

【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 Kimmig (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:

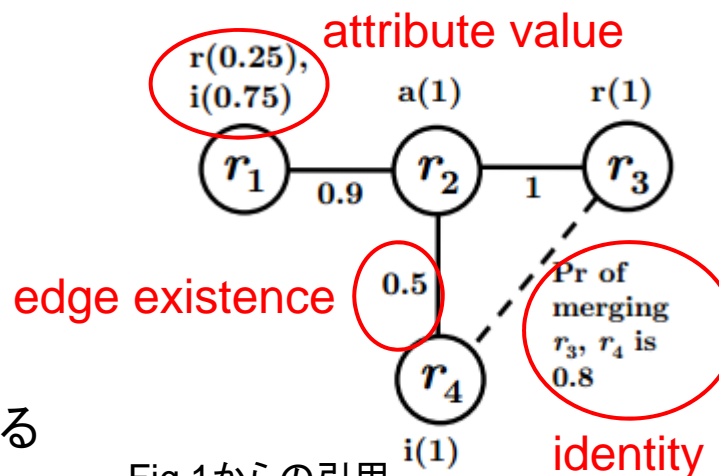
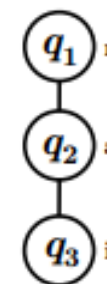


Fig.1からの引用

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** ($\Pr(s.l)$)

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

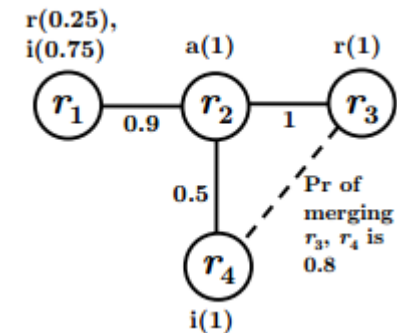
- ▶ **エッジの存在の不確実性の定義** ($\Pr((s_1, s_2).e)$)

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

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:

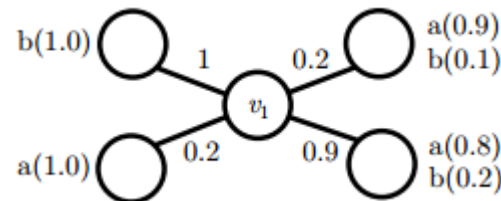
Cardinality a|bを含めて近傍ノードの個数

Partial Probability Upperbound

a=0.9; b=1.0

Full Probability Upperbound:

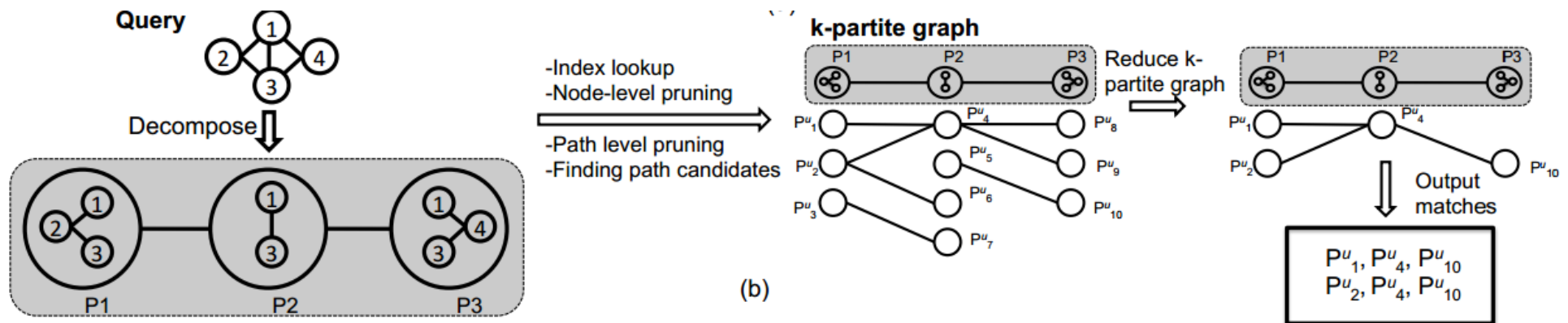
a=0.9X0.8; b=1X1.0



σ	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

【ICDE 2014勉強会】

Session 25: Uncertain and Probabilistic Data

担当：石川佳治(名大)

Some figures are copied from ICDE 2014 proceedings.

User-Driven Refinement of Imprecise Queries

- ▶ Bahar Qarabaqi, Mirek Riedewald (Northeastern U., USA)

- ▶ 目的: 曖昧な検索条件を対話的に洗練

- ▶ 特徴: 確率的なアプローチ

- ▶ 例: 鳥の観察画像の検索

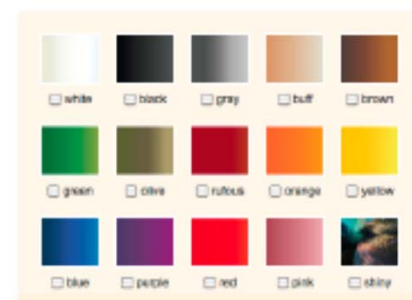

- ▶ 複数の属性: 羽の色, 大きさ, . . .

- ▶ ユーザは特定の画像を検索したいのではなく, 意図する鳥の種別 (例: ツグミ) で絞りこみたい

- ▶ ユーザは属性に対し**信念確率**を付与可能

- ▶ 例: 意図している鳥の多くが赤い羽根を持つとき, hasRedColorWing属性に (yes = 80%, no = 20%) という確率分布を指定

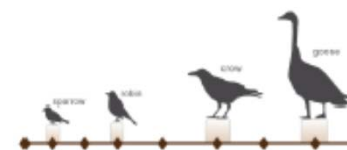
What color is the wing?
Select 1-3 colors.



I'm sure Pretty Sure Not Sure

Select the approximate size of the bird.

Click on a point on the scale below



I'm sure Pretty Sure Not Sure

インタフェースのイメージ

- ▶ 属性 $X_2 = \text{Size}$ と属性 $X_4 = \text{MainColor}$ に分布を指定した後の状況
- ▶ 類似度によるランキング(中)
- ▶ 未指定の属性で結果を改善しそうなものの順位
- ▶ 既指定の属性の感度 (sensitivity)

感度: 確率分布の値をいじるとどのくらい敏感に反応するか

Attribute	Quality Improvement	Rank	Species	Image	Specified Attribute	Sensitivity Score
X_3 :ShapeGroup	83.27	1	C_{245} :Eastern Bluebird		X_2 :Size	206.88
X_{14} :BillLength	81.12	2	C_{211} :Blue Jay		X_4 :MainColor	18.01
X_{11} :WingColor	74.98	3	C_{223} :Barn Swallow			
X_1 :Time	69.47	4	C_{212} :Western Scrub-Jay			
X_7 :Location	65.24	5	C_{233} :Red-breasted Nuthatch			
X_5 :BreastColor	64.02	6	C_{242} :Blue-grey Gnatcatcher			
X_6 :BreastPattern	63.18	7	C_{221} :Tree Swallow			
X_9 :BackColor	57.79	8	C_{331} :Indigo Bunting			
X_{18} :LegColor	49.35	9	C_{246} :Western Bluebird			
X_8 :BellyPattern	48.81	10	C_{210} :Steller's Jay			

技術的な部分

- ▶ 感度分析 (Sensitivity Analysis)
 - ▶ 取りうる確率分布の値の変化で、**ランクが最大限どの程度変化するか**を評価: この論文で頑張っている部分 (面倒)
 - ▶ ランクの違いはMinkowski距離で評価
 - ▶ 確率に関する性質をうまく使って効率的に判定
- ▶ 追加条件のランキング (画面の左側)
 - ▶ **改善度の期待値** (expected improvement) を評価
 - ▶ 注目する画像 (エンティティ) を高く, 他をより低く評価するのがよい
 - ▶ エントロピーの概念を用いた定式化
 - ▶ 実直な評価手法では対話性に難あり
 - ▶ 決定木のアンサンブルによる近似的な評価