

 $\rightarrow$  Key challenge: Which queries admit efficient (polynomial time) computation of marginal probabilities for their 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. [Benjelloun'06, Antova'09]

				0	Jrdors			Linei	tem	
	Custome	r						disc	ckey	φ
ckev	name	φ	Okey	Скеу	uate	$\downarrow \Psi$	1	0.1	1	Z1
1	lee		1	1	1995-01-10	<i>y</i> 1	1	0.0	1	
T	Joe	×1	2	1	1996-01-09	Va	T	0.2	1	22
2	Dan	<i>x</i> <sub>2</sub>	-	-	1004 11 11	52	3	0.4	2	Z3
			3	2	1994-11-11	<i>Y</i> 3	3	0.1	2	Z4

Query asking for the dates of discounted orders shipped to customer 'Joe': Customer(*ckey*, Joe), Orders(*okey*, *ckey*, *date*), Lineitem(*okey*, *disc*, *ckey*)

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

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.

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	$P(x_1) \cdot P(y_1) \cdot P(z_1) \cdot P(z_2)$

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

• The truth table is exponential in the number of variables.

Two ideas:

[O.'08+]

- 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.

- 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, more technical part of the tutorial (not covered here):

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

Why Probabilistic Databases?

Probabilistic Data Models

The Query Evaluation Problem

Query Evaluation: Complexity and Algorithms

Live Demo with MayBMS

#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  $\phi$ , 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&Suciu'04]
- Quantified queries (division, set comparisons) [Fink&O.'11]
- UCQ: unions of conjunctive queries [Dalvi&Suciu'12]
- RNCQ: ranking NCQ [O.&Wen'12]
- 1RA<sup>-</sup>: NCQ's relational algebra counterpart extended with negation [Fink&O.'14]

### 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.

#### Examples

Hierarchical queries:

- $\exists_{\mathbf{A}} \exists_{\mathbf{B}} [ (R(\mathbf{A}) \land S(\mathbf{A}, \mathbf{B})) \land \neg T(\mathbf{A}, \mathbf{B}) ]$
- $\exists_{A} \exists_{B} [(R(A) \land T(B)) \land \neg (U(A) \land V(B))]$
- $\exists_{A} \exists_{B} \Big[ (M(A) \land N(B)) \land \neg [(R(A) \land T(B)) \land \neg (U(A) \land V(B))] \Big]$

#### Examples

Hierarchical queries:

•  $\exists_{\mathbf{A}} \exists_{\mathbf{B}} [ (R(\mathbf{A}) \land S(\mathbf{A}, \mathbf{B})) \land \neg T(\mathbf{A}, \mathbf{B}) ]$ 

• 
$$\exists_{\mathbf{A}} \exists_{\mathbf{B}} [(\mathbf{R}(\mathbf{A}) \land T(\mathbf{B})) \land \neg (U(\mathbf{A}) \land V(\mathbf{B}))]$$

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

Non-hierarchical queries:

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

#### Hardness Proof Idea

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

- Given a non-hierarchical  $1RA^-$  query Q and
- Any positive bipartite DNF formula  $\Psi$  over disjoint sets  ${\bf X}$  and  ${\bf Y}$  of random variables.
- #Ψ 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 Ψ.

#### [Grädel'98, Dalvi& Suciu'04]

Input formula and query:

- $\Psi = x_1y_1 \lor x_1y_2 \lor x_2y_1$  over sets  $\mathbf{X} = \{x_1, x_2\}, \mathbf{Y} = \{y_1, y_2\}$
- $Q = \exists_{\mathbf{A}} \exists_{\mathbf{B}} \left[ R(\mathbf{A}) \land S(\mathbf{A}, \mathbf{B}) \land T(\mathbf{B}) \right]$

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	<b>Α Β</b> 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>	$x_1 y_1 \top$	<i>x</i> <sub>1</sub> <i>y</i> <sub>1</sub> <b>x</b> <sub>1</sub> <i>y</i> <sub>1</sub>	() Ψ
<i>x</i> <sub>2</sub> <b>x</b> <sub>2</sub>	<i>Y</i> 2 <b>Y2</b>	$x_1 y_2 \top$	<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>	

#### [Grädel'98, Dalvi& Suciu'04]

Input formula and query:

- $\Psi = x_1y_1 \lor x_1y_2 \lor x_2y_1$  over sets  $\mathbf{X} = \{x_1, x_2\}, \mathbf{Y} = \{y_1, y_2\}$
- $Q = \exists_{\mathbf{A}} \exists_{\mathbf{B}} \left[ R(\mathbf{A}) \land S(\mathbf{A}, \mathbf{B}) \land T(\mathbf{B}) \right]$

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	<b>Α Β</b> 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>	$x_1 y_1 \top$	<i>x</i> <sub>1</sub> <i>y</i> <sub>1</sub> <b>x</b> <sub>1</sub> <i>y</i> <sub>1</sub>	() Ψ
<i>x</i> <sub>2</sub> <b>x</b> <sub>2</sub>	<i>Y</i> 2 <b>Y2</b>	$x_1 y_2 \top$	<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! 39/48

#### A Surprising Example of Hardness Reduction

Input formula and query:

- $\Psi = x_1y_1 \lor x_1y_2$  over sets  $X = \{x_1\}, Y = \{y_1, y_2\}$
- $Q = \exists_{\mathbf{A}} \Big[ R(\mathbf{A}) \land \neg \exists_{\mathbf{B}} \big( T(\mathbf{B}) \land S(\mathbf{A}, \mathbf{B}) \big) \Big]$

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	Т	5	$T \wedge S$	$\exists_{B}(T \land S)$	$R \wedge$	$\neg \exists_{B}(T \land S)$
ΑΦ	ΒΦ	<b>Α Β</b> Φ	ΑΒΦ	ΑΦ	Α	φ
1 ⊤	$x_1 \neg x_1$	$1 x_1 \top$	$1 x_1 \neg \mathbf{x}_1$	$1 \ \neg x_1 \lor \neg y_1$	1	x1y1
2 T	$y_1 \neg \mathbf{y_1}$	$1 y_1 \top$	$1 y_1 \neg \mathbf{y}_1$	$2 \ \neg \textbf{x_1} \lor \neg \textbf{y_2}$	2	$x_1y_2$
	<i>y</i> <sub>2</sub> ¬ <b>y</b> <sub>2</sub>	$2 x_1 \top$	$2 x_1 \neg x_1$			
		$2 y_2 \top$	2 <i>y</i> <sub>2</sub> ¬ <b>y</b> <sub>2</sub>			

#### A Surprising Example of Hardness Reduction

Input formula and query:

- $\Psi = x_1y_1 \lor x_1y_2$  over sets  $X = \{x_1\}, Y = \{y_1, y_2\}$
- $Q = \exists_{\mathbf{A}} \Big[ R(\mathbf{A}) \land \neg \exists_{\mathbf{B}} \big( T(\mathbf{B}) \land S(\mathbf{A}, \mathbf{B}) \big) \Big]$

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	$T \wedge S$	$\exists_{\mathbf{B}}(T \land S)$	$R \wedge $	$\neg \exists_{\mathbf{B}}(T \land S)$
ΑΦ	ΒΦ	ΑΒΦ	ΑΒΦ	ΑΦ	Α	φ
_1 ⊤	$x_1 \neg x_1$	$1 x_1 \top$	$1 x_1 \neg \mathbf{x_1}$	$\boxed{1 \ \neg x_1 \lor \neg y_1}$	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>	$2 x_1 \top$	$2 x_1 \neg \mathbf{x_1}$			
		$2 y_2 \top$	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!

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.

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. & Huang'08] It admits linear-size OBBDs.
- Not read-once for queries with negation
  - It admits OBBDs of size linear in the database size <u>but</u> exponential in the query size.

[Fink& O'14]

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<sub>1</sub> ∨ · · · ∨ Q<sub>k</sub>.
- 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& Suciu'12]

# 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	$R \wedge T$	$R \wedge T \wedge \neg (U \wedge V)$
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} \end{array} $	$\begin{array}{cccc} 1 & 1 & r_{1}t_{1}\neg(u_{1}v_{1}) \\ 1 & 2 & r_{1}t_{2}\neg(u_{1}v_{2}) \\ 2 & 1 & r_{2}t_{1}\neg(u_{2}v_{1}) \\ \end{array}$

## 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	$R \wedge T$	$R \wedge T \wedge \neg (U \wedge V)$
ΑΦ	ВΦ	ΑΦ	ВΦ	ΑΒΦ	ΑΒΦ
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 $r_1t_1$ 1 2 $r_1t_2$ 2 1 $r_2t_1$ 2 2 $r_2t_2$	$\begin{array}{rrrr} 1 & 1 & r_1 t_1 \neg (u_1 v_1) \\ 1 & 2 & r_1 t_2 \neg (u_1 v_2) \\ 2 & 1 & r_2 t_1 \neg (u_2 v_1) \\ 2 & 2 & r_2 t_2 \neg (u_2 v_2) \end{array}$

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.
- $Q_1 \vee Q_2$  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:



Why Probabilistic Databases?

Probabilistic Data Models

The Query Evaluation Problem

Query Evaluation: Complexity and Algorithms

Live Demo with MayBMS

### The MayBMS/SPROUT Probabilistic Database Management System



http://maybms.sourceforge.net

- Extension of open-source PostgreSQL relational DBMS
- Available manual, worked-out examples, and publications
- 3-min video demo: http://www.cs.ox.ac.uk/dan.olteanu/maybms.mp4

We now go over a few simple examples from the manual.

# Thank you!