

FUNCTIONAL DEPENDENCIES WITH PREDICATES: WHAT MAKES THE g_3 -ERROR EASY TO COMPUTE?

Simon Vilmin¹

Pierre Fausse -- Giovagnoli^{2,3}

Jean-Marc Petit²

Vasile-Marian Scuturici²

¹ LIS, Aix-Marseille Université

² LIRIS, INSA Lyon

³ Compagnie Nationale du Rhône

Data vs. Domain Knowledge

F	E	P
2.5	10.1	22.9
2.7	10.4	23.2
2.6	10.3	23.0
2.5	10.2	23.3
2.6	10.1	23.1
2.6	10.3	22.9

(Unique) counterexample

Data from a hydropower turbine:

incoming flow F ($\text{m}^3 \cdot \text{s}^{-1}$)

elevation E of the waterfall (m)

power P produced (MW)

Domain knowledge:

P is determined by E and F ,

i.e. $P = f(E, F)$

almost!

Question: is knowledge supported by data?

Some Database Terminology

attribute E with values in its domain, $\text{dom}(E)$

relation scheme R ,
a set of attributes

r	F	E	P
t_1	2.5	10.1	22.9
t_2	2.7	10.4	23.2
t_3	2.6	10.3	23.0
t_4	2.5	10.2	23.3
t_5	2.6	10.1	23.1
t_6	2.6	10.3	22.9

tuple t_4 , maps each attribute
to a value in its domain

$$t_4[F] = 2.5$$

relation r over R ,
set of tuples over R

Domain Knowledge and Functional Dependencies

Question: is knowledge supported by data?



P determined by E and F



Function $f(E, F) = P$



Functional Dependency (FD) $EF \rightarrow P$



Question: does the FD hold in the relation?

Functional Dependencies

Syntax

DEF: a Functional Dependency (FD) over R is an expression $X \rightarrow A$ where $X \subseteq R$ and $A \in R$.

Semantic

DEF: given r and $X \rightarrow A$ over R , $X \rightarrow A$ holds in r , $r \models X \rightarrow A$, if $\forall t_1, t_2 \in r$, $t_1[X] = t_2[X]$ entails $t_1[A] = t_2[A]$

r_1	B	C	A	D
t_1	0	0	1	a
t_2	0	0	2	b
t_3	0	1	1	c

$r_1 \not\models BC \rightarrow A$

(t_1, t_2)
Counterexample

r_2	B	C	A	D
t_1	0	1	1	a
t_2	0	1	1	b
t_3	1	2	3	c

$r_2 \models BC \rightarrow A$

FDs vs. Real Life

must hold on the **whole** dataset → **outliers?**



$\forall t_1, t_2 \in \tau, t_1[X] = t_2[X]$ entails $t_1[A] = t_2[A]$

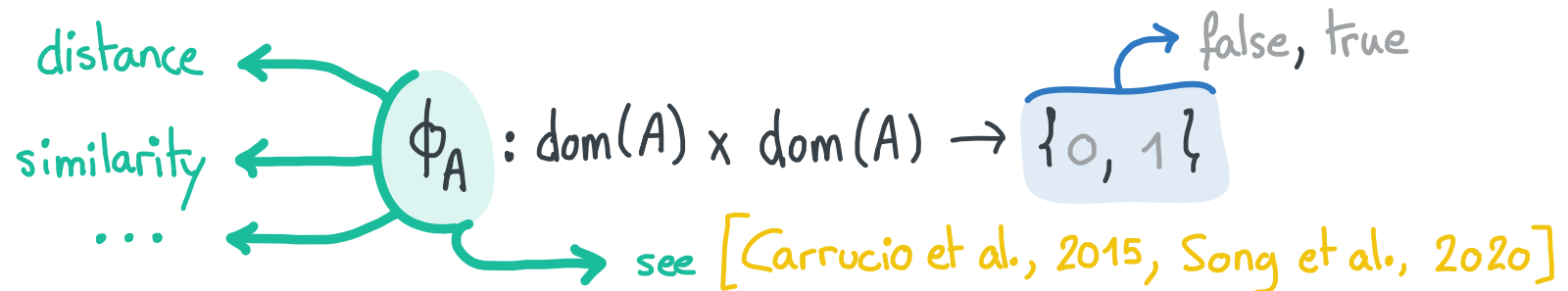
Compare values with **mathematical equality**
→ **imprecisions? other comparison criteria?**



Ideas: - replace equality with (binary) predicates
- use a **coverage measure**

Predicates to relax equality

binary predicate ϕ_A for $A \in R$: predicate to compare values in $\text{dom}(A)$



Relation scheme with predicates (R, Φ) : Φ set of predicates, one for each $A \in R$

DEF: given r and $X \rightarrow A$ over (R, Φ) , $X \rightarrow A$ holds in r wrt Φ , written $r \models_{\Phi} X \rightarrow A$, if $\forall t_1, t_2 \in r$,

$$\bigwedge_{B \in X} \phi_B(t_1[B], t_2[B]) = 1 \text{ implies } \phi_A(t_1[A], t_2[A]) = 1$$

Semantic

The g_3 -error with Predicates

g_3 -error : coverage measure for FDs with equality [Kivinen, Mannila, 1995]

minimum proportion of tuples to remove from r to satisfy $X \rightarrow A$

adapted to predicates [Foures--Giovagnoli et al., 2022]

DEF: Let (R, Φ) be a relation scheme with predicates and let $r, X \rightarrow A$ be a relation and a FD over (R, Φ) . The g_3 -error of $r, X \rightarrow A$ wrt Φ is:

$$g_3^{\Phi}(r, X \rightarrow A) = 1 - \frac{\max(\{|s| : s \subseteq r, s \models_{\Phi} X \rightarrow A\})}{|r|}$$

size of the largest subrelation satisfying $X \rightarrow A$

Back to the Example

Γ	F	E	P
t_1	2.5	10.1	22.9
t_2	2.7	10.4	23.2
t_3	2.6	10.3	23.0
t_4	2.5	10.2	23.3
t_5	2.6	10.1	23.1
t_6	2.6	10.3	22.9

(t_1, t_4) real counterexample
 (t_3, t_6) no longer a "false positive"
 $\{t_1, t_3, t_6\}$ maximum subrelation satisfying $EF \rightarrow P$
 $g_3^\phi(\Gamma, EF \rightarrow P) = 0.5$

$$\phi_P = \phi_E = \phi_F \quad \phi_P(x, y) = 1 \iff |x - y| \leq 0.1$$

$$\Phi = \{\phi_P, \phi_E, \phi_F\}$$

Back to our Problem

Question: is knowledge supported by data?

↓ function \leftrightarrow FDs

Question: does the FD hold in the relation?

↓ $g_3 + \Phi$

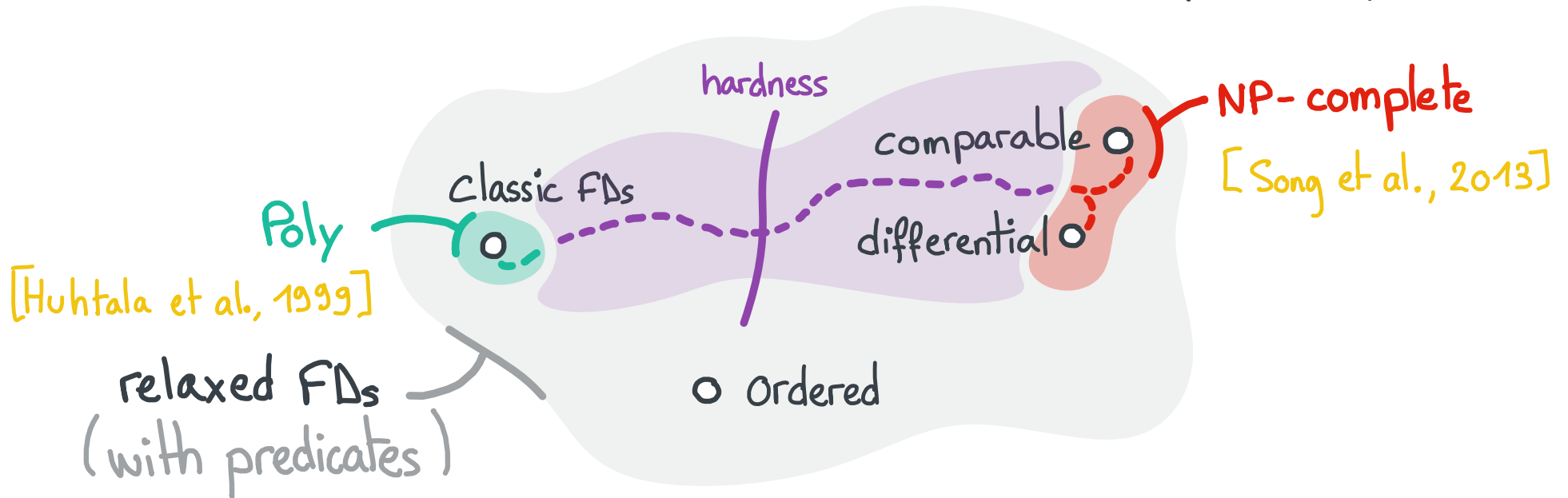
PROB: Error Validation Problem with Predicates (EVPP)

In: a relation scheme with predicates (R, Φ) , r and

$X \rightarrow A$ over (R, Φ) , $k \in \mathbb{R}$

Out: YES if $g_3^{\Phi}(r, X \rightarrow A) \leq k$, NO otherwise

The complexity of EVPP



Question: what makes EVPP tractable or not?

Idea: study predicate properties

$$(ref) \phi_A(x, x) = 1 \quad (sym) \phi_A(x, y) = 1 \Rightarrow \phi_A(y, x) = 1$$

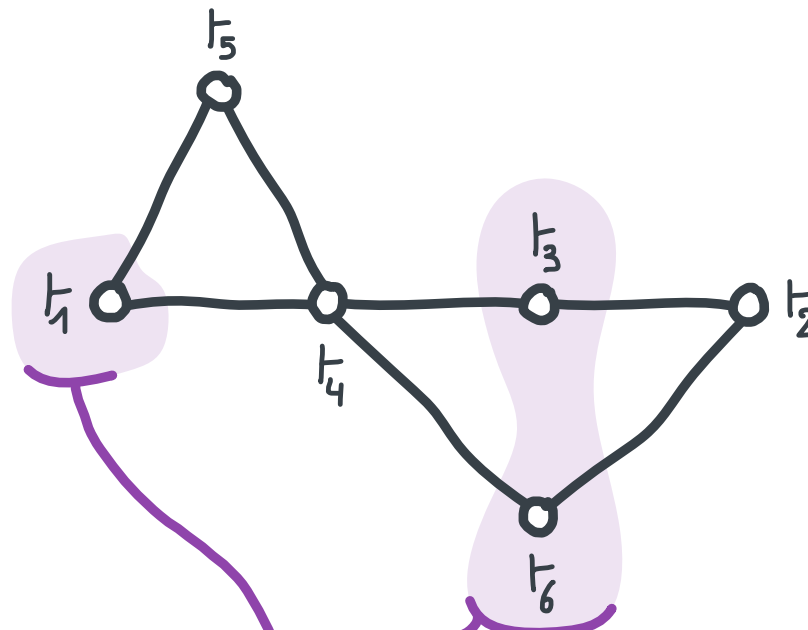
$$(tra) \phi_A(x, y) = \phi_A(y, z) = 1 \Rightarrow \phi_A(x, z) = 1$$

$$(asym) \phi_A(x, y) = \phi_A(y, z) = 1 \Rightarrow x = y$$

Conflict-graphs

Conflict-graph $CG_{\Phi}(r, EF \rightarrow P)$
 [Bertossi, 2011]

r	F	E	P
t_1	2.5	10.1	22.9
t_2	2.7	10.4	23.2
t_3	2.6	10.3	23.0
t_4	2.5	10.2	23.3
t_5	2.6	10.1	23.1
t_6	2.6	10.3	22.9



$s \models_{\Phi} EF \rightarrow P \iff s$ independent set of $CG_{\Phi}(r, EF \rightarrow P)$

EVPP and Maximum Independent Sets

Question: what makes EVPP tractable or not?

↓ CG_{Φ}

PROB: Maximum Independent Set (MIS)

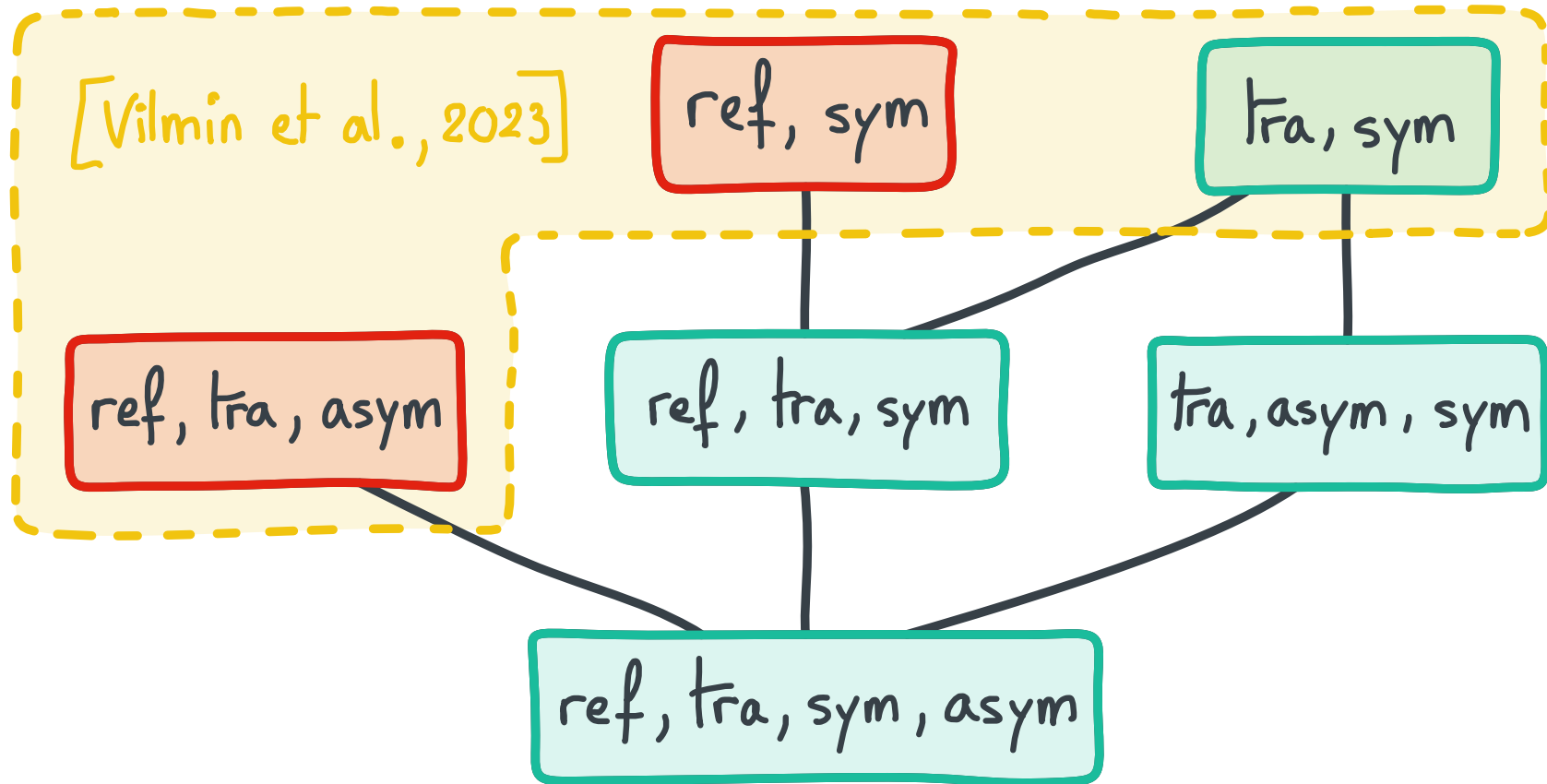
In: a graph $G = (V, E)$, $k \in \mathbb{N}$

Out: YES if there exists an ind. set $I \subseteq V$
of G s.t. $|I| \geq k$, NO otherwise

↓

Answer: The structure of CG_{Φ} imposed by Φ

Overview of our Results

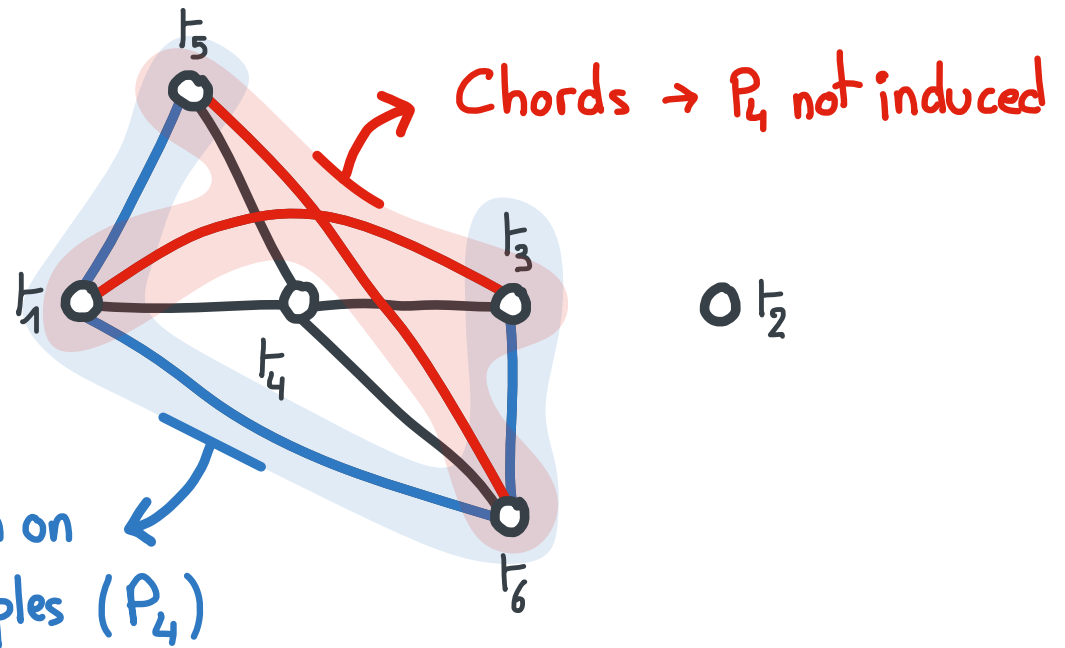


EVPP NP-C

EVPP Poly

Γ_a and sym

Γ	F	E	P
t_1	2.5	10.1	22.9
t_2	2.7	10.4	23.2
t_3	2.6	10.3	23.0
t_4	2.5	10.2	23.3
t_5	2.6	10.1	23.1
t_6	2.6	10.3	22.9



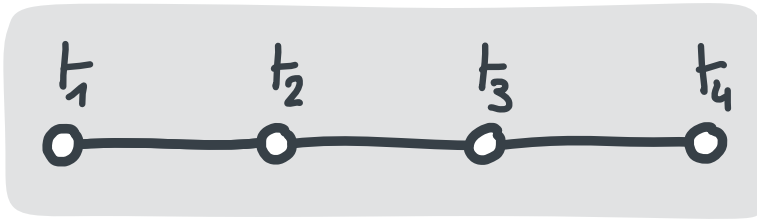
$$\phi_F(x, y) = 1 \iff 2.5 \leq x, y \leq 2.6$$

$$\phi_E(x, y) = 1 \iff 10.1 \leq x, y \leq 10.3$$

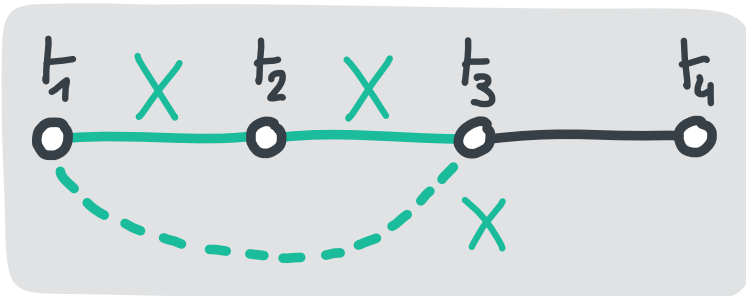
$$\phi_P(x, y) = 1 \iff 23.0 \leq x, y \leq 23.2$$

Γ_{ra} and sym, Ideas

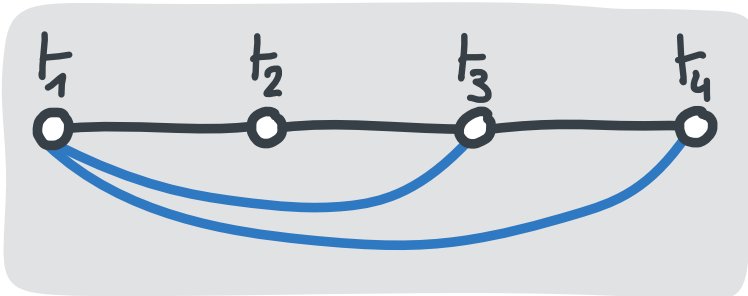
P_4 in $CG_{\mathbb{F}}(r, X \rightarrow A)$



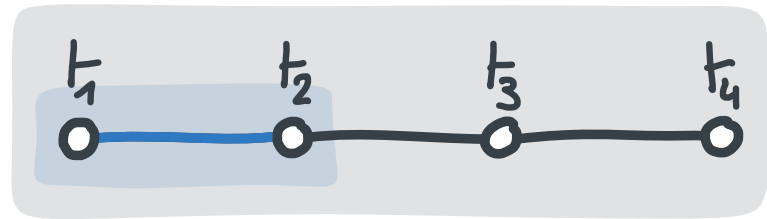
t_1, \dots, t_4 agree on X



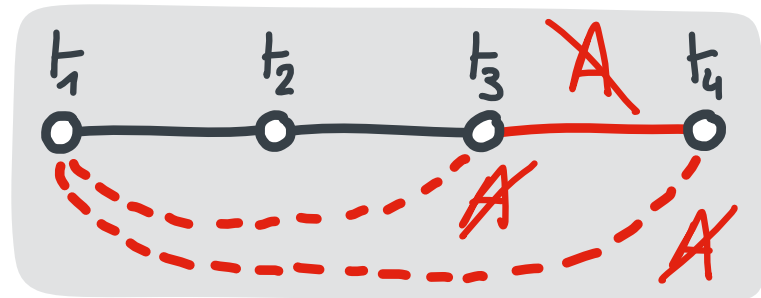
P_4 is not induced



t_1, t_2 agree on X but not on A



t_1, t_3 OR t_1, t_4 disagree on A



Co-graph, MIS poly

$\Gamma_{ra}, \text{sym}: CG_{\mathbb{F}}$ is P_4 -free

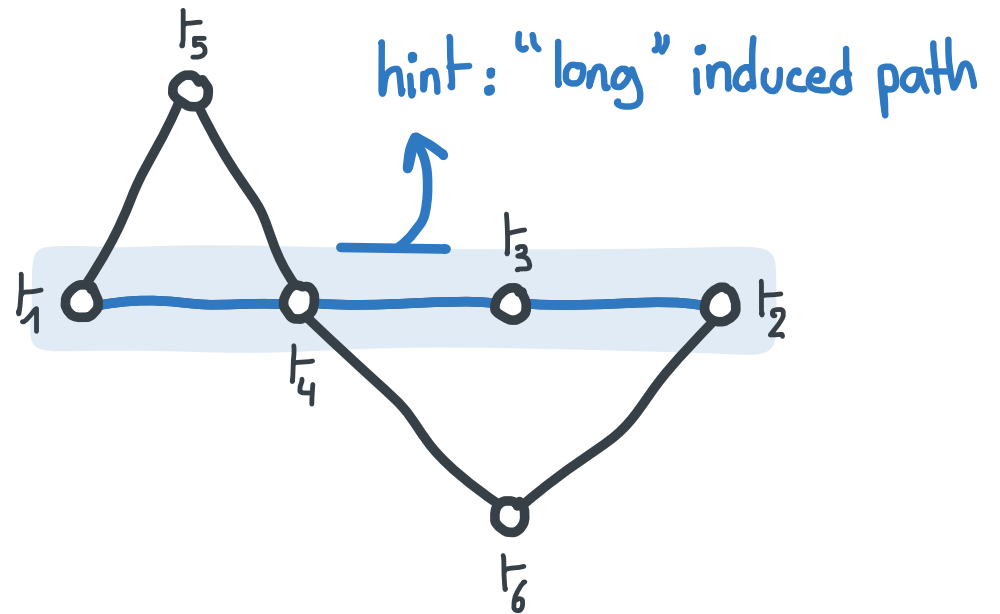
\rightarrow EVPP poly

ref and sym

ref, sym: $CG_{\mathbb{F}}$ can be any graph \rightarrow EVPP hard

\rightarrow MIS hard

Γ	F	E	P
t_1	2.5	10.1	22.9
t_2	2.7	10.4	23.2
t_3	2.6	10.3	23.0
t_4	2.5	10.2	23.3
t_5	2.6	10.1	23.1
t_6	2.6	10.3	22.9



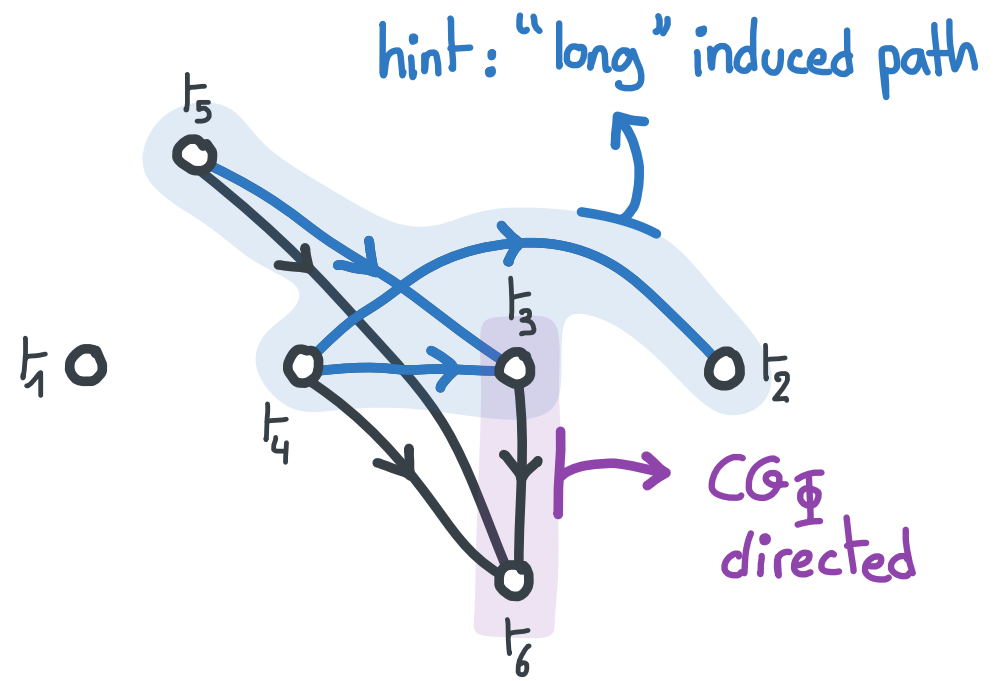
$$\phi_P = \phi_E = \phi_F \quad \phi_P(x, y) = 1 \iff |x - y| \leq 0.1$$

tra, ref and asym

tra, ref, asym: CG_{Φ} can be any 2-subdivision graph \rightarrow EVPP hard

MIS hard

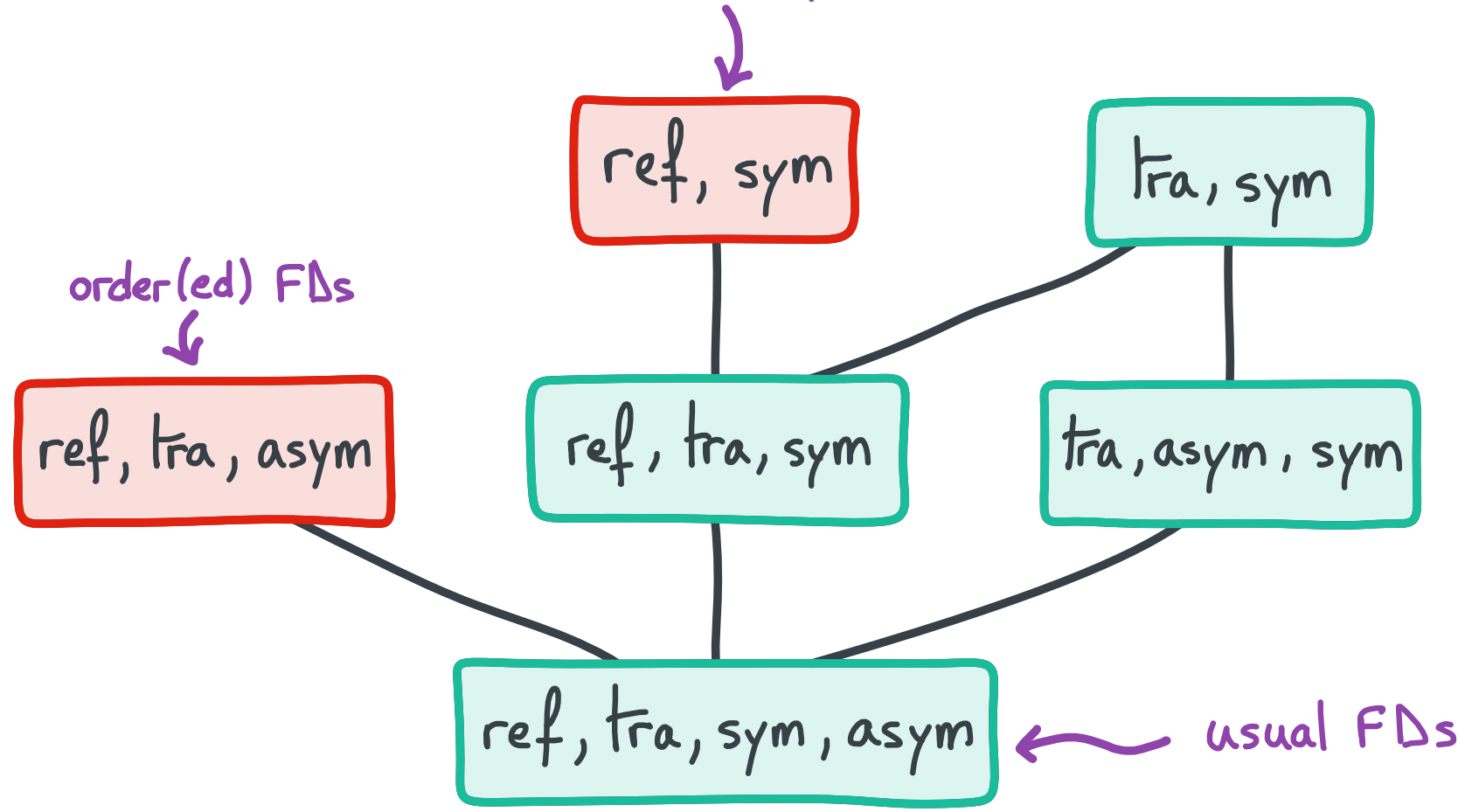
Γ	F	E	P
t_1	2.5	10.1	22.9
t_2	2.7	10.4	23.2
t_3	2.6	10.3	23.0
t_4	2.5	10.2	23.3
t_5	2.6	10.1	23.1
t_6	2.6	10.3	22.9



$$\phi_P = \phi_E = \phi_F \quad \phi_P(x, y) = 1 \iff x \leq y$$

Hierarchy

metric, similarity, comparable FDs



EVPP NP-C

EVPP Poly

Conclusion

EVPP: estimate the g_3 -error of a functional dependency with predicates

- can be used to confront experts knowledge against data

[Faure -- Giovagnoli, 2022]

- Complexity depends on the properties of predicates and the underlying conflict-graph [Bertossi, 2011]

Main results [Vilmin et al., 2023]

- having sym and tra \Rightarrow EVPP poly
- dropping sym or tra \Rightarrow EVPP NP-complete

Further research:

- Practical algorithms for special cases?
- Connections with repairs of sets of FDs? [Livshits et al., 2020]

References

L. Bertossi

Database repairing and consistent query answering
Synthesis Lectures on Data Management, 2011

[Bertossi, 2011]

L. Caruccio, V. Deufemia, G. Giuseppe

Relaxed Functional Dependencies—a survey of approaches
IEEE Transactions on Knowledge and Data Engineering, 2015

[Caruccio et al., 2015]

P. Faure -- Giovagnoli, J.-M. Petit, V.-M. Scutarici

Assessing the Existence of a Function in a Dataset with the g_3 indicator
IEEE International Conference on Data Engineering, 2022

[Faure -- Giovagnoli et al., 2022]

Y. Huhtala, J. Kärkkäinen, P. Porkka, H. Toivonen

TANE: An efficient algorithm for discovering functional and approximate dependencies
The Computer Journal, 1999

[Huhtala et al., 1999]

J. Kivinen, H. Mannila

Approximate inference of functional dependencies from relations
Theoretical Computer Science, 1995

[Kivinen, Mannila, 1995]

References

E. Livshits, B. Kimelfeld, R. Sudeepa

Computing optimal repairs for functional dependencies

ACM Transactions on Database Systems, 2020

[Livshits et al., 2020]

S. Song, L. Chen, P. Yu

Comparable dependencies over heterogeneous data

The VLDB journal, 2013

[Song et al., 2013]

S. Song, F. Gao, R. Huang, C. Wang

Data Dependencies over Big Data: A Family Tree

IEEE Transactions on Knowledge and Data Engineering, 2020

[Song et al., 2020]

S. Vilmin, P. Faure--Giovagnoli, J.-M. Petit, V.-M. Scuturici

Functional Dependencies with Predicates: What Makes the g_3 -error easy to compute?

LNCS, ICCS, 2023

[Vilmin et al., 2023]