

【ICDE 2014勉強会】

**Session 25:  
Uncertain and Probabilistic Data**

担当：胡(名大)

Some figures are copied from ICDE 2014 proceedings.

# Subgraph Pattern Matching over Uncertain Graphs with Identity Linkage Uncertainty

- ▶ Walaa Eldin Moustafa (U. Maryland), Angelika Kimming (KU Leuven),
- ▶ Amol Deshpande, Lise Getoor (U. Maryland)
- ▶ 不確実グラフへのクエリ
- ▶ 背景:

## 情報抽出と統合によるグラフデータの不確実性:

- ▶ 例: 同じ実世界エンティティは異なるデータソースに引用され、統合されたとき、不確実性を生み出す

## 3種類の不確実性

- ▶ **identity uncertainty** (同一性)
- ▶ **attribute value uncertainty**  
(ラベルの値の不確実性)
- ▶ **edge existence uncertainty**  
(エッジは存在するかどうか)

動機: 共同でこれらの不確実性を対処する方法はまだ不足している

クエリグラフ:

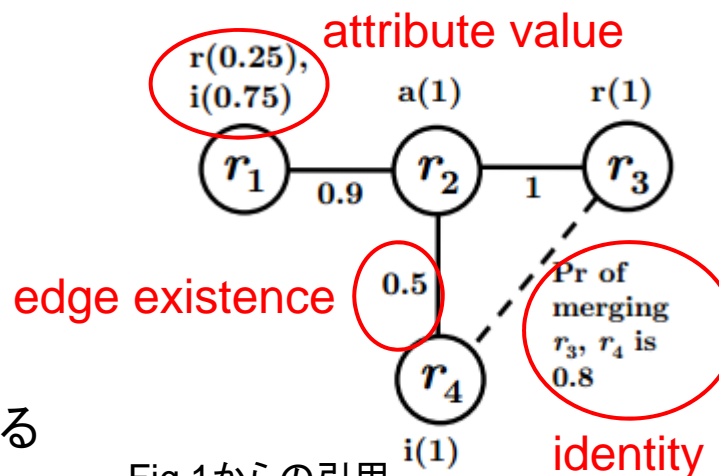
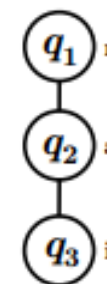


Fig.1からの引用

# PEGの提案

- ▶ **PEG** (probabilistic entity graph) : エンティティレベルで不確実グラフの分布を定義する

- ▶ **同一不確実性の定義** ( $f^N(s_1.n=v_1, \dots, s_k.n=v_k)$ )

$$f^N(s_1.n=v_1, \dots, s_k.n=v_k) = \begin{cases} p^s(s_i.x=T) & \text{if } v_i=T \text{ and, for all } j \neq i, v_j=F \\ 0 & \text{otherwise.} \end{cases}$$

- ▶ **ラベル値の不確実性の定義** ( $\Pr(s.l)$ )

$$\Pr(s.l) = [m^\Sigma(\{p^r | r \in s\})](s.l)$$

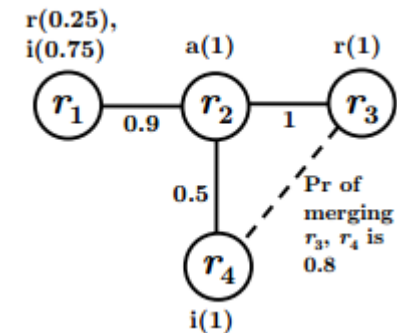
- ▶ **エッジの存在の不確実性の定義** ( $\Pr((s_1, s_2).e)$ )

$$\Pr((s_1, s_2).e) = [m^{\{T,F\}}(\{p^{(r_1, r_2)} | r_i \in s_i\})]((s_1, s_2).e)$$

例:

$$\Pr(s_1.l=r) = 0.25 ; \Pr(s_1.l=i) = 0.75$$

$$\Pr((s_1, s_2).e=T) = 0.9$$



例:

$$1) f^N(s_1.n=T, \text{other}=F) = 1$$

$$2) f^N(s_3.n=T, \text{other}=F) = 0.25$$

$$3) f^N(s_{34}.n=T, \text{other}=F) = 0.5$$

# サブグラフパターンマッチングアルゴリズム

手法: コンテキストウェアアパスインデキシング  
候補結合による削減

## Offline Phase

PEG



Compute  
⇒

Path Index

Key	Value
(a,a), 0.9	$P_{11}^u, P_{21}^u, P_{31}^u$
(a,b), 0.9	$P_{43}^u$
(b,b), 0.9	$P_{51}^u, P_{52}^u, P_{79}^u$
...	...

Context Information

$c(v, \sigma)$		a	b		
$ppu(v, \sigma)$	$v_1$	4	3		
	$v_2$	2	1		
$fpu(v, \sigma)$	$v_1$	$v_2$	$v_3$	0	5
	$v_2$	$v_3$	0.8	0.75	
	$v_3$	0.56	0.5		

(a)

- ▶ コンポーネントの確率を事前計算
- ▶ パスだけでなく、そのパスの近傍の文脈情報も含めてのインデックス

例:

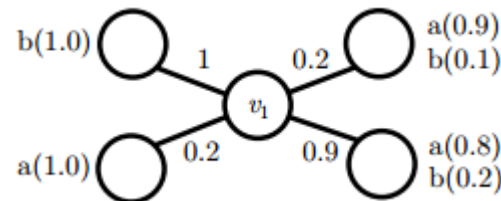
Cardinality a|bを含めて近傍ノードの個数

Partial Probability Upperbound

a=0.9; b=1.0

Full Probability Upperbound:

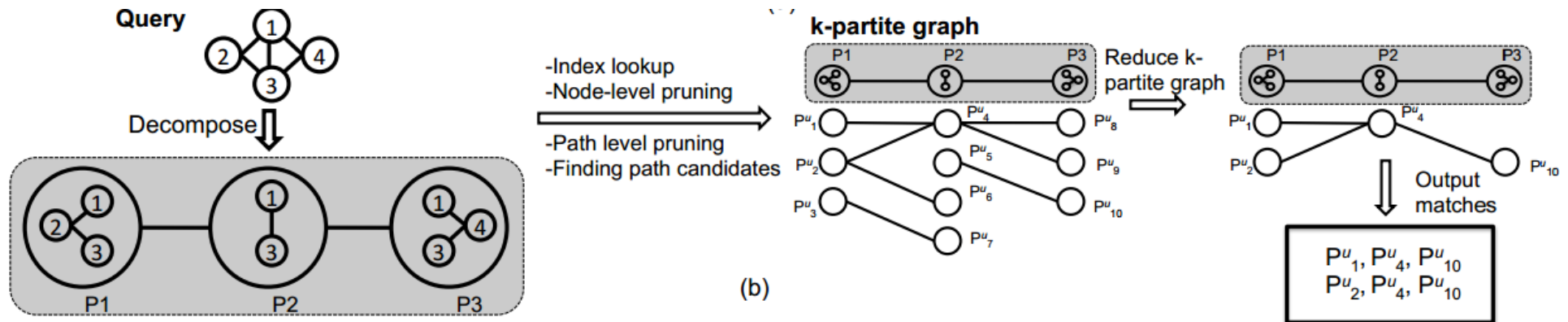
a=0.9X0.8; b=1X1.0



$\sigma$	a	b
$c(v_1, \sigma)$	3	3
$ppu(v_1, \sigma)$	0.9	1.0
$fpu(v_1, \sigma)$	0.72	1.0

# サブグラフパターンマッチングアルゴリズム

## Online Phase



- 1)クエリからの分解
- 2)分解したクエリ毎にパス候補を探す
- 3)結合できる候補パスを探し出す
- 4)k-partite graphに基づく候補の削減
- 5)フルクエリのマッチを探し出す