

# Machine Learning on graphs. Link Prediction

I. Makarov & L.E. Zhukov

**BigData Academy MADE from Mail.ru Group**

**Network Science**



## 1 Link Prediction

- Similarity-based
- Matrix Factorization
- Random walks
- Other approaches and challenges

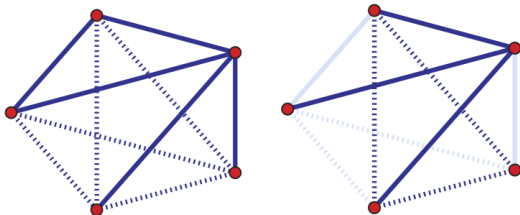
## 2 Graph Embeddings

- Problem statement
- Structural graph embeddings (simple models)

- Node classification (attribute inference)
- Link prediction (missing/hidden links inference)
- Community detection (clustering nodes in graph)
- Graph visualization (cluster projections)

- **Link prediction.** A network is changing over time. Given a snapshot of a network at time  $t$ , predict edges added in the interval  $(t, t')$
- **Link completion** (missing links identification). Given a network, infer links that are consistent with the structure, but missing (find unobserved edges)
- **Link reliability.** Estimate the reliability of given links in the graph.
  
- Predictions: link existence, link weight, link type

# Link prediction



- Graph  $G(V,E)$
- Number of "missing edges":  $|V|(|V| - 1)/2 - |E|$
- In sparse graphs  $|E| \ll |V|^2$ , Prob. of correct random guess  $O(\frac{1}{|V|^2})$

## Link prediction by proximity scoring

- 1 For each pair of nodes compute proximity (similarity) score  $c(v_1, v_2)$
- 2 Sort all pairs by the decreasing score
- 3 Select top  $n$  pairs (or above some threshold) as new links
- 4 Quality measurements - precision  $TP/(TP + FP)$ , precision at top  $N$

# Local similarity indices

Local neighborhood of  $v_i$  and  $v_j$

- Number of common neighbors:

$$s_{ij} = |\mathcal{N}(v_i) \cap \mathcal{N}(v_j)|$$

- Jaccard's coefficient:

$$s_{ij} = \frac{|\mathcal{N}(v_i) \cap \mathcal{N}(v_j)|}{|\mathcal{N}(v_i) \cup \mathcal{N}(v_j)|}$$

- Resource allocation:

$$s_{ij} = \sum_{w \in \mathcal{N}(v_i) \cap \mathcal{N}(v_j)} \frac{1}{|\mathcal{N}(w)|}$$

Adamic/Adar:

$$s_{ij} = \sum_{w \in \mathcal{N}(v_i) \cap \mathcal{N}(v_j)} \frac{1}{\log |\mathcal{N}(w)|}$$

- Preferential attachment:

$$s_{ij} = k_i \cdot k_j = |\mathcal{N}(v_i)| \cdot |\mathcal{N}(v_j)|$$

or

$$s_{ij} = k_i + k_j = |\mathcal{N}(v_i)| + |\mathcal{N}(v_j)|$$

- Clustering coefficient:

$$s_{ij} = CC(v_i) \cdot CC(v_j)$$

or

$$s_{ij} = CC(v_i) + CC(v_j)$$



- Local Path Index:

$$s_{lp} = A^2 + \alpha A^3$$

- High-order LPI:

$$s_{lp(n)} = \sum_{i=2}^n \alpha^{i-2} A^i$$

or

$$s_{ij} = CC(v_i) + CC(v_j)$$

# Path based methods

Paths and ensembles of paths between  $v_i$  and  $v_j$

- Shortest path:

$$s_{ij} = -\min_s \{path_{ij}^s > 0\}$$

- Katz score:

$$s_{ij} = \sum_{s=1}^{\infty} \beta^s |paths^{(s)}(v_i, v_j)| = \sum_{s=1}^{\infty} (\beta A)_{ij}^s = (I - \beta A)^{-1} - I$$

- Personalized (rooted) PageRank:

$$PR = \alpha(D^{-1}A)^T PR + (1 - \alpha) \cdot (e_i + e_j)$$

Liben-Nowell and Kleinberg, 2003

- Expected number of random walk steps:

hitting time:  $s_{ij} = -H_{ij}$

commute time  $s_{ij} = -(H_{ij} + H_{ji})$

normalized hitting/commute time  $s_{ij} = -(H_{ij}\pi_j + H_{ji}\pi_i)$

- SimRank:

$$\text{SimRank}(v_i, v_j) = \frac{C}{|\mathcal{N}(v_i)| \cdot |\mathcal{N}(v_j)|} \sum_{m \in \mathcal{N}(v_i)} \sum_{n \in \mathcal{N}(v_j)} \text{SimRank}(m, n)$$

Liben-Nowell and Kleinberg, 2003

# Community based methods

- Within-inter community/cluster of  $v_i, v_j \in C$

$$\sum_{w \in \mathcal{N}(v_i) \cap \mathcal{N}(v_j)} \frac{|\{w \in C\}|}{|\{w \notin C\}|}$$

- Common neighbors with community information,  $v_i, v_j \in C$ ,  $f(w) = 1$  if  $w \in C$

$$|\mathcal{N}(v_i) \cap \mathcal{N}(v_j)| + \sum_{w \in \mathcal{N}(v_i) \cap \mathcal{N}(v_j)} f(w)$$

- Resource allocation index with community information (soundarajan-hopcroft),  $v_i, v_j \in C$ ,  $f(w) = 1$  if  $w \in C$

$$\sum_{w \in \mathcal{N}(v_i) \cap \mathcal{N}(v_j)} \frac{f(w)}{|\mathcal{N}(w)|}$$

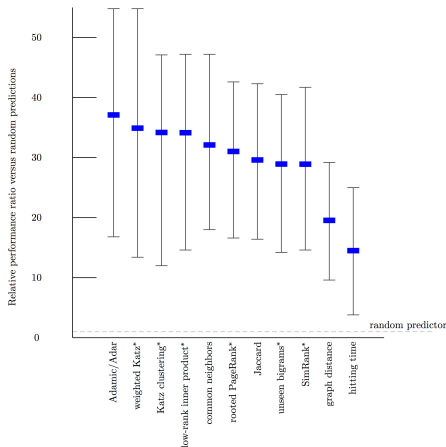
# Low-rank approximations

- Low-rank approximation (truncated SVD)

$$A = \sum_k^n U_k S_k V_k^T \rightarrow \sum_k^r U_k S_k V_k^T = A', r < n$$

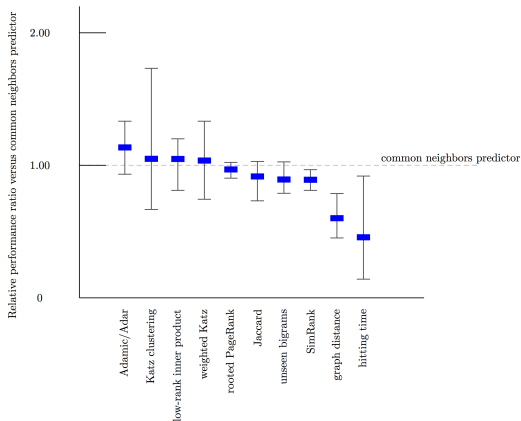
$$\begin{pmatrix} \hat{X} \\ \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & \\ \vdots & \vdots & \ddots & \\ x_{m1} & & & x_{mn} \end{pmatrix} \\ m \times n \end{pmatrix} \approx \begin{pmatrix} U \\ \begin{pmatrix} u_{11} & \dots & u_{1r} \\ \vdots & \ddots & \\ u_{m1} & & u_{mr} \end{pmatrix} \\ m \times r \end{pmatrix} \begin{pmatrix} S \\ \begin{pmatrix} s_{11} & 0 & \dots \\ 0 & \ddots & \\ \vdots & & s_{rr} \end{pmatrix} \\ r \times r \end{pmatrix} \begin{pmatrix} V^T \\ \begin{pmatrix} v_{11} & \dots & v_{1n} \\ \vdots & \ddots & \\ v_{r1} & & v_{rn} \end{pmatrix} \\ r \times n \end{pmatrix}$$

# Evaluation of scoring prediction



Ratio of predictor performance over the baseline, averaged 5 datasets

# Evaluation of scoring prediction



Ratio of predictor performance over the baseline, averaged 5 datasets

Liben-Nowell and Kleinberg, 2007

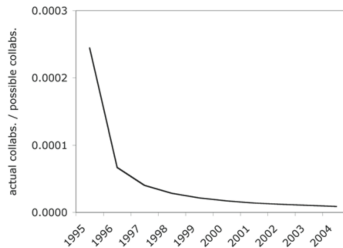
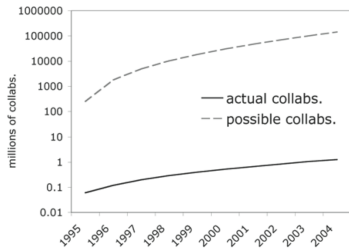
Challenging classification problem:

- Computational cost of evaluating of very large number of possible edges (quadratic in number of nodes)
- Highly imbalanced class distribution: number of positive examples (existing edges) grows linearly and negative quadratically with number on nodes



# Prediction difficulty

## Actual and possible collaborations between DBLP authors



## Extreme class imbalance

from Rattigan and Jensen, 2005

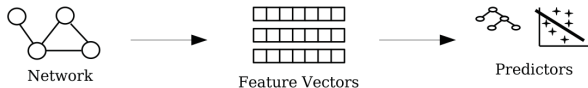
# Link prediction with supervised learning

Supervised learning:

- ① Features generation
- ② Model training
- ③ Testing (model application)

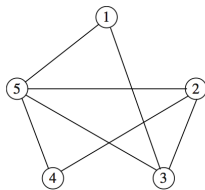
Features:

- Topological proximity features
- Aggregated features
- Content based node proximity features

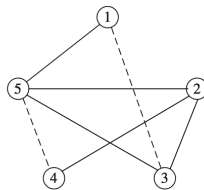


# Simple evaluation

Simple "hold out set" evaluation



Whole graph



Training graph

- Precision and Recall, F-measure

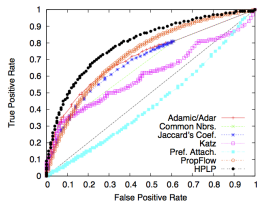
$$Precision = \frac{TP}{TP + FP}, \quad Recall = \frac{TP}{TP + FN}$$

$$F = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall}$$

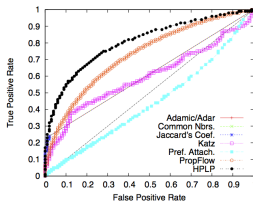
- True positive rate (TPR), False positive rate (FPR), ROC curve, AUC

$$TPR = \frac{TP}{TP + FN}, \quad FPR = \frac{FP}{FP + TN}$$

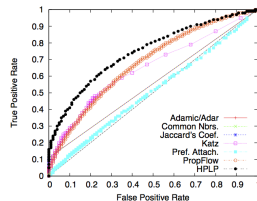
# ROC curves



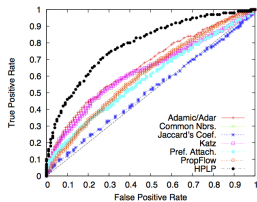
(a) phone  $n = 2$



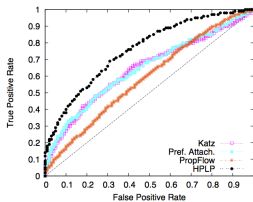
(b) phone  $n = 3$



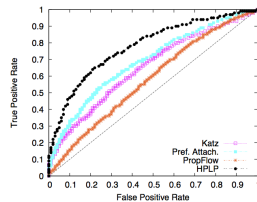
(c) phone  $n = 4$



(d) condat  $n = 2$



(e) condat  $n = 3$



(f) condat  $n = 4$

from Lichtenwalter, 2010

# Training and testing

## Evaluation for evolving networks

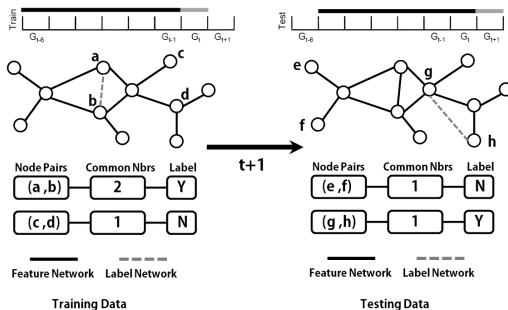


image from Y. Yang et.al, 2014

- Local model, Markov random fields [Wang, 2007]
- Hierarchical probabilistic model [Clauset, 2008]
- Probabilistic relations models:
  - Bayesian networks [Getoor, 2002]
  - relational Markov networks [Tasker, 2003, 2007]

## 1 Link Prediction

- Similarity-based
- Matrix Factorization
- Random walks
- Other approaches and challenges

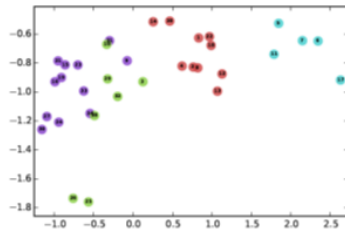
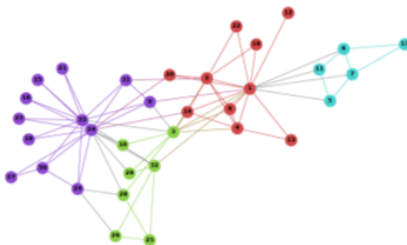
## 2 Graph Embeddings

- Problem statement
- Structural graph embeddings (simple models)



# Graph Embeddings

- Necessity to automatically select features
- Reduce domain- and task- specific bias
- Unified framework to vectorize network
- Preserve graph properties in vector space
- Similar nodes  $\rightarrow$  close embeddings

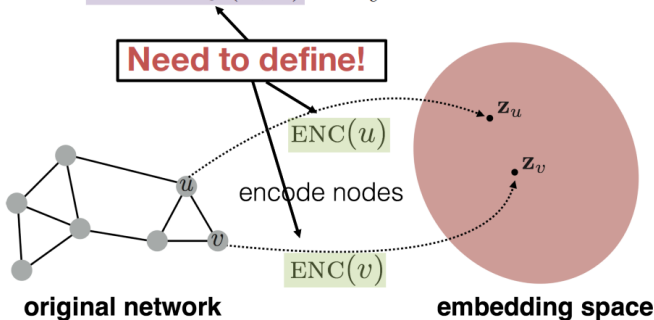


<sup>1</sup><http://snap.stanford.edu/proj/embeddings-www/>

# Graph Embeddings

- Define **Encoder**
- Define **Similarity**/graph feature to preserve graph properties
- Define similarity/distance in the embedding space
- **Optimize** loss to fit embedding with similarity computed on graph

**Goal:**  $\text{similarity}(u, v) \approx \mathbf{z}_v^\top \mathbf{z}_u$



# Structural Graph Embeddings

- Embedding look-up (each node - separate vector)
- Different similarity measures (adjacency, common neighbours, distances, exact function, etc.)
- Quadratic optimization for MSE loss
- Fast models via random walks

# First-order Proximity

- Similarity between  $u$  and  $v$  is  $A_{uv}$
- MSE Loss
- Variant of Matrix Decomposition

The diagram illustrates the Mean Squared Error (MSE) loss for graph proximity. The loss  $\mathcal{L}$  is defined as the sum over all node pairs  $(u, v) \in V \times V$  of the squared difference between the embedding similarity  $\mathbf{z}_u^\top \mathbf{z}_v$  and the weighted adjacency matrix element  $\mathbf{A}_{u,v}$ .

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \| \mathbf{z}_u^\top \mathbf{z}_v - \mathbf{A}_{u,v} \|^2$$

Annotations in the diagram:

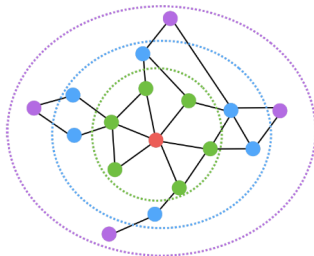
- $\mathcal{L}$ : loss (what we want to minimize)
- $\sum_{(u,v) \in V \times V}$ : sum over all node pairs
- $\mathbf{z}_u^\top \mathbf{z}_v$ : embedding similarity
- $\mathbf{A}_{u,v}$ : (weighted) adjacency matrix for the graph

from Lescovec et al., 2018

- Pros:
  - Use SGD for scalable optimization
  - Matrix factorization (SVD) or decomposition (QR) may be applicable
- Cons:
  - Quadratic complexity
  - Large embeddings space
  - No indirect graph properties are preserved

# Multi-order Proximity

- Similarity of neighborhoods of  $u$  and  $v$  via indices or  $k$ -hop paths
- Direct optimization of exact similarity metric

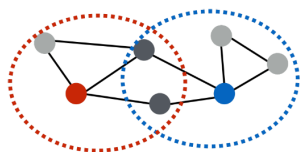


- **Red:** Target node
- **Green:** 1-hop neighbors
  - $\mathbf{A}$  (i.e., adjacency matrix)
- **Blue:** 2-hop neighbors
  - $\mathbf{A}^2$
- **Purple:** 3-hop neighbors
  - $\mathbf{A}^3$

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \|\mathbf{z}_u^\top \mathbf{z}_v - \mathbf{A}_{u,v}^k\|^2$$

# Multi-order Proximity

- Similarity score  $S_{uv}$  as Jaccard/Common Neighbours, etc. (HOPE)



$$\mathcal{L} = \sum_{(u,v) \in V \times V} \left\| \boxed{\mathbf{z}_u^\top \mathbf{z}_v} - \boxed{S_{u,v}} \right\|^2$$

embedding similarity

multi-hop network similarity (i.e., any neighborhood overlap measure)

- Weighted k-hop paths with different k (GraRep)

$$\tilde{\mathbf{A}}_{i,j}^k = \max \left( \log \left( \frac{(\mathbf{A}_{i,j}/d_i)}{\sum_{l \in V} (\mathbf{A}_{l,j}/d_l)^k} \right)^k - \alpha, 0 \right)$$

node degree

constant shift

from Lescovec et al., 2018

- Even worse complexity

- Similarity between  $u$  and  $v$  is probability to co-occur on a random walk
- Sample each vertex  $u$  neighborhood  $N_R(u)$  (multiset) by short random walks via strategy  $R$
- Optimize similarity considering independent neighbor samples via MLE (remind Word2Vec)

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} -\log(P(v|\mathbf{z}_u))$$

from Lescovec et al., 2018



- $P(v|z_u)$  is approximated via softmax over similarity  $z_u^T \cdot z_v$

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} -\log \left( \frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)} \right)$$


- Problem in second  $\Sigma$  over all nodes
- Hard to find optimal solution

# Negative Sampling

- Use *Negative Sampling* to approximate denominator

$$\log \left( \frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)} \right)$$

random distribution  
over all nodes


$$\approx \log(\sigma(\mathbf{z}_u^\top \mathbf{z}_v)) - \sum_{i=1}^k \log(\sigma(\mathbf{z}_u^\top \mathbf{z}_{n_i})), n_i \sim P_V$$

from Lescovec et al., 2018

- Sample in proportion to node degree
- Experiment with  $k$  to impact negative prior and robustness
- No need to sample non-connected edges — same as random

# Feature representation

- How to construct pair of nodes representation having node embeddings?
- Will it be more efficient than  $\sigma(z_i^t \cdot z_j)$

Symmetry operator	Definition
Average	$\frac{f_i(u) + f_i(v)}{2}$
Hadamard	$f_i(u) \cdot f_i(v)$
Weighted- $L_1$	$ f_i(u) - f_i(v) $
Weighted- $L_2$	$(f_i(u) - f_i(v))^2$
Neighbor Weighted- $L_1$	$\left  \frac{\sum_{w \in N(u) \cup \{u\}} f_i(w)}{ N(u)  + 1} - \frac{\sum_{t \in N(v) \cup \{v\}} f_i(t)}{ N(v)  + 1} \right $
Neighbor Weighted- $L_2$	$\left( \frac{\sum_{w \in N(u) \cup \{u\}} f_i(w)}{ N(u)  + 1} - \frac{\sum_{t \in N(v) \cup \{v\}} f_i(t)}{ N(v)  + 1} \right)^2$

DOI: [10.7717/peerj-cs.172/table-2](https://doi.org/10.7717/peerj-cs.172/table-2)

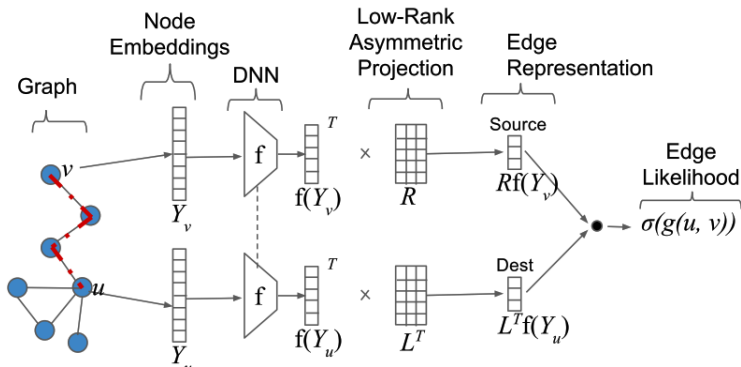
- How efficient simple solution?
- Works for undirected networks
- Samples neighbor information for low cost
- Not stable across different datasets ( $L_1$  works in general better than  $L_2$ )
- For weighted networks it is better to solve binary classification stacked with regression rather than directly solve link regression problem

from Makarov et al., 2019

# Directed network link prediction

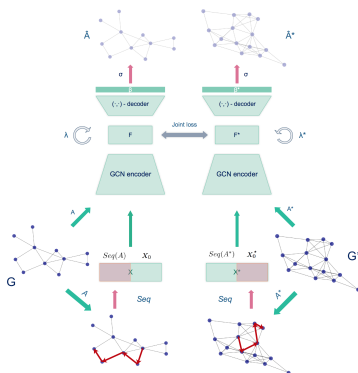
- When order matters, how to build classifier (see HOPE also)?
- Concat works not good probably - use asymmetric encoding via bi-linear form of compressed embeddings

$$M = LR, \quad g(u, v) = f(Y_u)^T M f(Y_v)$$



# Self-supervised learning via Line graph

- Edge-vertex dual (Line) graph allows to build dual representation and learn any edge embedding function
- Joint constraints on original and Line graph under bijective closure with agglutination of nodes embeddings in dual representation



from Makarov et al., 2021

- D. Liben-Nowell and J. Kleinberg. The link prediction problem for social networks. *Journal of the American Society for Information Science and Technology*, 58(7):1019-1031, 2007
- R. Lichtenwalter, J. Lussier, and N. Chawla. New perspectives and methods in link prediction. *KDD 10: Proceedings of the 16th ACM SIGKDD*, 2010
- M. Al Hasan, V. Chaoji, S. Salem, M. Zaki, Link prediction using supervised learning. *Proceedings of SDM workshop on link analysis*, 2006
- M. Rattigan, D. Jensen. The case for anomalous link discovery. *ACM SIGKDD Explorations Newsletter*. v 7, n 2, pp 41-47, 2005
- M. Al. Hasan, M. Zaki. A survey of link prediction in social networks. In *Social Networks Data Analytics*, Eds C. Aggarwal, 2011.

# References

- B. Perozzi, R. Al-Rfou, and S. Skiena. "Deepwalk: Online learning of social representations." In Proceedings of the 20th ACM SIGKDD international conference , pp. 701-710. 2014.
- Mutlu, Ece C., and Toktam A. Oghaz. "Review on graph feature learning and feature extraction techniques for link prediction." arXiv preprint arXiv:1901.03425 (2019).
- Makarov, Ilya, Olga Gerasimova, Pavel Sulimov, and Leonid E. Zhukov. "Dual network embedding for representing research interests in the link prediction problem on co-authorship networks." PeerJ Computer Science 5 (2019): e172.
- S. Abu-El-Haija, B. Perozzi, and R. Al-Rfou. "Learning edge representations via low-rank asymmetric projections." In Proceedings of the 2017 ACM CIKM conference, pp. 1787-1796. 2017.
- H. Cai, V.W. Zheng, and K.C.C. Chang. "A comprehensive survey of graph embedding: Problems, techniques, and applications." IEEE Transactions on Knowledge and Data Engineering 30, no. 9: 1616-1637, 2018