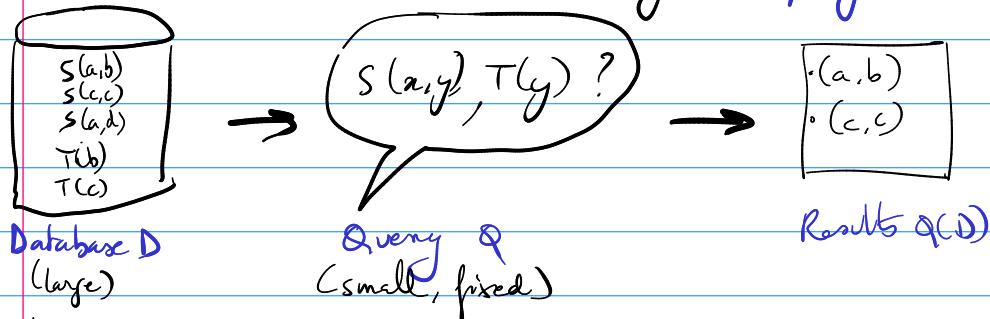


The (\ln) Tractability of Query Evaluation over Probabilistic Databases

Introduction and Background

Central problem in database theory : query evaluation



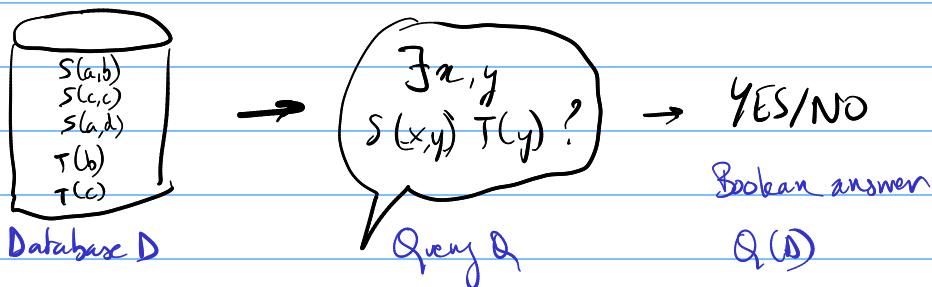
We study the data complexity of this problem, where Q is fixed

\Rightarrow What is the complexity, given D , of computing $Q(D)$?

In this talk we will only distinguish two complexity regimes: \rightarrow tractable (in PTIME)

\rightarrow intractable (NP-hard or worse)

Up to a PTIME transformation (listing each possible output), we can assume that the query is Boolean



For most database query languages (CQL, UCQL, Datalog...) this problem is always tractable

In this talk: we study a probabilistic version of this problem, where the data is uncertain

\rightarrow Practical motivations: real-world data is often uncertain, noisy, etc

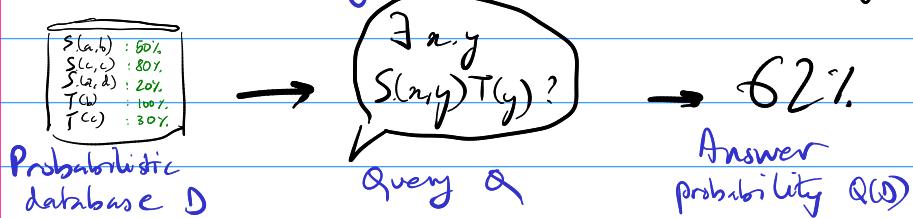
\rightarrow Theoretical motivations: the corresponding counting problem is interesting

Probabilistic database :

$S(a,b)$: 50%
$S(c,c)$: 80%
$S(a,d)$: 20%
$T(b)$: 100%
$T(c)$: 30%

- each fact is annotated by a probability
- shows how likely the fact is to exist
- we assume independence across facts

The Probabilistic Query Evaluation problem for a query Q ($PQE(Q)$)



Unfortunately the problem $PQE(Q)$ can now become intractable

⇒ For which queries Q can we tractably solve $PQE(Q)$?

Connections to :

- Graphical models : inference in Bayesian networks, Markov networks...
- Knowledge compilation : the study of tractable formalisms to represent Boolean functions
- Model counting of Boolean formulas, counting complexity...

Structure of the talk

- Preliminaries
- Upper bounds
 - For tractable self-join-free conjunctive queries
 - Via provenance circuits
 - For tractable UCQs
 - Intensional-extensional conjecture
- Lower bounds
 - #P-hardness for intractable self-join-free conjunctive queries
 - Extending to uniform reliability
 - (#P-hardness of intractable UCQs)
 - #P-hardness of unbounded UCQs
- Conclusion

Preliminaries : database theory

- Relational signature σ : predicate names and arity

R with arity 2, S with arity 2
 $a, b, c \dots$

- Constants

- Fact F over σ : predicate and constants

$R(a), S(a, b), S(c, c)$

- Instance I : finite set of facts

also called: database, structure, model ...

$\{R(a), S(a, d)\}$

- Variables

$x, y, z \dots$

- Atom over σ : like a fact but with variables

$R(x)$

- Conjunctive query (CQ): existentially quantified conjunction of atoms

$\exists x, y R(x) \wedge S(x, y)$
or just $R(a), S(a, y)$

- Important special case: Self-Join-Free CQ (SJFCQ)

$R(x) S(x, y)$ ok

→ every relation is used at most once

$S(x, y) S(y, z)$ not ok

- Union of conjunctive queries (UCQ): disjunction of CQs

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

equivalently: positive Boolean first-order formula

Note: we can reuse variables across disjuncts:

$S(x, y) \vee T(y) \quad \text{vs} \quad S(x, y) \vee T(z)$

- Homomorphism from a CQ Q to an instance I
witnesses that I satisfies Q : $I \models Q$

$Q: R(x) S(x, y) \quad I = \{R(a) S(a, a)\}$

$h(x) := a, h(y) := a$

Also characterizes CQ entailment: Q implies Q' ($\forall I, I \models Q \Rightarrow I \models Q'$)

iff Q' has a homomorphism into Q

Preliminaries: probabilistic databases

- Tuple-independent database (TID): $I \in (I, \pi)$:

- an instance I

- a function $\pi: I \rightarrow [0,1]$

$$S(a,b) : 50\%$$

$$S(c,d) : 80\%$$

$$S(a,d) : 20\%$$

$$T(b) : 100\%$$

$$T(c) : 30\%$$

- Possible world: a subinstance $I' \subseteq I$

A TID is a concise representation of a probability distribution over its possible worlds, where each fact has the indicated probability and we assume independence.

$$S(a,b) : 50\% \rightarrow 50\%$$

$$S(c,d) : 80\% \rightarrow 20\%$$

$$S(a,d) : 20\% \rightarrow 20\%$$

$$T(b) : 100\% \rightarrow 100\%$$

$$T(c) : 30\% \rightarrow 70\%$$

Formula for $I' \subseteq I$:

$$\Rightarrow \text{prob of possible world: } 50\% \times 20\% \times 20\% \times 100\% \times 70\%$$

$$\Pr_I(I') = \prod_{F \in I'} \underbrace{\pi(F)}_{\text{kept facts}} \times \prod_{F \in I \setminus I'} \underbrace{(1 - \pi(F))}_{\text{discarded facts}}$$

- Probability of a query Q on a TID $I = (I, \pi)$

$$\Pr_I(Q) = \sum_{\substack{I' \subseteq I \\ I' \models Q}} \Pr_I(I')$$

$$S(a,b) : 50\%$$

$$S(c,d) : 80\%$$

$$S(a,d) : 20\%$$

$$T(b) : 100\%$$

$$T(c) : 30\%$$

Query: $S(x,y) \wedge T(y)$

Proba:

$$1 - (1 - 50\%) \times (1 - 80\%) \times 30\%$$

$$= 62\%$$

- Probabilistic Query Evaluation problem for query Q : $PQE(Q)$:

- input: TID I

- output: compute $\Pr_I(Q)$.

- Special case: Uniform Reliability for Q , $UR(Q)$:

- input: instance I

- output: compute how many subinstances satisfy Q : $|\{I' \subseteq I \mid I' \models Q\}|$

$\Rightarrow UR(Q)$ is $PQE(Q)$ in the case where all probabilities are $\frac{1}{2}$ (or 0) up to factor $\frac{1}{2^{|I|}}$.

Goal: For which queries Q is $PQE(Q)$ tractable?

Tractable cases for SJFCQs: examples

Example 1: $Q_1 = S(x, y)$

How to solve $\text{PDE}(Q_1)$?

$$\rightarrow \Pr_D(Q_1) = \Pr_D(\text{some } S\text{-fact in } D)$$

$$= \Pr_D \left(\bigvee_{\substack{F \text{ S-fact} \\ \in D}} F \text{ is kept} \right) \xrightarrow{\text{Independent OR}}$$

$$= 1 - \Pr_D \left(\bigwedge_{\substack{F \text{ S-fact} \\ \in D}} F \text{ is discarded} \right) \xrightarrow{\text{independent AND}}$$

$$= 1 - \prod_{\substack{F \text{ S-fact} \\ \in D}} \Pr_D(F \text{ is discarded})$$

$$\begin{aligned} & 1 - (1 - 50\%) \\ & \times (1 - 80\%) \\ & \times (1 - 20\%) \\ & = 92\% \end{aligned}$$

Example 2: $Q_2 = S(x, y), S'(x', y')$.

\rightarrow We know how to compute $\Pr_D(S(x, y))$ (above) and $\Pr_D(S'(x', y'))$ (ditto)

$$\text{and: } \Pr_D(Q_2) = \Pr_D(S(x, y) \text{ AND } S'(x', y'))$$

\uparrow
independent AND

\rightarrow uses self-join-freeness

\rightarrow uses the fact that there are no stored variables

$$\text{so } \Pr_D(Q_2) = \Pr_D(S(x, y)) \times \Pr_D(S'(x', y'))$$

Example 3: $Q_3 = R(x) S(x, y)$

\rightarrow The variable x occurs in each atom so it partitions the database facts

$$\text{and: } \Pr_D(Q_3) = \Pr_D \left(\bigvee_{\substack{a \text{ constant} \\ \in D}} R(a) \wedge S(a, x) \right)$$

\nwarrow independent OR

$$= 1 - \prod_{\substack{a \in D \\ \text{acb}}} \left(1 - \Pr_D(R(a) \wedge S(a, x)) \right)$$

\rightarrow uses self-join-freeness
 \rightarrow each database fact is "relevant" for at most one disjunct depending on the value of its first element

$$\text{now } \Pr_D(R(a) \wedge S(a, x)) \text{ is } \#(R(a)) \times \Pr_D \left(\bigvee_{\substack{b \text{ constant} \\ \in D}} S(a, b) \right) \xrightarrow{\text{independent OR}}$$

$$= \#(R(a)) \times \Pr_D \left(\bigvee_{\substack{b \text{ constant} \\ \in D}} S(a, b) \right)$$

$$\text{so } \Pr_D(Q_3) = 1 - \prod_{\substack{a \in D \\ \text{R(a) fact}}} \left(1 - \#(R(a)) \times \left(1 - \prod_{\substack{b \in D \\ \text{S(a,b) fact}}} \left(1 - \#(S(a, b)) \right) \right) \right)$$

which is PTIME...

Tractable cases for SJFCQ : tractability criterion

Def: A hierarchical expression is a (Boolean) CQ built from:

- atoms

$$R(x)$$

- conjunction

$$\Theta_1 \wedge \Theta_2$$

- existential quantification

$$\exists x \ Q(x, y_1, \dots, y_n)$$

where the quantified variable occurs in each atom under the quantification

\Rightarrow A CQ is hierarchical if it can be written as a hierarchical expression

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

$\exists x (R(x) \wedge \exists y (S(x, y) \wedge \underline{T(y)}))$ not ok

Example: On an arity-1 signature, all SJFCQs admit hierarchical expressions.

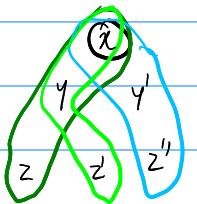
$$(\exists x R_1(x) \wedge R_2(x)) \wedge (\exists y R_3(y) \wedge R_4(y)) \wedge \dots$$

On an arity-2 signature, the hierarchical SJFCQs consist of (possibly multiple) star patterns



Each star pattern has a "center" and all facts of the pattern involve the center

On a higher-arity signature, hierarchical SJFCQs are more complex to represent



$$\begin{aligned}
 & (\exists x (\exists y (\exists z R_1(xyz)) \wedge (\exists z' R_2(xyz')))) \\
 & \quad \wedge (\exists y (\exists z'' R_3(xyzz'') \wedge (\exists w R_4(uvw))) \wedge \exists w R_5(uvw))
 \end{aligned}$$

The hierarchical SJFCQs are all tractable:

Thm: Let Q be a SJFCQ. If Q is hierarchical then $PQE(Q)$ is in PTIME.

Tractable cases for SJFCQ: tractability proof

Note: We extend the definition of hierarchical expressions to allow atoms featuring constants.

Claim: A hierarchical expression can be decomposed by applying the rules:

- conjunction $\mathcal{Q}_1 \wedge \dots \wedge \mathcal{Q}_n$ of hierarchical expressions that share no variable
- existential quantification $\exists x \mathcal{Q}$ where x occurs in all atoms of \mathcal{Q} and $\mathcal{Q}[x/a]$ is a hierarchical expression
- atoms involving only constants

Proof: Immediate (top-down rephrasing of the bottom-up definition)

Now the probability of a hierarchical expression can be computed via the rules:

- conjunction: $\Pr_0(\mathcal{Q}_1 \wedge \dots \wedge \mathcal{Q}_n) = \Pr_0(\mathcal{Q}_1) \times \dots \times \Pr_0(\mathcal{Q}_n)$
independent AND
→ uses self-join-freeness
- existential quantification: $\Pr_0(\exists x \mathcal{Q}) = \Pr_0\left(\bigvee_a \mathcal{Q}[x/a]\right)$
independent OR
→ uses self-join-freeless
→ uses the fact that x occurs in all atoms
- atoms $\Pr_0(F) = \pi(F)$ for F a fact (atom with only constants)

⇒ This computes the probability of a hierarchical SJFCQ in linear time in the input TID assuming that arithmetic operations take constant time

Let us see this result in a different way, using Boolean circuits

- define the problem of computing the probability of a circuit
- define tractable circuit classes for the problem
- explain how to reduce PQE to compute the probability of a provenance circuit
- explain why, for SJFCQs, we can make the circuit tractable

Circuits, probabilities, and tractable circuits

- Boolean circuit : DAG where the vertices (gates) are :

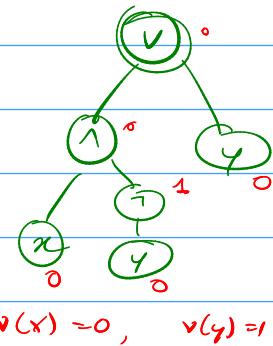
 - Variables :

 - Boolean operators \wedge, \vee, \neg (indegree 1)

 - Distinguished output gate

- Valuation of variables X : function $v: X \rightarrow \{0, 1\}$

 - gives an evaluation of all gates of the circuit



- Probability computation on a circuit C :

Given a probabilistic valuation $\pi: X \rightarrow [0, 1]$

$$\begin{aligned}\pi(x) &: 20\% \\ \pi(y) &: 70\% \\ \text{proba.} &: 76\%\end{aligned}$$

what is the probability that the circuit evaluate to 1, assuming independence?

→ This problem is intractable (SAT reduces to it)

- Deterministic decomposable circuit (d-D) :

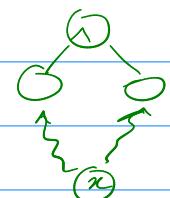
 - each AND gate is decomposable :



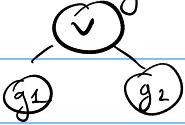
no variable has a path in the circuit to g_1 and g_2

→ the gate g_1 depends on x_1 and g_2 depends on x_2
with $x_1 \cap x_2 = \emptyset$

We do not have :



 - each OR gate is deterministic :



no valuation makes both g_1 and g_2 true

→ they are mutually exclusive

We do not have :



valuation π

Prop : We can perform probability computation on a d-D circuit in linear time (assuming constant-time arithmetics).

Proof : Compute the probability on each gate bottom-up

• Variable	• NOT	• AND	• OR
------------	-------	-------	------

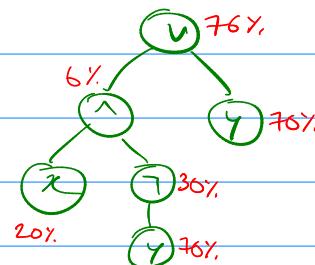
(x) $\pi(x)$

(g) $1 - \pi(g)$

(g_1) $\pi(g_1)$ (g_2) $\pi(g_2)$
independent

• OR

(g_1) $\pi(g_1)$ (g_2) $\pi(g_2)$
 $\pi(g_1) + \pi(g_2)$
mutually exclusive



Remark: does not work for general circuits, because of correlations



Provenance Circuits

• Provenance of a query Q on an instance I :

Boolean function $\text{Prov}(Q, I)$ having as variables the facts of I such that for any valuation $v: I \rightarrow \{0, 1\}$,

letting $I_v = \{F \in I \mid v(F) = 1\}$ be the subinstance of facts mapped to 1 we have $I_v \models Q$ iff $\text{Prov}(Q, I)$ evaluates to 1 under v

→ The provenance describes the subinstances of I that satisfy Q

• Provenance circuit of Q on I : Boolean circuit that expresses $\text{Prov}(Q, I)$

Thm: "We can perform PQE via provenance circuits":

Let $J = (I, \pi)$ be a TID, let Q be a query, let C be a provenance circuit of Q on I : the probability $\Pr_J(Q)$ is precisely the answer of probability evaluation of C with π

Proof: By definition

See Green,
Karrarachchi,
Tannen, PODS'03

Thm: For any fixed UCA Q , given an instance I , we can compute in PTIME a provenance circuit

Proof: write it in Disjunctive Normal Form (DNF)

with a disjunct for every image of a homomorphism of Q in I

This can be useful, e.g., to explain why the query holds but does not help for PQE unless the circuit is in a tractable class (here, d-Ds).

$I = \{S(a,b), S(c,c), S(a,a), T(b), T(c)\}$

$Q: S(x,y) \wedge T(y)$

homomorphism 1:

$x \mapsto a, y \mapsto b$
⇒ disjunct $S(a,b) \wedge T(b)$

homomorphism 2:

$x \mapsto c, y \mapsto c$
⇒ disjunct $S(c,c) \wedge T(c)$

Via Olteanu,
Huang, SUM'08

Thm: For any hierarchical SJF(Q) Q , given any instance I ,

we can compute in linear time a d-D provenance circuit of Q on I

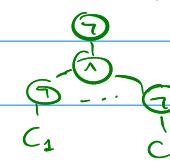
→ We get the tractability of PQE as a corollary.

Proof: Recall which rules we used:

- conjunction: $\Pr_a(Q_1 \wedge \dots \wedge Q_n)$ independent AND
- existential: $\Pr_a(\bigvee_a Q[\exists a])$ independent OR
- facts: variable gate



The subcircuits C_i depend on facts from different relations, so this n is decomposable



(Using de Morgan's law)

The subcircuits C_i depend on different facts (depending on the value of the root variable), so this n is decomposable

Note: We can get circuits in more restricted formalisms (read-once formulas, OBDDs...)

Tractable cases for UCQs: independent \wedge and \vee

We go from STFCQs to UCQs \rightarrow union in addition to conjunction
 \nwarrow repeated relation symbols

Def: A UCQ is an expression built with atoms, \wedge , \vee , and existential quantification
 $(\exists x R(x) \wedge \exists y T(y)) \vee \exists x y S(x,y)$

Standard way to write a UCQ: union of CQs:

$$(\exists x y \underline{R(x)} \wedge \underline{T(y)}) \vee \exists x y \underline{S(x,y)}$$

We decompose the CQs into connected components based on shared variables (called components)

We get $(C_{11} \wedge C_{12} \wedge \dots) \vee (C_{21} \wedge C_{22} \wedge \dots) \vee \dots$ where the C_{ij} are components (with quantifiers)

This is called the DNF representation of the UCQ

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

Far less standard representation: CNF representation of the UCQ

Obtained using de Morgan's law

We get $(C'_{11} \vee C'_{12} \vee \dots) \wedge (C'_{21} \vee C'_{22} \vee \dots) \wedge \dots$

$$(\underbrace{\exists x R(x) \vee \exists y S(x,y)}_{\text{---}}) \wedge (\underbrace{\exists y T(y) \vee \exists x S(x,y)}_{\text{---}})$$

For STFCQs we simplified \wedge as independent AND

\rightarrow no longer possible in general: repeated relation symbols

If it is possible, we do it:

Rule: If we can write Q as $Q' \wedge Q''$ where Q' and Q'' are Boolean queries that share no relation symbols, then $Pr(Q) = Pr(Q') \times Pr(Q'')$.

$$\text{Ex: } Pr((\exists x R(x) \vee \exists y T(y)) \wedge \exists x y S(x,y)) = Pr(\exists x R(x) \vee \exists y T(y)) \times Pr(\exists x y S(x,y))$$

Rule: If we can write Q as $Q' \vee Q''$ where Q' and Q'' are Boolean queries that share no relation symbols, then

$$Pr(Q) = 1 - (1 - Pr(Q'))(1 - Pr(Q''))$$

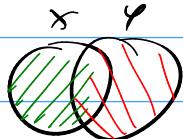
\rightarrow can be tested by putting Q in CNF, and in DNF, and comparing connected symbol-components.

Tractable case for UCQs : inclusion-exclusion

Say we have a UCQ where we cannot apply these two rules

Idea : use inclusion-exclusion principle on the CNF to get no unions of components
(and then see what we can do ...)

$$\text{Inclusion-exclusion (sets)} : |X \cup Y| = \underline{|X|} + \underline{|Y|} - \underline{|X \cap Y|}$$



$$\text{Inclusion-exclusion (probabilistic)} : \Pr(Q \vee Q') = \Pr(Q) + \Pr(Q') - \Pr(Q \wedge Q')$$

→ independence NOT required!

$$\text{Inclusion-exclusion (dual; for CNF)} : \Pr(Q \wedge Q') = \Pr(Q) + \Pr(Q') - \Pr(Q \vee Q').$$

$$\begin{aligned} \text{Example} : & \Pr((\exists x R(x) \vee \exists y S(y)) \wedge (\exists y S(y) \vee \exists z T(z))) \\ &= \Pr(\underbrace{\exists x R(x) \vee \exists y S(y)}_{\text{independent AND}}) + \Pr(\underbrace{\exists y S(y) \wedge \exists z T(z)}_{\text{independent AND}}) \\ &\quad - \Pr(\underbrace{\exists x R(x) \vee \exists y S(y) \vee \exists y S(y) \vee \exists z T(z)}_{\text{redundant (see later)}}) \\ &\quad \qquad \qquad \qquad \text{independent OR} \end{aligned}$$

OK, which case remains?

→ A union of components $c_1 \vee \dots \vee c_n$ which is not an independent OR

Idea : Remember hierarchical expressions!

If some c_i does not have a root variable, we fail (intractable)

Otherwise, $\exists x c_1 \vee \exists x c_2 \vee \dots \vee \exists x c_n \equiv \exists x [c_1 \vee c_2 \vee \dots \vee c_n]$

↑ now x occurs in all atoms!
↑ specific to union!

Idea : Can we write this as $\bigvee_a [c_1 \vee \dots \vee c_n] L(a)$?

Tractable cases for UCQs : separator variables

$$\begin{aligned} \text{Example: } & (\exists_{xy} R(x), S(x,y)) \vee (\exists_{xy} R'(x), S(x,y)) \\ & \equiv \exists x (\exists_y R(x) S(x,y) \vee \exists_y R'(x) S(x,y)) \\ & \equiv \bigvee_{a \in \text{domain}} \exists_y R(a) S(a,y) \vee \exists_y R'(a) S(a,y) \end{aligned}$$

Is this an independent OR? \Rightarrow YES, a database fact contributes to at most one disjunct
 So we can simply compute $\Pr(\exists_y \underline{R(a)} \underline{S(a,y)} \vee \exists_y \underline{R'(a)} \underline{S(a,y)})$

Components: only shared variable count

$$\begin{aligned} \text{This is } & \Pr((R(a) \vee R'(a)) \wedge (R(a) \wedge \cancel{\exists_y S(a,y)}) \wedge (R'(a) \wedge \cancel{\exists_y S(a,y)})) \wedge (\exists_y S(a,y) \vee \cancel{\exists_y S(a,y)}) \\ & = \Pr((R(a) \vee R(a')) \wedge \cancel{\exists_y S(a,y)}) \quad : \text{ok} \\ & \qquad \qquad \text{independent OR} \qquad \text{independent AND} \end{aligned}$$

So $(\exists_{xy} R(x), S(x,y)) \vee (\exists_{xy} R'(x), S(x,y))$ is tractable

Counter-example: $(\exists_{xy} R(x), S(x,y)) \vee (\exists_{xy} S(y,x), T(x))$

$$\equiv \exists x (\exists_y R(x) S(x,y)) \vee (\exists_y S(y,x), T(x)) \quad \text{instance}$$

$$\equiv \bigvee_{a \in \text{domain}} \exists_y R(a) S(a,y) \vee \exists_y S(y,a), T(a)$$

On the instance: $(\exists_y R(u) S(u,y) \vee \exists_y S(y,u), T(u)) \vee (\exists_y R(v) S(v,y) \vee \exists_y S(y,v), T(v))$

PROBLEM:
 disjoint atoms are not independent !!

In fact, this query is intractable... the order of variables matters!

Def: A separator variable for a union of components $c_1 \vee \dots \vee c_n$

is a root variable x such that we can write $\exists x c_1 \vee \dots \vee c_n$

(formally, we have chosen one root variable per component: x occurs in all atoms)

such that for every relation S there is at least one position i s.t.

x occurs at position i in every S -atom

\Rightarrow guarantees that $\bigvee_a Q[x/a]$ is really an independent OR

Tractable cases for UCQs: Summary of the algo so far

- If Q is a fact $Q = R(a_1, \dots, a_n)$ (atom with no variables) then
look up its probability [base]
- If you can write $Q = Q' \sqcap Q''$ with Q' and Q'' syntactically independent (Q in CNF), then $\Pr(Q) = \Pr(Q') \times \Pr(Q'')$ [Independent AND]
- If you can write $Q = Q' \vee Q''$ with Q' and Q'' syntactically independent (Q in DNF) then $\Pr(Q) = 1 - (1 - \Pr(Q'))(1 - \Pr(Q''))$ [Independent OR]
- If you have $Q = Q_1 \sqcap \dots \sqcap Q_n$ in CNF ($n > 1$) then

$$\Pr(Q) = \Pr(Q_1) + \Pr(Q_2) + \dots + \Pr(Q_n)$$
 - $\Pr(Q_1 \vee Q_2) = \Pr(Q_1 \vee Q_3) = \dots$
 - + ...
 - ...
- If you have $Q = c_1 \vee \dots \vee c_k$ as a union of components and we can write Q as $\exists x \ c_1 \vee \dots \vee c_k$ with x a separator variable then return $1 - \prod_{a \in \text{domain}} 1 - \Pr(Q[a/x])$ [Independent project]
- Else FAIL

Making the algorithm complete: 2 small fixes and one big fix

Small fix 1: minification

At every step, make the query minimal by removing redundant disjuncts/conjuncts
ex: $R(x) \vee R(x) S(xy) T(y)$

Small fix 2: ranking

Transformation on instance and database to make sure all atoms "use variables in the same order"

ex: $R(xy) R(yx)$ PTIME but no separator variable

→ Fix an order on the domain of instance I

→ Replace $R(a,a)$ by $R_=(a)$, for $a \neq b$ replace $R(a,b)$ by $R_<(a,b)$, for $a \geq b$ replace $R(a,b)$ by $R_>(a,b)$

→ Consider all possible rankings on the query variables $x = y, x < y, x \geq y$

→ Get disjuncts $R_<(xy) R_>(x,y) \vee R_=(x)$

Note: $R(x) S(xy) \vee S(xy) T(y)$ is ranked

Tractable cases for UCs : Cancellations in inclusion-exclusion

Big fix : In inclusion-exclusion, some terms may cancel out.

Recall that $\Pr(Q_1 \cup Q_2 \cup \dots \cup Q_n) = \Pr(Q_1) + \Pr(Q_2) + \dots + \Pr(Q_n)$

$$= \Pr(Q_1 \cup Q_2) - \dots$$
$$+ \dots$$

Some of these queries may end up being equivalent and cancel out.

ex : $(Q_1 \cup Q_3) \cap (Q_1 \cup Q_2) \cap (Q_1 \cup Q_3) \cap (Q_1 \cup Q_2 \cup Q_3)$

written $23 \cap 03 \cap 13 \cap 012$

Inclusion-exclusion: $\Pr(23 \cap 03 \cap 13 \cap 012) =$

$$\begin{aligned} & \Pr(23) + \Pr(03) + \Pr(13) + \Pr(012) \\ & - \Pr(023) - \Pr(123) - \Pr(0123) - \Pr(013) - \Pr(012) - \Pr(013) \\ & + \Pr(0123) + \Pr(0123) + \Pr(0123) + \Pr(0123) \\ & - \Pr(0123) \end{aligned}$$

It may be the case that $\Pr(0123)$ is intractable to compute!

Hence, fix: In inclusion-exclusion, regroup the equivalent queries
and recurse only on those with non-zero coefficient

→ Now the algorithm is complete

Tractable cases for UCQs : provenance circuits

Recall the case of SJFCQs :

Then for any hierarchical SJFCQ Q , given any instance I , we can compute in linear time a d-D provenance circuit of Q on I
 → We get the tractability of PQE as a corollary.

Is the same true for UCQs?

Recall the rules :

- independent AND : $Q = Q' \wedge Q''$ $\xrightarrow{\text{indep}}$ ok (decomposable AND)
- independent OR : $Q = Q' \vee Q''$ $\xrightarrow{\text{indep}}$ ok (decomposable AND + negation)
- independent project : $Q = \bigvee_a Q[x/a]$ $\xrightarrow{\text{indep}}$ ok (dito)
- inclusion-exclusion (with cancellations) $\Pr(Q_1 \wedge Q_2 \wedge \dots \wedge Q_n) = \Pr(Q_1) + \dots - \dots + \dots$

⇒ problem! This is not a logical operation but arithmetic on the probabilities

In simple cases, we can do it at a logical level:

$$Q_1 \wedge Q_2 \equiv \neg \left(\neg \left(Q_1 \vee \neg \left(Q_1 \wedge Q_2 \right) \right) \vee \neg Q_2 \right)$$

\uparrow mutually exclusive
 \uparrow mutually exclusive

→ expresses $\Pr(Q_1 \wedge Q_2)$ as a function of $\Pr(Q_1)$, $\Pr(Q_2)$, $\Pr(Q_1 \vee Q_2)$

$$\begin{aligned} \Pr(Q_1 \wedge Q_2) &= 1 - (1 - (\Pr(Q_1) + 1 - \Pr(Q_1 \vee Q_2)) + 1 - \Pr(Q_2)) \\ &= \Pr(Q_1) + 1 - \Pr(Q_1 \vee Q_2) - 1 + \Pr(Q_2) \\ &= \Pr(Q_1) + \Pr(Q_2) - \Pr(Q_1 \vee Q_2) \end{aligned}$$

We can apply this recursively, but does not handle cancellations

⇒ Open problem!

See Monet,
PODS'20

Intensional-extensional conjecture: There are UCQs Q such that $\text{PQE}(Q)$ is tractable

but we cannot build a d-D provenance circuit for Q on an input instance in PTIME

(I hope the conjecture can be disproved)

Remark: known to be true for weaker formalisms than d-Ds :

- FBDDs and dec-DNFs and DLDDs
- Structured d-DNFs

Also open for d-DNNFs (d-Ds with negations at leaves), open if $d\text{-D} \equiv d\text{-DNFs}$
 Some characterizations known for weaker formalisms, e.g., OBDDs (inversion-free UCQs)

Beame, Li, Ray,
Suri, PODS'17
Bova, Sjeider,
PODS'17

Stand-alone rephrasing of the open problem

Take a join-semilattice (V, \leq)

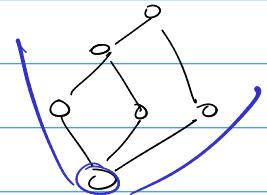
Define the Möbius function for $v \in V$

$$\mu(v) = 1 - \sum_{w < v} \mu(w)$$

Let $S \subseteq V$ be the set of vertices v having $\mu(v) \neq 0$

Inductively define the reachable subsets:

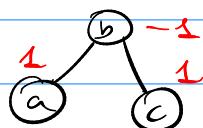
- For each $v \in S$, its upset $\{w \in V \mid v \leq w\}$ is reachable
- If S_1 and S_2 are reachable and disjoint then $S_1 \sqcup S_2$ is reachable
- If S_1 and S_2 are reachable and $S_2 \subseteq S_1$ then $S_1 \setminus S_2$ is reachable



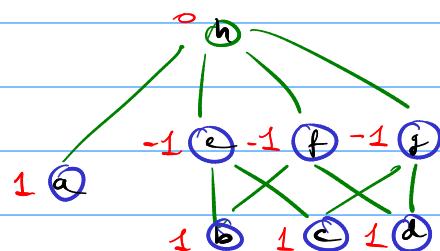
The join-semilattice (V, \leq) is good if V is reachable.

\Rightarrow Are all join-semilattices good?

ex:

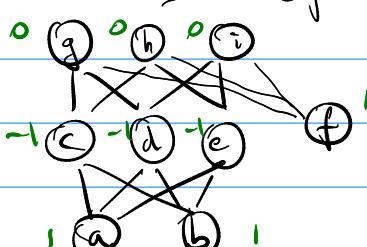


$$\begin{aligned} &\{\{a, b\}, \{\{b\}, \{c, b\}\} \\ &\{\{a\}\} = \{\{a, b\} \setminus \{\{b\}\}\} \\ &\{\{a, b, c\}\} = \{\{a\}\} \cup \{\{c, b\}\} \end{aligned}$$



$$\hat{b} - \hat{e} + \hat{c} - \hat{f} + \hat{a} - \hat{g} + \hat{d} = V$$

False on DAGs: (found by brute-force search)



Lower bounds for non-hierarchical SJFCQs : basics

- **#P** : class of counting problems expressible as the number of accepting paths of a polynomial-time nondeterministic Turing machine

#SAT : count the satisfying assignments

→ A counting analogue of NP : NP asks "is the count ≥ 1 ?"

→ but there are PTIME problems whose counting analogue is #P-hard (see later)

Remark : If (non-probabilistic) query evaluation for a query Q is in PTIME, then $PQE(Q)$ is in #P

Proof : The Turing machine first guesses a valuation (need to bias by the probabilities)
then evaluates the query on the possible world

(Technically in FP^{#P} because the output of PQE is a fraction and we must normalize by the denominator)

Provost, Ball,
SICOMP'83
• #P-hard problem: every problem in #P reduces to it by a PTIME (Turing) reduction

Our initial #P-hard problem : the #PP2DNF problem aka counting independent sets in bipartite graphs

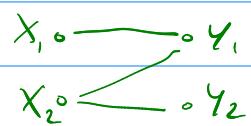
Input :

- Variables partitioned into two sets $X = \{X_1, \dots, X_n\}$ and $Y = \{Y_1, \dots, Y_m\}$
- Clauses $X_{i_1} \wedge Y_{j_1} \vee \dots \vee X_{i_k} \wedge Y_{j_k}$
 - no negation
 - one X -variable and one Y -variable per clause

Output : How many satisfying assignments for the formula $(X_1 \wedge Y_1) \vee \dots \vee (X_n \wedge Y_m)$?

Example : formula $(X_1 \wedge Y_1) \vee (X_2 \wedge Y_1) \vee (X_2 \wedge Y_2)$

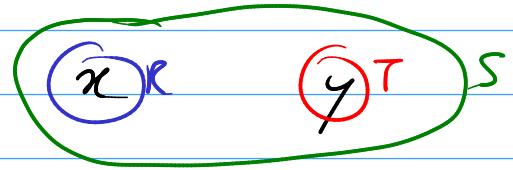
→ 8 satisfying assignments out of 16



Lower bounds for non-hierarchical SJFCQs: R-S-T queries

Let's first show hardness for the simple non-hierarchical SJFCQ:

$Q: R(x), S(x,y), T(y)$



Theorem: PQE for the SJFCQ $Q; R(x) S(x,y) T(y)$ is #P-hard

Proof: consider a #PPDNF instance $\{x_1, \dots, x_n\}, \{y_1, \dots, y_m\}, x_i \wedge y_j, \dots, x_k \wedge y_k$

Create the TID:

$$\begin{array}{ccc} \frac{1}{2} R(a_1) & \xrightarrow{\quad} & T(b_1) \frac{1}{2} \\ : & \cancel{\xrightarrow{\quad}} & : \\ \frac{1}{2} R(a_n) & & T(b_m) \frac{1}{2} \end{array}$$

plus facts $S(a_{ip}, b_{jp})$ with prob 1 for $1 \leq p \leq k$

→ there is a probability-preserving bijection
between possible worlds and valuations

→ a valuation satisfies the formula

(x_{ip} and y_{ip} both true for some $1 \leq p \leq k$)

iff the possible world satisfies Q

(contains acts $R(a_{ip}), S(a_{ip}, b_{jp}), T(b_{jp})$ for some $1 \leq p \leq k$)

Let us extend to other simple non-hierarchical SJFCQs:

Def: An R-S-T query is a query of the form:

$Q_{rst}(x,y) : R_r(x), \dots, R_r(x), S_s(x,y), \dots, S_s(x,y), T_t(y), \dots, T_t(y)$

for some $r,s,t > 0$.

Corollary: PQE(Q) is #P-hard for any R-S-T query

Proof: Same reduction as before, add facts $R_2(a), \dots, R_r(a)$

$S_2(a,b), \dots, S_s(a,b)$

$T_2(b), \dots, T_t(b)$

for all elements a, b , with prob 1.

Lower bounds for non-hierarchical SJFCQs: hard patterns

We want to show NP-hardness for all non-hierarchical SJFCQs

Lemma: A SJFCQ Q is hierarchical iff the following holds:

Letting $\text{atoms}(x)$ for a variable x of Q be
the set of atoms that contain x ,

For any two variables x and y ,
the sets $\text{atoms}(x)$ and $\text{atoms}(y)$ are either:

- disjoint : $\text{atoms}(x) \cap \text{atoms}(y) = \emptyset$
- comparable : $\text{atoms}(x) \subseteq \text{atoms}(y)$ or $\text{atoms}(y) \subseteq \text{atoms}(x)$

Q : $R(x), S(x,y), T(z)$

$$\text{atoms}(x) = \{R(x), S(x,y)\}$$

$$\text{atoms}(y) = \{S(x,y)\}$$

$$\text{atoms}(z) = \{T(z)\}$$

Q : $R(x), S(x,y), T(y)$

$$\text{atoms}(x) = \{R(x), S(x,y)\}$$

$$\text{atoms}(y) = \{T(y), S(x,y)\}$$

intersect but not comparable

Proof :

- hierarchical \Rightarrow condition, by induction

- One atom A : true, $\text{atoms}(x) = \{A\}$ for all vars

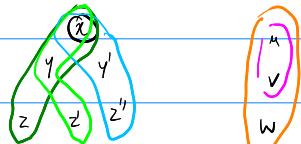
- $\exists x Q(x, y_1, \dots, y_n)$: true, atoms is unchanged

- $Q_1 \wedge Q_2$: if a variable x occurs both in Q_1 and Q_2 , then it must occur in all atoms
(because it will be quantified later)

so $\text{atoms}(x)$ and $\text{atoms}(y)$ cannot be a counterexample for such a variable x

- for x occurring only in Q_1 and y occurring only in Q_2 , $\text{atoms}(x) \cap \text{atoms}(y) = \emptyset$

- for x and y occurring only in Q_1 , conclude by induction hypothesis
(ditto for x and y only in Q_2)



- condition \Rightarrow hierarchical, also by induction

Take the variables x_1, \dots, x_n such that

$\text{atoms}(x_1), \dots, \text{atoms}(x_n)$ are maximal under inclusion

We must have $\text{atoms}(x_i) \cap \text{atoms}(x_j) = \emptyset$

or $\text{atoms}(x_i) = \text{atoms}(x_j)$ for all i, j

(otherwise, by the condition, they are comparable,
and different, contradicting maximality)

Write each atom with no variables as a conjunct

The maximal sets cover all other atoms:

each distinct set is a conjunct where

you quantify the variables achieving the set, and reverse

$$x: A_1 A_2 A_3 A_4 \quad u: A_g$$

$$y: A_1 A_2 A_3 A_4 \quad v: A_g$$

$$z: A_1 A_2 \quad w: A_3$$

$$z': A_1$$

$$\exists x \exists y (A_4 \wedge (\exists z A_2 \wedge (\exists z' A_1)) \wedge (\exists w A_3)) \wedge (\exists u \exists v A_g)$$

Corollary: If Q is a non-hierarchical SJFCQ, it has two variables x, y such that

- $\text{atoms}(x) \cap \text{atoms}(y) \neq \emptyset$ So there are atoms with x and y Call these atoms $R(x)$

- $\text{atoms}(x) \not\subseteq \text{atoms}(y)$

- $\text{atoms}(y) \not\subseteq \text{atoms}(x)$

- atoms with x but not y

- atoms with y but not x

$$R(x,y)$$

$$R(y)$$

(non-empty)

Lower bound for non-hierarchical SJFCQ : reducing from R-S-T

Theorem: Let \mathcal{Q} be a SJFCQ with nonempty atom sets

$R(x)$ atoms using x but not y

$S(x,y)$ atoms using x and y

$T(y)$ atoms using y but not x

for some variables x,y . Let $r = |R|$, $s = |S|$, $t = |T|$.

Then $\text{PQE}(\mathcal{Q}_{\text{rst}})$ reduces in PTIME to $\text{PQE}(\mathcal{Q})$

Proof: Take a TID J on which we want to solve $\text{PQE}(\mathcal{Q}_{\text{rst}})$.

Let c be a new constant.

Build J' from J by changing the facts:

- For the p -th atom $R(x, z_1, \dots, z_n)$ of R
replace each fact $R_p(a)$ by $R(a, c, \dots, c)$
keeping the same probability
- Do similarly for R, S, T
- For each atom in $Z = Q \setminus (R \cup S \cup T)$
letting Z be the relation, add a fact $Z(c)$ with proba 1.

$$\begin{array}{ccccccccc} R(x), R'(x, w, x), S(x, w, y) & T(w, y) & U(w, w) \\ \overbrace{R}^{\frac{1}{r}} & \overbrace{S}^{\frac{1}{s}} & \overbrace{T}^{\frac{1}{t}} & \overbrace{U}^{\frac{1}{s}} \end{array}$$

$R(a), R'_1(a), S(a, b), T_1(b)$
becomes
 $R(a), R'_1(a, c, a), S(a, c, a), T_1(c, b)$
 $U(c, c)$

Correctness: there is a probability-preserving bijection between the possible worlds of J and J'

- if a possible world of J satisfies \mathcal{Q}_{rst} then the translation of any match witnesses that the corresponding world of J' satisfies \mathcal{Q}
- if a possible world of J' satisfies \mathcal{Q} , restricting to $R \cup S \cup T$ and taking the preimage of a match, the corresponding world of J satisfies \mathcal{Q}

Dalvi, Suciu,
VLDB'07

Conclusion: For any SJFCQ \mathcal{Q} :

- either \mathcal{Q} is hierarchical and we can solve $\text{PQE}(\mathcal{Q})$ in PTIME on any instance I by building a d-D circuit representing $\text{Prob}(\mathcal{Q}, I)$
- or $\text{PQE}(\mathcal{Q})$ is #P-hard

#P-hardness of Uniform Reliability for Non-hierarchical SJFCQ

Remember Uniform Reliability (UR): all facts have probability $\frac{1}{2}$

We know it is tractable for hierarchical SJFCQs (like PQE)

But is it also #P-hard for non-hierarchical SJFCQs?

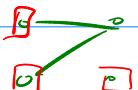
(Reduction uses probs $\frac{1}{2}, 0$, and 1)

Theorem: UR(Q) is #P-hard for any non-hierarchical SJFCQ Q

- Also works for other fixed probabilities than $\frac{1}{2}$ (if $0 < p < 1$)
- Example use case of the interpolation method used in many hardness proofs

By the previous techniques it essentially suffices to show:

Theorem: UR(Q) is #P-hard for R(a) S(x,y) T(y)



We will reduce from counting the independent sets of a bipartite graph (\leftarrow #PPLEAF)

Let $G = (R \cup T, S)$ be the input and X be the set of independent sets: we want $|X|$

Principle of the interpolation method:

1. Partition the subsets of RUT based on some properties

$$2^{RUT} = X_0 \sqcup X_1 \sqcup \dots \sqcup X_n$$

where $X = X_{i_1} \sqcup \dots \sqcup X_{i_k}$

2. Code G into TIDs with gadgets having a degree of freedom

$$J_0, J_1, \dots, J_n$$

3. Relate the answer of the target problem, $\Pr_{\mathcal{I}_d}(Q)$ for each d,

to the cardinalities $|X_0|, \dots, |X_n|$ as a linear equation

4. Show invertibility of the matrix of the system

NP-hardness of UR: 1. Partitioning

Consider the bipartite graph $G = (R \cup T, S)$

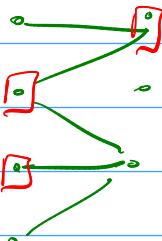
and subsets of vertices $R' \subseteq R$, $T' \subseteq T$

The "important quantities" of R' and T' are:

- number of edges $\rightarrow c_{\infty}$

- number of edges $\square \text{ or } \square \rightarrow c_{\infty}$

- number of edges $\square \rightarrow c_{11}$



$$c_{\infty} = 1$$

$$c_{11} = 3$$

$$c_{\infty} = 1$$

Note that $c_{\infty} + c_{11} + c_{\infty} = |S|$ (number of edges), so 2 "degrees of freedom"

Write X_S the set of vertex subsets with quantities as indicated

($|X_S| = 0$ if $c_{\infty} + c_{11} + c_{\infty} \neq |S|$)

We have: $2^{|R \cup T|} = \bigcup_S X_S$

The independent sets are precisely in the X_S where $c_{11} = 0$

$$X = \bigsqcup_{c_{\infty}, c_{11}} X_{c_{\infty}, c_{11}, 0}$$

$$|X| = \sum_{c_{\infty}, c_{11}} |X_{c_{\infty}, c_{11}, 0}|$$

So let us show how to compute all the X_S

#P-hardness of UR : 2-Coding

It is easy to code the bipartite graph  to a PDE instance $R(a_1) \xrightarrow{s} T(b_1)$
 $R(a_2) \xrightarrow{s} T(b_2)$

but all facts now have probability $1/2$... what can vary?

Idea : Code each edge by p copies of a "back-and-forth path"

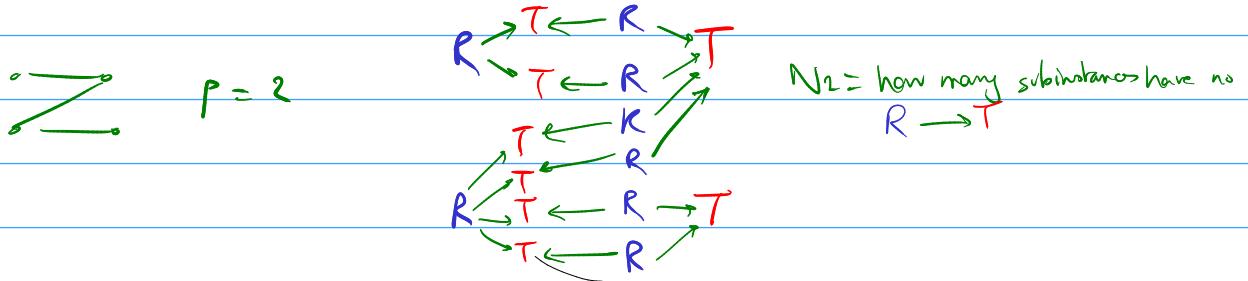


For a choice of p , call T_p the result of coding the input G in this way

#P-hardness of UR: 3. Relating

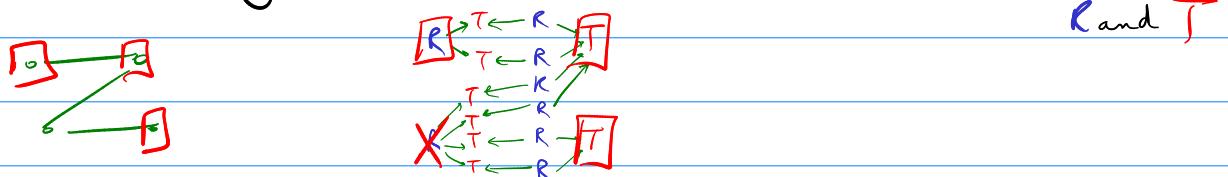
$2^{|J_p|}$ - oracle answer

Given a TID J_p , our UR oracle tells us how many subinstances violate Q: N_p



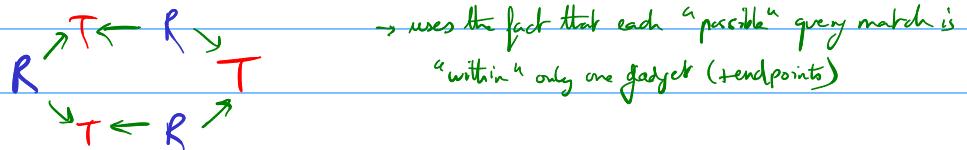
How to relate N_p and the $|X_2|$?

- In J_p , when choosing a possible world, we must decide if we keep/discard the endpoints R and T



$$N_p = \sum_{\substack{R' \subseteq R \\ T' \subseteq T}} N_p(R', T')$$

- Now, $N_p(R', T')$ is the product of the number of violating possible worlds for each gadget



- Each such number depends only on p (here $p=2$) and whether the endpoints are kept or not

(4 possible types: $\square = \square$ $\square \times \square$ $\times = \square$ $\times = \times$)

Call $N_{p,bb}$ this number, for $b, b' \in \{0, 1\}$ the type

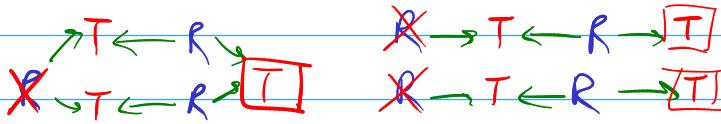
by symmetry $N_{p,00} = N_{p,11}$

- How many gadgets of each type we have is given by the "important quantities" $c_{00}(R', T')$, $c_{01}(R', T')$, $c_{10}(R', T')$

$$N_p(R', T') = N_{p,00}^{c_{00}} \times N_{p,01}^{c_{01}} \times N_{p,11}^{c_{11}}$$

#P-hardness of UR : 3. Relating (cont'd)

- Let's compute $N_{p_{bb'}}$. First note that $N_{p_{bb'}} = (N_{bb'})^P$



- Let's compute!

$$N_{00} \quad \cancel{X} \rightarrow T \leftarrow R \rightarrow \cancel{X} \quad \frac{2 \times 2 \times (8-1)}{2} = 28 = 4 \times 7$$

$$N_{01} \quad \cancel{X} \rightarrow T \leftarrow R \rightarrow \boxed{T} \quad \frac{2 \times (8+3)}{R \text{ missing}} = 22 = 2 \times 11$$

(or $\boxed{R} \rightarrow T \leftarrow R \rightarrow \cancel{X}$ — symmetry)

$$N_{11} \quad \boxed{R} \rightarrow T \leftarrow R \rightarrow \boxed{T} \quad \begin{matrix} 3 \times 3 \\ \text{middle } S \\ \text{missing} \end{matrix} + \begin{matrix} 4 \\ \text{both } R \\ \text{and } T \text{ missing} \end{matrix} + \begin{matrix} 2 \times 2 \\ \text{one } R \\ \text{present} \\ \text{one } T \\ \text{missing} \end{matrix} = 17$$

Now let's put the equations together

$$\begin{aligned} N_p &= \sum_{\substack{R \in R \\ T \in T}} N_p(R, T) = \sum_{\substack{R \in R \\ T \in T}} N_{p_{00}}^{c_{00}} \times N_{p_{01}}^{c_{01}} \times N_{p_{11}}^{c_{11}} \\ &= \sum_{\substack{R \in R \\ T \in T}} \left(N_{p_{00}}^{c_{00}} \times N_{p_{01}}^{c_{01}} \times N_{p_{11}}^{c_{11}} \right)^P \end{aligned}$$

The sum argument only depends on the $c_{bb'}$ \rightarrow inject the $|X_{c_{bb'}}|$!

$$N_p = \sum_{c_{00}, c_{01}, c_{11}} |X_{c_{00}c_{01}c_{11}}| \times \left(N_{p_{00}}^{c_{00}} \times N_{p_{01}}^{c_{01}} \times N_{p_{11}}^{c_{11}} \right)^P$$

This is in fact the equation system:

$$\begin{pmatrix} N_0 \\ N_1 \\ N_2 \\ \vdots \\ N_B \end{pmatrix} = \begin{pmatrix} \alpha(0,0,0)^0 & \dots & \alpha(|S|,|S|,|S|)^0 \\ \vdots & \ddots & \vdots \\ \alpha(0,0,0)^B & & \alpha(|S|,|S|,|S|)^B \end{pmatrix} \times \begin{pmatrix} |X_{000}| \\ |X_{001}| \\ |X_{011}| \\ \vdots \\ |X_{111}| \end{pmatrix}$$

where $B = (|S|+1)^3 - 1$

and

$$\alpha(ijk) = N_{p_{00}}^i N_{p_{01}}^j N_{p_{11}}^k$$

#P-hardness of UR : 4. Showing invertibility

The matrix M has coefficients $M_{p,(i,j,k)} = \alpha(ijk)^p$

We recognize a Vandermonde matrix on values $\alpha(ijk)$

Lemma : A Vandermonde matrix is invertible iff its "base values" are all different

$$\begin{pmatrix} v_0^0 & v_1^0 \\ v_0^1 & v_1^1 \end{pmatrix}$$

$$v_0 \neq v_1$$

Proof : If we could write a linear combination of the rows $\lambda_0 R_0 + \dots + \lambda_n R_n = 0$
 then, looking at each column, each base value $(\lambda_0, \dots, \lambda_{n-1}) \neq (0, \dots, 0)$
 is a different root of the polynomial $\lambda_0 x^0 + \dots + \lambda_{n-1} x^{n-1}$ of degree $n-1$
 \Rightarrow impossible, a non-zero polynomial of degree $n-1$ has at most $n-1$ roots

\Rightarrow are all the $\alpha(ijk)$ values distinct ??

Let's assume $\alpha(ijk) = \alpha(i'j'k')$

$$\text{so } N_\infty^{i'} \times N_0^{j'} \times N_1^{k'} = N_\infty^{i''} \times N_0^{j''} \times N_1^{k''}$$

$$\text{and } 4^{i''} 7^{i'} \times 2^{j''} 11^{j'} \times 17^{k''} = 4^{i''} 7^{i'} \times 2^{j''} 11^{j'} \times 17^{k'}$$

Now by uniqueness of the prime number decomposition we get $(ijk) = (i'j'k')$

\rightarrow so the matrix M of the system is invertible

\rightarrow so we can use the equation system to recover

the desired values $|X_{ijk}|$ from the oracle results N_p
 (invert the matrix in PTIME)

\rightarrow and we get $|X| = \sum_{ij} |X_{ij}|$ showing #P-hardness

Note : here we used symmetry and got lucky with the numbers

General proof : more elaborate asymmetric coding plus rounding argument plus
 lemma showing that the "determinant" $N_{00}N_{11} - N_{01}N_{10}$ is $\neq 0$.

#P-hardness for UCQs where the algorithm fails

Recall that we had a complicated algorithm to compute the probability of UCQs
→ which failed when it could not find a separator variable

Thm: If the algorithm fails on a query \mathcal{Q} then $\text{PQE}(\mathcal{Q})$ is #P-hard.

Very technical proof. High level ideas:

- Focus on queries on which the algorithm fails immediately : immediately unsafe query
If the algorithm eventually fails on a query \mathcal{Q} , then there is a query \mathcal{Q}' on which it fails immediately s.t. $\text{PQE}(\mathcal{Q}')$ reduces to $\text{PQE}(\mathcal{Q})$.
- Rewrite immediately unsafe queries to make them leveled:
allows us to assume that the domain is k -partite and each atom $R(x_1, \dots, x_k)$ maps variable x_i to elements of one part only
→ show that the leveling is still immediately unsafe
→ we can assume that queries are leveled
- Define forbidden queries, the minimal immediately unsafe leveled queries, under some order
(sequence of relation arities, and number of atoms)
→ The simpler queries to which a forbidden query “rewrites” are always safe by minimality
→ Show that forbidden queries have been simplified to have only two levels,
ie, any leveled immediately unsafe query can be simplified further
→ A forbidden query is a disjunction of components with atoms $R(x), S(x,y), T(y)$ on a partitioned variable set $X \cup Y$
- Show that forbidden queries are #P-hard : different types (1-1, 1-2, 2-1, 22)
→ Very technical; uses the interpolation method

What about uniform reliability:

Thm: If $\text{PQE}(\mathcal{Q})$ is #P-hard for a UCQ \mathcal{Q} , this still holds with probas $\{0, \frac{1}{2}, 1\}$

Further, $\text{UR}(\mathcal{Q})$ is #P-hard for final type-1 queries.

Open problem: Is $\text{UR}(\mathcal{Q})$ #P-hard for every UCQ \mathcal{Q} where $\text{PQE}(\mathcal{Q})$ is #P-hard?

#P-hardness of PQE for unbounded UCQ^∞

We now extend our study of PQE to queries beyond UCQ_S .

We keep a key property: closure under homomorphisms

Def: A query Q is *closed under homomorphisms* if for any instances I and I' ,
if $I \models Q$ and I has a homomorphism to I' ($I \xrightarrow{\text{hom}} I'$)
then $I' \models Q$.

In particular: UCQ_S , but also languages like Datalog or regular path queries (RPQs)

Def: A regular path query Q on an arity-2 signature σ is a regex over σ

An instance I satisfies Q if there is a "path" of facts

$R_1(a_1, a_2) R_2(a_2, a_3) \dots R_n(a_n, a_{n+1})$ such that the word $R_1 \dots R_n$

is in the language of the regex Q

We write UCQ^* the queries closed under homomorphism.

Why? Because they are like UCQ_S with infinitely many disjuncts

$\text{RS}^*T : R(x_1 x_2) T(x_2 x_3) \vee R(x_1 x_2) S(x_2 x_3) T(x_3 x_4) \vee R(x_1 x_2) S(x_2 x_3) S(x_3 x_4) T(x_4 x_5)$

Def: A UCQ^∞ is *unbounded* if it is not expressible as a UCQ

(infinitely many non-redundant disjuncts)

ex: RS^*T

Rem: A UCQ^∞ is unbounded iff it has infinitely many minimal models

i.e. $I \models Q$ but for any $I' \not\models I$, $I' \not\models Q$

For technical reasons we restrict to arity-2 signatures (probabilistic graphs)

Work in progress: lifting the restriction

Theorem: For any unbounded $\text{UCQ}^\infty Q$ on an arity-2 signature, $\text{PQE}(Q)$ is #P-hard.

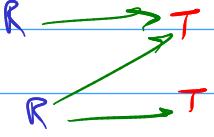
Work in progress: what about UR?

Proof structure: first look at bipartite R-S-T instances then arbitrary instances

$\#P$ -hardness for unbounded UCQ $^\infty$: bipartite R-S-T instances

Def: A bipartite RST instance I consists of:

- a bipartite domain $\text{dom}(I) = U \sqcup V$,
- R -facts on a subset $U' \subseteq U$ and T -facts on a subset $V' \subseteq V$
- S -facts on a subset $E \subseteq U \times V$



Def: A bipartite RST TID is a TID whose underlying instance is bipartite RST

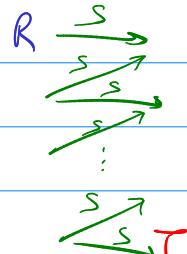
We will reduce from hard UCQ $^\infty$ (bounded or unbounded) on bipartite RST instances

$$Q_1: R(x) S(xy) T(y) \quad \#P\text{-hard}$$

$$Q_2: R(x) S(xy) S(x'y) S(x'y') T(y')$$

⋮

$$Q_n: R(x_1) S(x_1y_1) S(x_2y_1) S(x_2y_2) \dots S(x_ny_n) T(y_n)$$



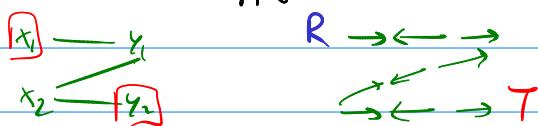
Lemma: For all $n > 0$, PQE(Q_n) on bipartite RST TIDs is $\#P$ -hard

Proof: Usual hardness proof but using paths $\rightarrow \leftarrow \rightarrow \leftarrow \rightarrow \dots \rightarrow$

Example for Q_2 : reduce from #PP2DNF



Notice that non-satisfying valuations can create matches of Q_m but $m > n$

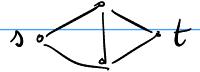


#P-hardness for unbounded \mathbb{Q}^∞ : bipartite R-S-T instances (contd)

Now consider the unbounded query $\mathbb{Q}_\infty = \bigvee_{n>0} \mathbb{Q}_n$. $R \xrightarrow{\leftarrow \rightarrow} T$

Lemma PQE(\mathbb{Q}_∞) on bipartite RST TIRs is #P-hard

Note: cannot do the usual reduction! (because of the \mathbb{Q}_m matches with $m > n$)



Def: #U-ST-CONN Undirected source-to-target connectivity problem

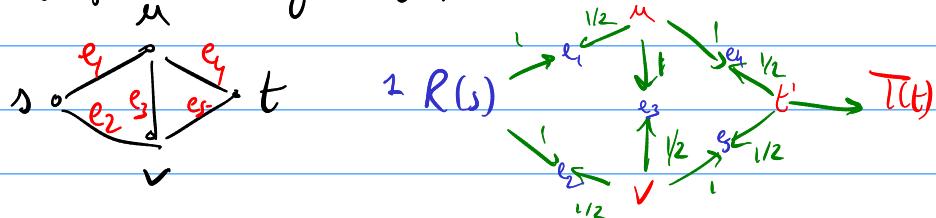
Input: undirected graph G , distinguished vertices s, t

Output: probability that s and t are connected (same connected component)
when all edges have proba $1/2$

Valiant
SICOMP'95

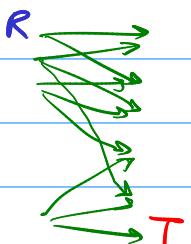
Prop: U-ST-CONN is #P-hard

Proof of lemma: given graph G code it to J :



Note: $\xrightarrow{1 \quad 1/2}$ or $\xrightarrow{1/2 \quad 1}$ is arbitrary

- one extra edge $R \rightarrow T$
- is bipartite



Correctness: Probability-preserving bijection between possible worlds

\Rightarrow if a possible world connects s and t then following any path gives a query match

\Rightarrow likewise, any match of \mathbb{Q}_∞ in a possible world of J gives a path

Now let's move to arbitrary graphs, reduce from bipartite RST TIRs.

$\#P$ -hardness for unbounded UCL^∞ : RST-patterns

We consider instances over a signature σ of arity 2 (edge colors).

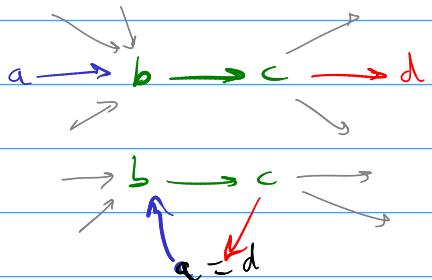
Edge orientation ignored for simplicity. Arity-1 facts also ignored.

Def: A RST-pattern in an instance I
is a choice of 3 facts

$$U(a,b)$$

$$V(b,c)$$

$$W(c,d)$$



Further, the 4 elements are distinct except maybe $a=d$.

We can use RST patterns to code an RST bipartite instance:

Def: Let I' be a bipartite RST instance

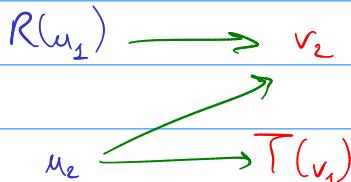
$\{u_1, \dots, u_n\}$ left vertices

$\{v_1, \dots, v_m\}$ right vertices

R-facts $R(u_1), \dots, R(u_n)$

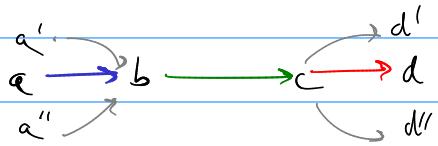
T-facts $T(v_1), \dots, T(v_m)$

S-facts $S(u_i, v_j)$ for $(u_i, v_j) \in E$



Let I be a graph instance

and a, b, c, d be an RST-pattern



The coding of I' in I is the graph instance

where we copy b to b_1, \dots, b_n

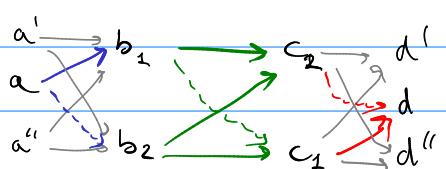
copy c to c_1, \dots, c_m

copy $U(a, b)$ to $U(a, b_1), \dots, U(a, b_n)$

copy $V(b, c)$ according to E

copy $W(c, d)$ to $W(c_1, d), \dots, W(c_m, d)$

copy all other facts "everywhere"



note: dashes --- means
"all facts except $V(b, c)$ "

To code a bipartite RST T_{10} , we give the proba of each $R(u_p)$ to $U(a, b_p)$
of each $T(v_q)$ to $W(c_q, d)$
of each $S(u_i, v_j)$ to its V -fact
and all other probas are 1.

Remark: If $I' \xrightarrow{\text{hom}} I''$ (on bipartite RST instances) [preserves homomorphisms]
then coding $(I', I) \xrightarrow{\text{hom}}$ coding (I'', I)

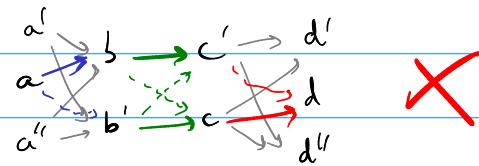
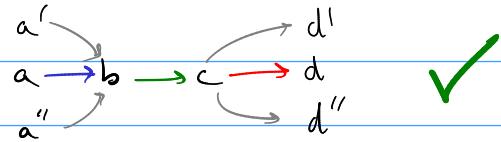
#P-hardness for unbounded UCQ[∞]: RST patterns (cont'd)

Consider a UCQ[∞] \mathcal{Q} . Fix an instance I and RST pattern $U(a,b) V(b,c) W(c,d)$

Def: The RST pattern is critical if:

- $I \models \mathcal{Q}$
- coding $(I', I) \not\models \mathcal{Q}$ for I' the bipartite RST instance

$$\begin{array}{c} R \\ \longrightarrow \\ \longrightarrow T \end{array}$$



Intuition: "the query cares about connecting \rightarrow and \rightarrow by $(\rightarrow \leftarrow)^{\pm}$ "
Two results to show

Thm: [Hardness of critical patterns]

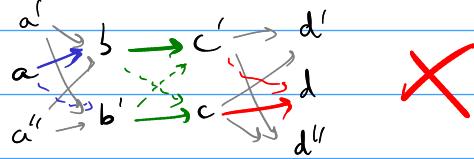
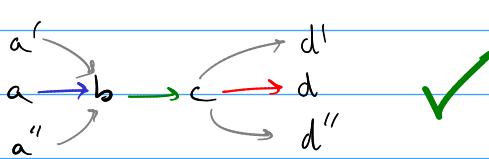
If we have a critical pattern for \mathcal{Q} then $PQE(\mathcal{Q})$ is #P-hard

Thm: [Finding a critical pattern]

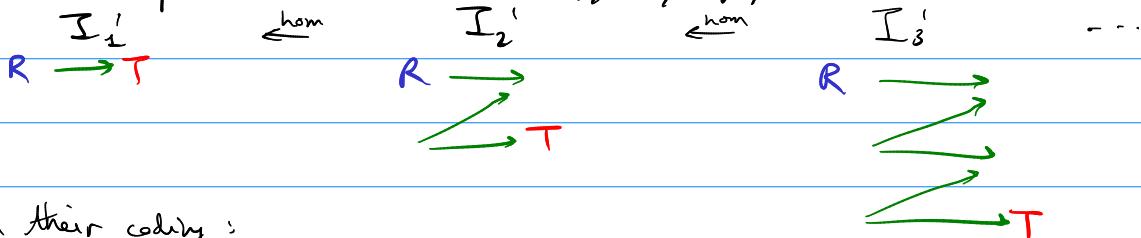
Every unbounded UCQ[∞] on graphs has a critical pattern.

$\#P$ -hardness for unbounded UCQ $^\infty$: hardness of critical patterns

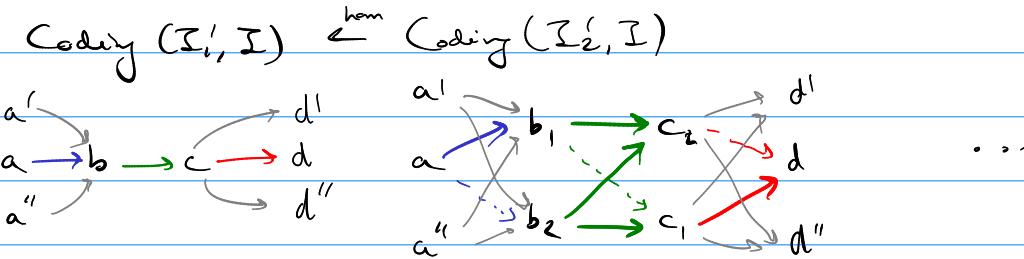
Fix the query Q , instance I , and critical pattern:



Consider the bipartite RST instances I'_1, \dots, I'_n, \dots



and their coding:



Two cases:

1. Coding (I'_i, I) satisfies the query ✓

Reduce from PQT(Q_{as}) for the connectivity query Q_{as} on a bipartite RST TID I' :

- A possible world of I' violating Q_{as} has a homomorphism to $R \xrightarrow{\quad} T$
so the world in the coding violates Q
- A possible world of I' satisfying Q_{as} has a homomorphism from some I'_i
so the world in the coding satisfies Q

2. There is a threshold i_0 , Coding (I'_i, I) satisfies the query for $i \leq i_0$ ✓
Coding (I'_i, I) violates the query for $i > i_0$ X

Reduce from PQE(Q_{i_0}) for the distance- i_0 query Q_{i_0} on a bipartite RST TID I' :

- A possible world of I' violating Q_{i_0} has a homomorphism to $R \xrightarrow{\quad} T$
 $\left. \begin{array}{c} \text{---} \\ \vdots \\ \text{---} \end{array} \right\} i_0+1$
- A possible world of I' satisfying Q_{i_0} has a homomorphism from $R \xrightarrow{\quad} T$
 $\left. \begin{array}{c} \text{---} \\ \vdots \\ \text{---} \end{array} \right\} i_0$

#P-hardness for unbounded UCQ[∞]: finding a critical pattern

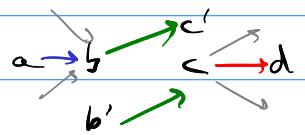
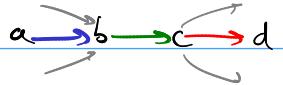
Fix an unbounded UCQ_{as} \mathcal{Q} and instance I with $I \vdash \mathcal{Q}$

First find a tight pattern then get a critical pattern

Def: The dissociation of an RST pattern

$U(a,b) V(bc) W(cd)$ in I is obtained by:

- copying b to b' and c to c'
- moving all facts $\{b,c\}$ including $S(b,c)$ to $\{b,c'\}$ and $\{b',c\}$



Note: this is not a coding!

Note: the resulting instance has a homomorphism back to I

Def: An RST pattern $U(a,b) V(bc) W(cd)$ in I is tight

if $I \vdash \mathcal{Q}$ but the dissociation of the pattern violates \mathcal{Q} .

Two remaining claims:

Thm: [Finding a tight pattern]

Any unbounded UCQ[∞] on graphs has a tight pattern

(This is where we really use the unboundedness of the query.)

Thm: [From tight to critical]

If a query has a tight pattern then it has a critical pattern

#P-hardness for unbounded UCQ^{oo}: finding a tight pattern

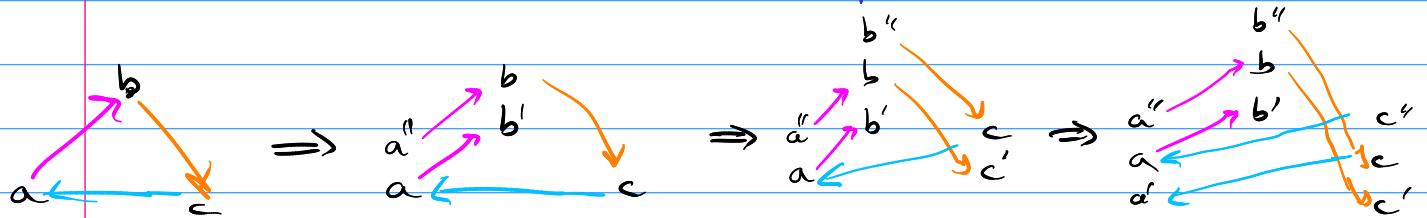
High-level idea:

- Take large minimal model of the query

- Iteratively dissociate edges

- Argue for termination

- Show that if the query is still true at the end
then this violates minimality



Why does it terminate?

Def: Leaf edge $\{a, b\}$

all facts involving b also involve a

which is

$$\begin{aligned} a'' &\rightarrow b \rightarrow c' \\ c'' &\rightarrow a \rightarrow b' \\ b'' &\rightarrow c \rightarrow a' \end{aligned}$$

Dissociation takes a non-leaf edge (center of an RST-pattern)

and replaces it by two leaf edges

\Rightarrow the number of non-leaf edges decreases

Cor: Any iterative dissociation process on an instance must terminate

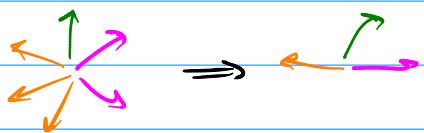
What are instances with no RST patterns left to dissociate? \rightarrow Only leaf edges



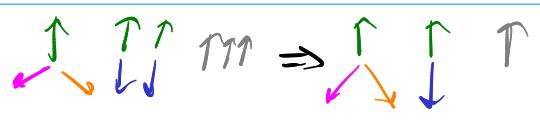
Lemma: An instance I with only leaf edges has a self-homomorphism to a constant-sized subset: call S_0 the maximal size.

- Finitely many possible fact combinations for an edge \rightarrow (depends on signature)
 \rightarrow let K be this number

- A connected component (star pattern) has a homomorphism to a subset of K edges (merge isomorphic branches)



- A disjoint union of star patterns has a homomorphism to a constant number of stars (merge isomorphic stars)



#P-hardness for unbounded UCQ[∞]: finding a tight pattern

(cont'd)

Let's conclude the proof: any unbounded UCQ[∞] \mathcal{Q} has a tight pattern.

As \mathcal{Q} is unbounded, it has a minimal model I of size $> S_0$.

$I \models \mathcal{Q}$

Iteratively dissociate until it is no longer possible: $I = I_0 \xleftarrow{\text{hom}} I_1 \xleftarrow{\text{hom}} \dots \xleftarrow{\text{hom}} I_n$

There are two possibilities:

1. $I_{i_0} \models \mathcal{Q}$ but $I_{i_0+1} \not\models \mathcal{Q}$ for some i

→ we have found a tight pattern

2. $I_n \models \mathcal{Q}$

→ let us show a contradiction

We know that I_n has only leaf edges

so $I_n \xrightarrow{\text{hom}} I'_n$ with $I'_n \subseteq I_n$ and $|I'_n| \leq S_0$ [use lemma]

as $I_n \models \mathcal{Q}$, we have $I'_n \models \mathcal{Q}$

as $I'_n \subseteq I_n$, we have $I'_n \xrightarrow{\text{hom}} I_n$

Thus, as $I_n \xrightarrow{\text{hom}} I_0 = I$ (homomorphism composition)

we have $I'_n \xrightarrow{\text{hom}} I$, with $|I'_n| \leq S_0$ and $I'_n \models \mathcal{Q}$

Take I' the image of I'_n in I by the homomorphism

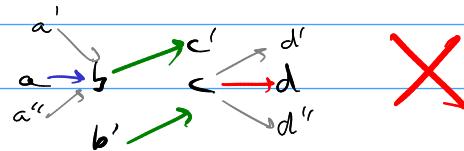
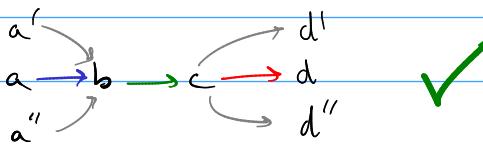
we have $I' \subseteq I$ with $|I'| \leq |I'_n| \leq S_0$ and $I' \models \mathcal{Q}$

→ This is impossible: it contradicts the minimality of I .

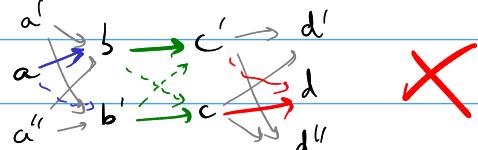
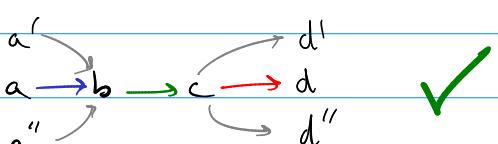
→ So we have found our tight pattern.

#P-hardness for unbounded UCD[∞]: from tight to critical

We know our query Q has a tight pattern



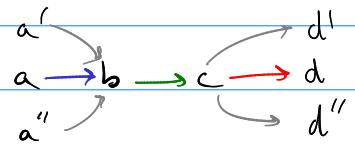
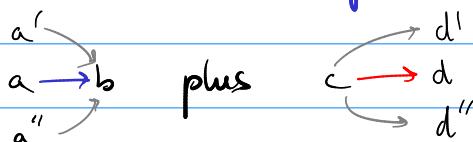
but what we need for hardness is a critical pattern



Idea: Minimality criterion :

Take a tight pattern minimizing the following weight (lexicographically) :

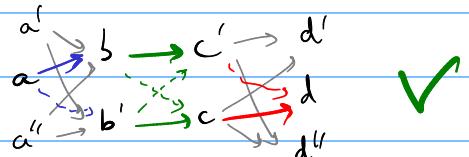
- the number of facts in $b \rightarrow c$
- the total number of incident facts in



Prop: Such a minimal tight pattern is in fact critical

for any choice of $U(a,b)$ $V(b,c)$ $W(c,d)$ in $a \rightarrow b$ $b \rightarrow c$ $c \rightarrow d$.

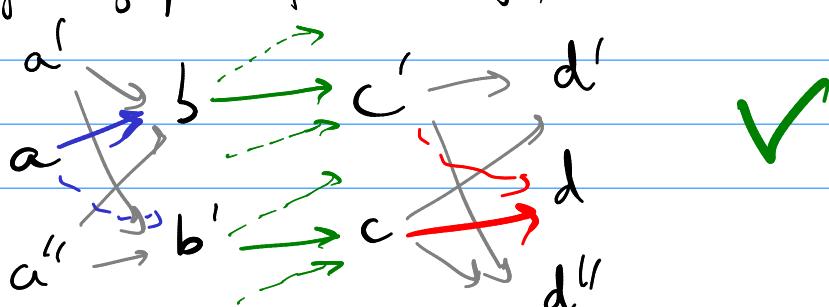
Proof: assume by contradiction that the query still holds when doing the coding



The edges $b \rightarrow c$ and $b' \rightarrow c'$ have one fact less than $b \rightarrow c$ in the original instance: $V(b,c)$ is missing.

→ so we can dissociate them without breaking the query
(or it would give a tight pattern of smaller weight)

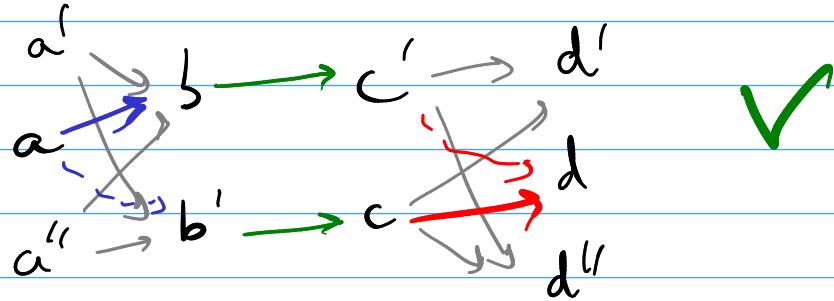
We get:



(cont'd)

#P-hardness for unbounded UCQ[∞]: from tight to critical

We can merge the \dashrightarrow intro with a homomorphism and Q still holds:



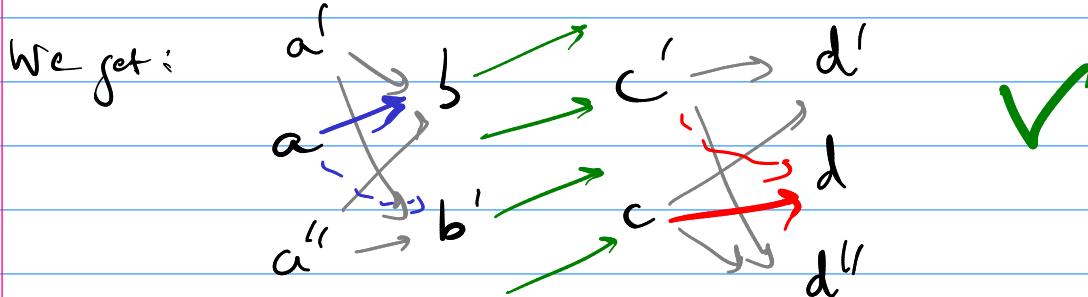
Now consider $a \dashrightarrow b' \rightarrow c \rightarrow d$
 $a'' \rightarrow b' \rightarrow c \rightarrow d''$

This has smaller weight than the minimal tight pattern:

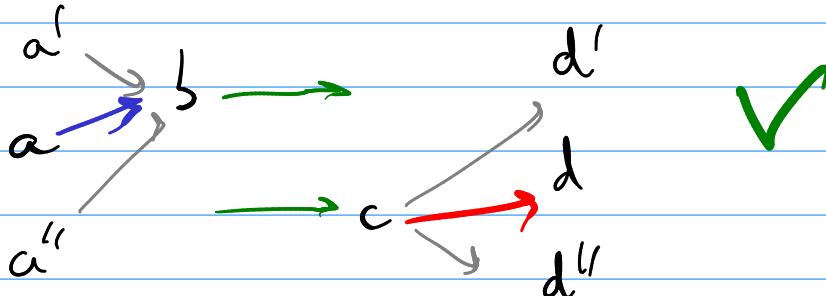
the incident fact $U(a, b)$ is missing from $a \dashrightarrow b$

→ so we can dissociate without breaking the query

Same reasoning on $b \rightarrow c'$



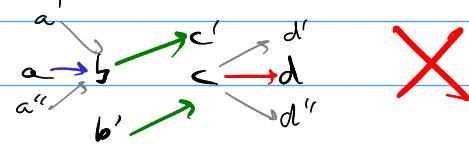
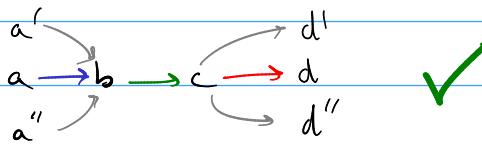
We homomorphically merge b and b', and c and c', and get:



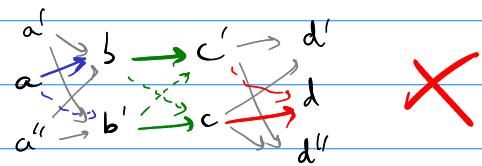
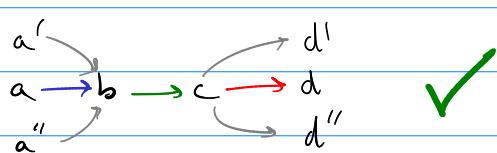
But this is the dissociation, which by assumption violates the query (tightness) \Rightarrow contradiction!
 So in fact the query was wrong all along, and we have found our critical pattern.

#P-hardness for unbounded UCD $^\infty$: recap

We have found a tight pattern (by iteratively dissociating)



We have shown that a minimal tight pattern is in fact critical



We have shown how to use a tight pattern to show hardness

(by reducing from PQE(Q_k) or PQE(Q_∞) on bipartite RST TIDS)

So we have indeed shown that PQE(Q) is #P-hard for our unbounded UCD $^\infty$

Bonus 1 : Treelike instances

Not many queries are tractable \Rightarrow restrict the shape of the data

Def: Tree width of a graph

Tree width of a relational instance \Rightarrow that of its Gaifman graph

We can allow expressive queries in Monadic second-order logic
 \rightarrow first-order logic (in particular $\forall \exists$, negation) plus quantification over sets

Courcelle,
Icc '90

Thm: Evaluating a fixed MSO query on an input instance
whose tree width is $\leq k$ for a constant k (treelike instance)
can be done in linear time data complexity

Proof: Encode the query to a tree automaton running on a tree encoding
of the tree decomposition

A. Bouchniak,
Szeider,
ICALP'15,
PODS'16

Thm: We can build in linear time a provenance d-D
of the query on the instance

Corollary: PQE for MSO queries over treelike instances is in linear time

Bonus 2: Lower bound on unbounded treewidth

What if the treewidth is not bounded? Can we use digewidth etc.?

⇒ Restrict to acyclic-2 instances for technical reasons

Def: Treewidth-constructible instance family

Infinite family \mathcal{I} of acyclic-2 instances such that, given $k \in \mathbb{N}$,
we can build in PTIME in k an instance $I \in \mathcal{I}$ such that $\text{treewidth}(I) \geq k$.

ex: All planar instances, all instances containing a triangle ---

Inhibition: Prevent that the growth in treewidth is "too slow" wrt the instance size

Def: For a query Q and family \mathcal{I} , $\text{PQE}(Q, \mathcal{I})$ is $\text{PQE}(Q)$

where the input TID must be an instance of \mathcal{I} (with arbitrary probabilities)

→ in particular we can close \mathcal{I} under subgraphs

A. Monet,
submitted,
2022
see A. Bousris,
Schnellhart, PODS'16

Theorem: For Q , the query testing if there are two incident facts (M_Q)
for any treewidth-constructible family \mathcal{I}
 $\text{PQE}(Q, \mathcal{I})$ is $\#P$ -hard under ZPP reductions

Key tool: finding arbitrary degree-3 planar graphs in a high-treewidth graph

Via Chekuri
and Chuzhoy,
JACM '16

Thm: There is $c \in \mathbb{N}$ such that the following holds:

Given a degree-3 planar graph H and a graph G with treewidth $|H|_c$,
we can embed G in H as a topological minor in ZPP.

Then show that Q_{it} is $\#P$ -hard on arbitrary subdivisions of an input graph

A. Monet,
Schnellhart
ICDT '18

Also lower bound on circuits: no polynomial-sized d-SDNNFs

Summary and outlook

What is the complexity of probabilistic query evaluation?

- On SSFCQs, well-understood
 - Hierarchical \Rightarrow tractable, via circuits
 - Non-hierarchical \Rightarrow intractable, even for UR
- On UCQs, understood, but...
 - PTIME algorithm which may fail
 - NP-hardness when the algorithm fails
- On UCQ $^\infty$, intractability except for UCQs
 - (?) Circuits?
 - (?) UR?
 - (?) Higherarity? (?) UR?

Not covered:

- Queries with negation [Fink, Olteanu, PODS'16]
- Queries with inequalities [Olteanu, Huang, SIGMOD'09]
- Provenance circuit bounds for other lineage classes [Jha, Suciu, TCS'13]
- Results on infinite domains [Carmeli, Grohe, Lindner, Strandkjaer, PODS'21]
- Support for updates [Bethke, Merg, PODS'21]
- Symmetric model counting [Beame, Van den Broeck, Gibkoff, Suciu, PODS'15]
- Practical implementation: see ProSQL [Senellart, ongoing]

Other directions for future research:

Livshits, Bertossi,
Kirelfeld, Sebag,
LMCS'21

- Approximations (beyond additive approximations)
- Reversibility of techniques? (enumeration, Shapley values)
- Connection to reasoning, e.g. probabilistic programming

Thanks for your attention! Questions?