

# On Embedding Uncertain Graphs

**Jiafeng Hu**, Reynold Cheng, Zhipeng Huang, Yixiang Fang and Siqiang Luo

Department of Computer Science

The University of Hong Kong

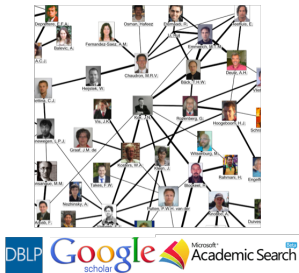
[jhu@cs.hku.hk](mailto:jhu@cs.hku.hk)

November 7, 2017

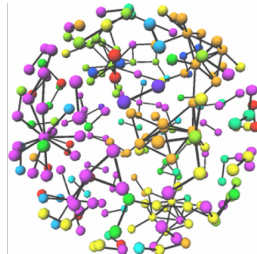
# Graphs are everywhere



**Social Network**



**Collaboration Network**

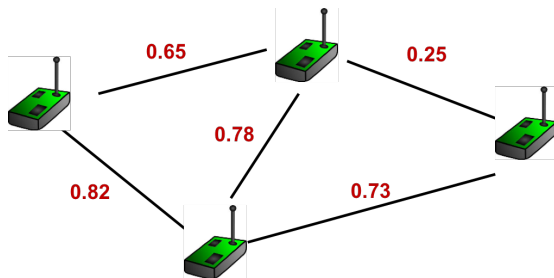


**Protein-Protein  
interaction Network**

# Uncertainty in Graph data

- Wireless sensor networks (WSNs)

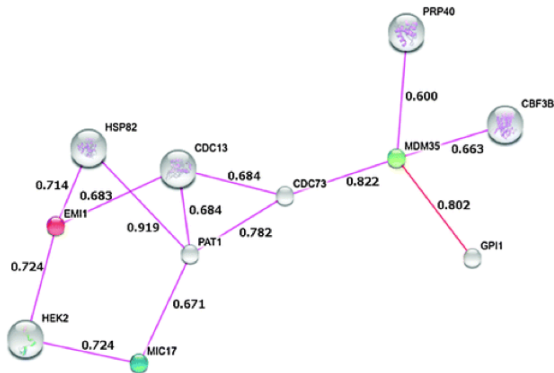
- ▶ Nodes: sensors
- ▶ Edges: wireless connectivity between sensors
- ▶ Uncertainties: probabilities of wireless connectivity



# Uncertainty in Graph data

## • Protein-Protein Interaction Networks (PPI)

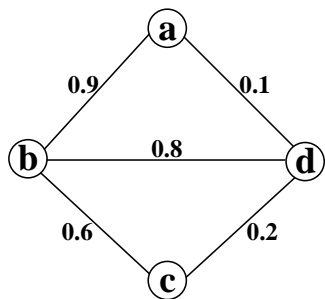
- ▶ Nodes: Proteins
- ▶ Edges: an interaction between proteins
- ▶ Uncertainties: probabilities of interactions between proteins derived from experimental evidence



Gabriele Cavallaro [Genome-wide analysis of eukaryotic twin CX<sub>9</sub>C proteins]

# Uncertainty in Graph data

Uncertain Graphs: each edge has an existence probability.



- Social Networks
- Traffic Networks
- Wireless Sensor Networks
- Protein-interaction Networks
- ...

# Possible World Semantics (PWS)

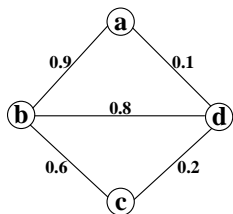
- Given an uncertain graph  $\mathcal{G}$ , a **possible world** of  $\mathcal{G}$  is a deterministic graph  $G = (V, E_G \subseteq E)$ . Assume the existence probabilities of edges are mutually **independent** [VLDB'10, KDD'10].

$$Pr[G] = \prod_{e \in E_G} P_e \prod_{e \in E \setminus E_G} (1 - P_e)$$

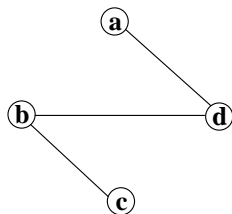
# Possible World Semantics (PWS)

- Given an uncertain graph  $\mathcal{G}$ , a **possible world** of  $\mathcal{G}$  is a deterministic graph  $G = (V, E_G \subseteq E)$ . Assume the existence probabilities of edges are mutually **independent** [VLDB'10, KDD'10].

$$Pr[G] = \prod_{e \in E_G} P_e \prod_{e \in E \setminus E_G} (1 - P_e)$$



(a) An uncertain graph  $\mathcal{G}$



(b) A possible world of  $\mathcal{G}$

$$\begin{aligned} Pr[G] &= P_{ad}P_{bd}P_{bc}(1 - P_{ab})(1 - P_{cd}) \\ &= 0.1 \times 0.8 \times 0.6 \times 0.1 \times 0.8 = 0.00384. \end{aligned}$$

## Tasks:

- Clustering [TKDE'13, ICDM'12]
- Classification [ICDM'09, SSDBM'14]
- k-NN queries [VLDB'10]
- ...



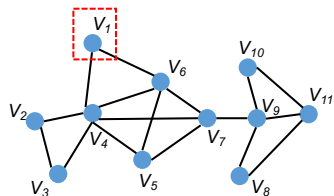
## Tasks:

- Clustering [TKDE'13, ICDM'12]
- Classification [ICDM'09, SSDBM'14]
- k-NN queries [VLDB'10]
- ...

## Shortcomings:

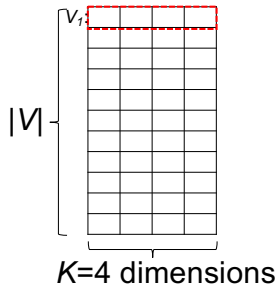
- High computational cost: expensive to compute similarities between nodes under PWS.
- Low adaptability: solutions are tailored for a particular mining task.

# Graph Embedding

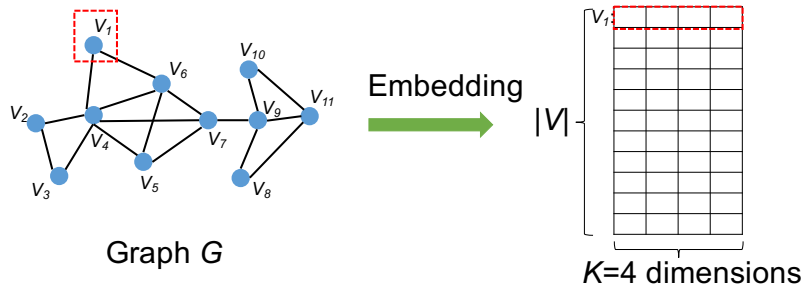


Graph  $G$

Embedding

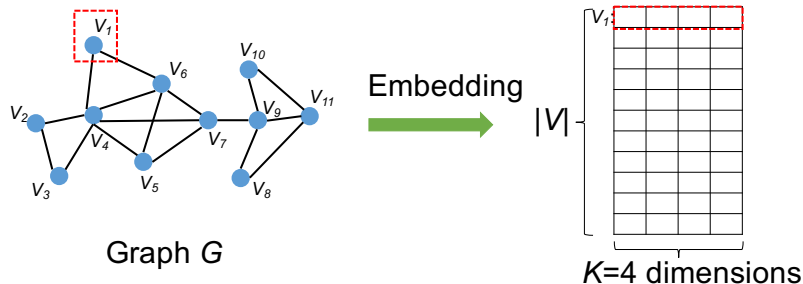


# Graph Embedding



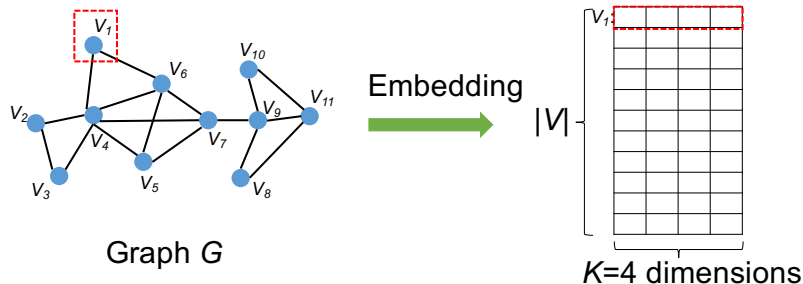
- Existing embedding solutions (for **deterministic** graphs): DeepWalk [KDD'14], LINE [WWW'15], node2vec [KDD'16], etc

# Graph Embedding



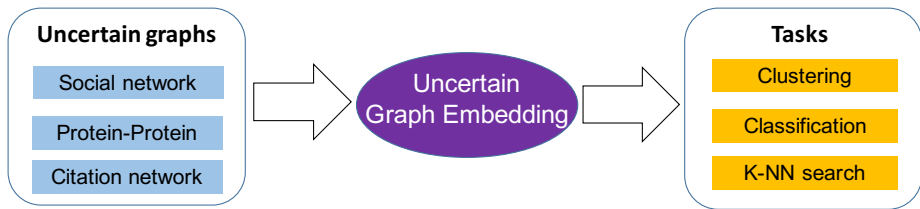
- Existing embedding solutions (for **deterministic** graphs): DeepWalk [KDD'14], LINE [WWW'15], node2vec [KDD'16], etc
- Not designed for uncertain graphs

# Graph Embedding

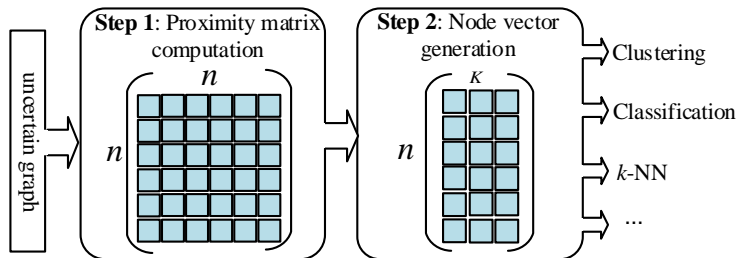


- Existing embedding solutions (for **deterministic** graphs): DeepWalk [KDD'14], LINE [WWW'15], node2vec [KDD'16], etc
- Not designed for uncertain graphs
- Simply remove the probabilities  $\rightarrow$  poor performance

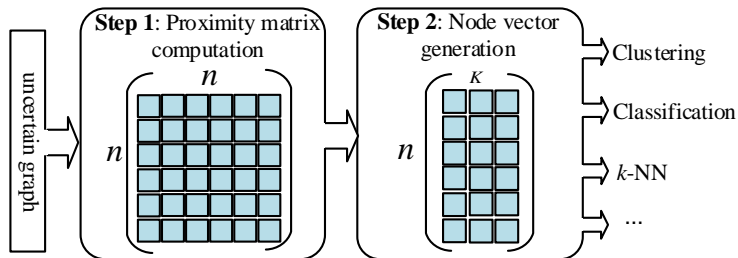
# Uncertain Graph Embedding



# URGE: UnceRtain Graph Embedding



# URGE: UnceRtain Graph Embedding

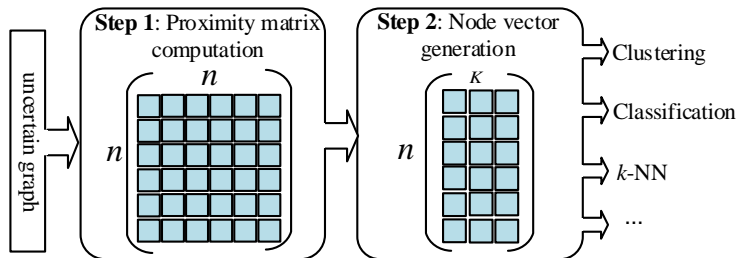


## Step 1: Proximity matrix $S$

- (second order) Expected Jaccard Similarity (EJS)
- (high order) Probabilistic Random Walk with Restart (PRWR)



# URGE: UnceRtain Graph Embedding

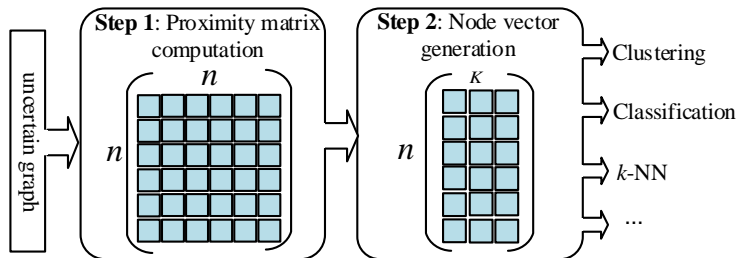


## Step 2: Objective function

$$\min_{\mathbf{U}, \tilde{\mathbf{U}}} \|\mathbf{S} - \mathbf{U}\tilde{\mathbf{U}}^T\|^2 + \frac{\lambda}{2} (\|\mathbf{U}\|^2 + \|\tilde{\mathbf{U}}\|^2)$$

- (Input)  $\mathbf{S} \in \mathbb{R}^{n \times n}$ : proximity matrix;
- (Input)  $\lambda$  controls the weight of the regularization term.
- (Output) matrices  $\mathbf{U} \in \mathbb{R}^{n \times K}$  and  $\tilde{\mathbf{U}} \in \mathbb{R}^{n \times K}$

# URGE: Uncertain Graph Embedding



## Step 2: Objective function

$$\min_{\mathbf{U}, \tilde{\mathbf{U}}} \|\mathbf{S} - \mathbf{U}\tilde{\mathbf{U}}^T\|^2 + \frac{\lambda}{2} (\|\mathbf{U}\|^2 + \|\tilde{\mathbf{U}}\|^2)$$

Negative sampling + asynchronous stochastic gradient algorithm (ASGD)

# How to compute the proximity matrix efficiently?

# Second-order Proximity: Expected Jaccard Similarity

Jaccard similarity between node  $u$  and  $v$  on a deterministic graph  $G$ :

$$S_{uv}^J = \frac{|N_G(u) \cap N_G(v)|}{|N_G(u) \cup N_G(v)|}$$

$N_G(x)$ : the neighbor set of node  $x$  on  $G$ .

# Second-order Proximity: Expected Jaccard Similarity

Jaccard similarity between node  $u$  and  $v$  on a deterministic graph  $G$ :

$$S_{uv}^J = \frac{|N_G(u) \cap N_G(v)|}{|N_G(u) \cup N_G(v)|}$$

$N_G(x)$ : the neighbor set of node  $x$  on  $G$ .

**Expected Jaccard Similarity (EJS):**

$$S_{uv}^{\text{EJS}} = \sum_{G \in \Omega(\mathcal{G})} (S_{uv}^J)_G Pr[G]$$

$\Omega(\mathcal{G})$ : the set of all possible worlds of  $\mathcal{G}$ .

# Computation of EJS

Lemma (A. Stuart, 1998; Z. Zhou, ICDM'13)

Given two nodes  $u$  and  $v$  of  $\mathcal{G}$ , let  $X_{uv} = |N_G(u) \cap N_G(v)|$  and  $Y_{uv} = |N_G(u) \cup N_G(v)|$ , where  $G$  is a possible world of  $\mathcal{G}$ . Then,

$$\mathbf{S}_{uv}^{\text{EJS}} = E \left[ \frac{X_{uv}}{Y_{uv}} \right] \approx \frac{E[X_{uv}]}{E[Y_{uv}]} - \frac{\text{Cov}(X_{uv}, Y_{uv})}{E[Y_{uv}]^2} + \frac{E[X_{uv}]\text{Var}(Y_{uv})}{E[Y_{uv}]^3}$$

- Z. Zou et al. study the problem of computing EJS between a pair of nodes  
→  $\mathcal{O}(d)$

- Z. Zou et al. study the problem of computing EJS between a pair of nodes  
→  $\mathcal{O}(d)$
- (Basic method) for each node  $u$  on  $\mathcal{G}$ , enumerate all 2-hop neighbors,  $v$ , and compute  $\mathbf{S}_{uv}^{\text{EJS}}$  →  $\mathcal{O}(nd^2 * d)$



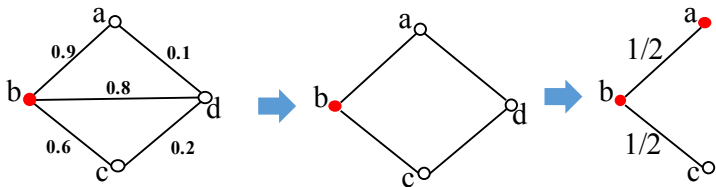
- Z. Zou et al. study the problem of computing EJS between a pair of nodes  
→  $\mathcal{O}(d)$
- (Basic method) for each node  $u$  on  $\mathcal{G}$ , enumerate all 2-hop neighbors,  $v$ , and compute  $\mathbf{S}_{uv}^{\text{EJS}} \rightarrow \mathcal{O}(nd^2 * d)$
- (**Our solution**) compute the EJS for all pair of nodes incrementally (i.e., the whole EJS matrix  $\mathbf{S}^{\text{EJS}}$ ) →  $\mathcal{O}(nd^2)$

# High-order Proximity: Probabilistic Random Walk with Restart

- Random walk (transition procedure) on uncertain graphs (for node  $u$ ) [VLDB'10]:
  - 1 generate a possible world  $G$  for  $u$ ;
  - 2 walk to a neighbor uniformly at random if there exists any neighbors of  $u$  in  $G$ , otherwise stay at  $u$ .

# High-order Proximity: Probabilistic Random Walk with Restart

- Random walk (transition procedure) on uncertain graphs (for node  $u$ ) [VLDB'10]:
  - 1 generate a possible world  $G$  for  $u$ ;
  - 2 walk to a neighbor uniformly at random if there exists any neighbors of  $u$  in  $G$ , otherwise stay at  $u$ .



# High-order Proximity: Probabilistic Random Walk with Restart

- Random walk (transition procedure) on uncertain graphs (for node  $u$ ) [VLDB'10]:
  - 1 generate a possible world  $G$  for  $u$ ;
  - 2 walk to a neighbor uniformly at random if there exists any neighbors of  $u$  in  $G$ , otherwise stay at  $u$ .
- **Probabilistic Transition Matrix, PTM** [VLDB'10, IS'15]:

$$\mathbf{W}_{uv} = \begin{cases} \prod_{(u,q) \in E} (1 - \mathbf{P}_{uq}), & u = v \\ \sum_{G \in \Omega(\mathcal{G}) \wedge (u,v) \in E_G} \frac{1}{d_u^G} Pr[G], & u \neq v \end{cases}$$

( $E_G$ : edge set of  $G$ ;  $d_u^G$ : out-degree of node  $u$  in  $G$ )

# High-order Proximity: Probabilistic Random Walk with Restart

- Random walk (transition procedure) on uncertain graphs (for node  $u$ ) [VLDB'10]:
  - 1 generate a possible world  $G$  for  $u$ ;
  - 2 walk to a neighbor uniformly at random if there exists any neighbors of  $u$  in  $G$ , otherwise stay at  $u$ .
- **Probabilistic Transition Matrix, PTM** [VLDB'10, IS'15]:

$$\mathbf{W}_{uv} = \begin{cases} \prod_{(u,q) \in E} (1 - \mathbf{P}_{uq}), & u = v \\ \sum_{G \in \Omega(\mathcal{G}) \wedge (u,v) \in E_G} \frac{1}{d_u^G} Pr[G], & u \neq v \end{cases}$$

( $E_G$ : edge set of  $G$ ;  $d_u^G$ : out-degree of node  $u$  in  $G$ )

- **Probabilistic Random Walk with Restart, PRWR**

$$\mathbf{S}^{\text{PRWR}} = (1 - \alpha)\mathbf{S}^{\text{PRWR}}\mathbf{W} + \alpha\mathbf{I}, \quad (\mathbf{I} \text{ is an identity matrix})$$

# Computation of PTM and PRWR

## Computation of PTM:

- (Basic method) an existing algorithm  $\rightarrow \mathcal{O}(nd^3)$  [IS'15]
- (Our method) further improvement, a hash-based method  $\rightarrow \mathcal{O}(hnd^2)$ ,  
 $h \ll d$ .

## Computation of PRWR:

- Monte Carlo method  $\rightarrow \mathcal{O}(nR\frac{1}{\alpha})$ 
  - ▶  $R$ : number of walkers;
  - ▶  $1/\alpha$ : expected length of random paths

# Experiments

Tasks: node clustering, node classification and  $k$ -NN search

Algorithms:

- Our algorithms:
  - ▶ URGE<sub>EJS</sub>: URGE algorithm based on EJS
  - ▶ URGE<sub>PRWR</sub>: URGE algorithm based on PRWR
- Existing embedding algorithms:
  - ▶ DeepWalk
  - ▶ LINE
  - ▶ node2vec <sub>$p$</sub>  <sup>$q$</sup>  (node2vec<sub>0.25</sub><sup>0.25</sup> and node2vec<sub>4</sub><sup>1</sup>)
- Existing non-embedding algorithms:
  - ▶ MCL (for deterministic graph clustering)
  - ▶ pKwikCluster (for uncertain graph clustering)
  - ▶ uBayes<sup>+</sup> (for uncertain graph classification)
  - ▶ MostProbPath (for uncertain graph  $k$ -NN)

# Task 1: Clustering

Dataset: 4 real uncertain Protein-Protein Interaction (PPI) networks<sup>1</sup>  
Ground truth: the complex-memberships lists from the MIPS database

Name	#Nodes	#Edges	Avg. Prob.
Krogan_core	2,708	7,123	0.68
Krogan_extend	3,672	14,317	0.42
Collins	1,622	9,074	0.78
Gavin	1,855	7,669	0.36

Table : Statistics of the PPI networks.

<sup>1</sup><http://www.nature.com/nmeth/journal/v9/n5/full/nmeth.1938.html>



# Task 1: Clustering (Cont'd)

- Metric: F1 score based on the confusion matrix (true positive, false positive, true negative and false negative)
- Hierarchical clustering in vector space (embedding-based algorithms)

# Task 1: Clustering (Cont'd)

- Metric: F1 score based on the confusion matrix (true positive, false positive, true negative and false negative)
- Hierarchical clustering in vector space (embedding-based algorithms)

Algorithm	Krogan_core	Krogan_extend	Collins	Gavin
DeepWalk	39.21	33.43	55.15	47.33
LINE	38.73	33.07	48.28	44.14
node2vec <sub>4</sub> <sup>1</sup>	39.30	33.06	52.42	46.17
node2vec <sub>0.25</sub> <sup>0.25</sup>	38.96	33.75	53.23	46.13
MCL	36.01	30.83	57.55	47.84
pKwikCluster	16.94	12.88	24.59	5.65
URGE <sub>EJS</sub>	38.39	30.08	55.61	<b>54.54</b>
URGE <sub>PRWR</sub>	<b>44.86</b>	<b>35.58</b>	<b>58.16</b>	52.59

Table : F1 scores (%) for clustering tasks.

# Task 1: Clustering (Cont'd)

- Metric: F1 score based on the confusion matrix (true positive, false positive, true negative and false negative)
- Hierarchical clustering in vector space (embedding-based algorithms)

Algorithm	Krogan_core	Krogan_extend	Collins	Gavin
DeepWalk	39.21	33.43	55.15	47.33
LINE	38.73	33.07	48.28	44.14
node2vec <sub>4</sub> <sup>1</sup>	39.30	33.06	52.42	46.17
node2vec <sub>0.25</sub> <sup>0.25</sup>	38.96	33.75	53.23	46.13
MCL	36.01	30.83	57.55	47.84
pKwikCluster	<b>16.94</b>	<b>12.88</b>	<b>24.59</b>	<b>5.65</b>
URGE <sub>EJS</sub>	38.39	30.08	55.61	<b>54.54</b>
URGE <sub>PRWR</sub>	<b>44.86</b>	<b>35.58</b>	<b>58.16</b>	52.59

Table : F1 scores (%) for clustering tasks.

## Task 2: Classification

Dataset:

- DBLP (co-authorship network): 45,583 edges, 14,376 papers, 4 classes
- Cora (citation network): 8,365 edges and 2,708 papers, 7 classes

\*Use the method proposed by P. Boldi et al. [VLDB'12] to do **obfuscation**.

Other setting:

- $k$  nearest neighbor classifiers (embedding-based algorithms).
- Varying training ratio ( $T_R$ ) from 20% to 80%
- Metric: Micro-F1, Macro-F1.

## Task 2: Classification (Cont'd)

Metric	Algorithm	20%	40%	60%	80%
Micro-F1(%)	DeepWalk	40.81	49.71	54.33	57.27
	LINE	40.87	47.96	53.26	56.52
	node2vec <sub>4</sub> <sup>1</sup>	41.25	51.29	55.67	59.27
	node2vec <sub>0.25</sub> <sup>0.25</sup>	40.16	49.33	53.53	56.98
	uBayes <sup>+</sup>	32.43	45.13	44.57	57.44
	URGE <sub>EJS</sub>	<b>58.00</b>	<b>63.31</b>	<b>66.36</b>	<b>69.45</b>
	URGE <sub>PRWR</sub>	52.16	58.08	61.12	63.97
Macro-F1(%)	DeepWalk	38.12	48.22	52.95	55.98
	LINE	39.02	46.37	51.70	55.08
	node2vec <sub>4</sub> <sup>1</sup>	39.42	49.21	53.93	57.69
	node2vec <sub>0.25</sub> <sup>0.25</sup>	38.11	47.64	51.93	55.37
	uBayes <sup>+</sup>	31.07	42.02	45.03	55.33
	URGE <sub>EJS</sub>	<b>55.48</b>	<b>61.45</b>	<b>64.49</b>	<b>67.50</b>
	URGE <sub>PRWR</sub>	49.86	56.41	59.64	62.42

Table : Results of classification on DBLP under different training ratio(%).

## Task 2: Classification (Cont'd)

Metric	Algorithm	20%	40%	60%	80%
Micro-F1(%)	DeepWalk	40.81	49.71	54.33	57.27
	LINE	40.87	47.96	53.26	56.52
	node2vec <sub>4</sub> <sup>1</sup>	41.25	51.29	55.67	59.27
	node2vec <sub>0.25</sub> <sup>0.25</sup>	40.16	49.33	53.53	56.98
	uBayes <sup>+</sup>	32.43	45.13	44.57	57.44
	URGE <sub>EJS</sub>	<b>58.00</b>	<b>63.31</b>	<b>66.36</b>	<b>69.45</b>
	URGE <sub>PRWR</sub>	52.16	58.08	61.12	63.97
Macro-F1(%)	DeepWalk	38.12	48.22	52.95	55.98
	LINE	39.02	46.37	51.70	55.08
	node2vec <sub>4</sub> <sup>1</sup>	39.42	49.21	53.93	57.69
	node2vec <sub>0.25</sub> <sup>0.25</sup>	38.11	47.64	51.93	55.37
	uBayes <sup>+</sup>	31.07	42.02	45.03	55.33
	URGE <sub>EJS</sub>	<b>55.48</b>	<b>61.45</b>	<b>64.49</b>	<b>67.50</b>
	URGE <sub>PRWR</sub>	49.86	56.41	59.64	62.42

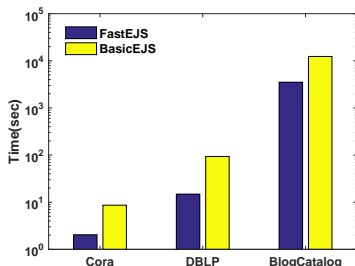
Table : Results of classification on DBLP under different training ratio(%).

# Efficiency

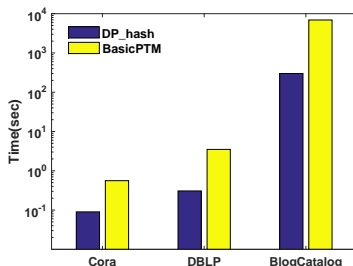
Datasets: DBLP, Cora, and BlogCatalog (relationships of the bloggers, 10K nodes and 665K edges)

Algorithms:

- EJS: BasicEJS vs FastEJS (**6+ times faster**)
- PTM: BasicPTM vs DP\_hash (**20+ times faster**)



(c) EJS



(d) PTM

# Conclusion

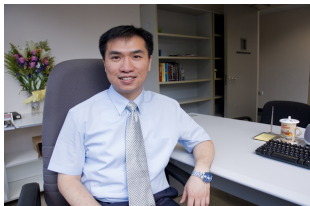
- Formulate the problem of uncertain graph embedding.
- Propose URGE, a proximity preserved embedding method for uncertain graphs.
- Develop efficient algorithms for two kinds of proximities (EJS and PRWR).
- Detailed evaluation on various tasks demonstrates the efficiency and effectiveness of the URGE solution.



# Our team



THE UNIVERSITY OF HONG KONG  
DEPARTMENT OF  
COMPUTER SCIENCE



Dr. Reynold Cheng ([ckcheng@cs.hku.hk](mailto:ckcheng@cs.hku.hk))



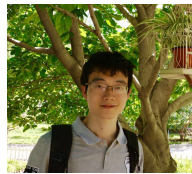
Jiafeng



Zhipeng



Yixiang



Siqiang

# Thanks!

## Q&A