Uncertainty and Robustness of Graph Embeddings

Aleksandar Bojchevski Technical University of Munich, Germany

Graph Embedding Day 2018 - Lyon

Neglected aspects of graph embeddings

Capturing uncertainty Robustness to noise Robustness to adversarial attacks



Neglected aspects of graph embeddings

Capturing uncertainty

Robustness to noise Robustness to adversarial attacks

Uncertainty and Robustness of Graph Embeddings - Bojchevski



Nodes are points in a low-dimensional space













Graph2Gauss - 3 key modeling ideas

1. Uncertainty



2. Personalized ranking



3. Inductiveness





Embed nodes as (Gaussian) distributions

Sources of uncertainty:

- Conflicting structure and attributes
- Heterogenous neighborhood
- Noise, outliers, anomalies,







For each node *i*: nodes in its (k)-hop neighborhood should be closer to *i* compared to nodes in its (k + 1)-hop neighborhood







For each node *i*: nodes in its (*k*)-hop neighborhood should be closer to *i* compared to nodes in its (k + 1)-hop neighborhood









For each node *i*: nodes in its (k)-hop neighborhood should be closer to *i* compared to nodes in its (k + 1)-hop neighborhood

Example: closer in terms of the KL Diveregence

KL is asymmetric \Rightarrow handles directed graphs





Personalized ranking implies pairwise constraints for node *i*

 $D_{KL}(\mathcal{N}_{j}||\mathcal{N}_{i}) < D_{KL}(\mathcal{N}_{j'}||\mathcal{N}_{i})$ $\forall j \in N_{i}^{(k)}, \forall j' \in N_{i}^{(k')}, \forall k < k'$

set of nodes in the k-hop neighborhood of node i







Generalize to unseen nodes by learning a mapping from features to embeddings





Graph2Gauss - 3 key modeling ideas

1. Uncertainty



2. Personalized ranking



3. Inductiveness



Learning with energy-based loss

$$E_{ij} = D_{KL}(\mathcal{N}_j || \mathcal{N}_i)$$
 $\mathcal{L} = \sum_{(i,j,j')} (E_{ij}^2 + \exp^{-E_{ij'}})$

Closer nodes should have lower energy Naively: $O(N^3)$ complexity

Node-anchored sampling strategy:

- For each node same one another node from every neighborhood
- Less than 4.2% triplets seen to match performance
- Lower gradient variance

Graph2Gauss is parameter/data efficient



Graph2Gauss captures uncertainty

Uncertainty correlates with diversity

Diversity: number of distinct classes in a node's k-hop neighborhood



Graph2Gauss captures uncertainty

Uncertainty reveals the

intrinsic latent dimensionality of the graph

Detected latent dimensions

 \approx number ground-truth communities



_____ Uncertainty and link prediction

Prune dimensions with high uncertainty

Maintaining link prediction performance



____ Graph2Gauss is effective for visualization



Neglected aspects of graph embeddings

Capturing uncertainty

Robustness to noise

Robustness to adversarial attacks

Uncertainty and Robustness of Graph Embeddings - Bojchevski





https://www.semanticscholar.org





Graph clustering

- Maximize within-cluster edges
- Minimize between cluster edges



Partition V into two sets C_1 and C_2 , such that the sum of the inter-cluster edge weights $\operatorname{cut}(C_1, C_2) = \sum_{v_1 \in C_1, v_2 \in C_2} w(v_1, v_2)$ is minimized



Drawbacks:

- Tends to cut small vertex sets from the rest of the graph
- Considers only inter-cluster edges, no intra-cluster edges

Ratio Cut: Minimize
$$\frac{cut(C_1,C_2)}{|C_1|} + \frac{cut(C_2,C_1)}{|C_2|}$$

Normalized Cut: Minimize
$$\frac{cut(C_1,C_2)}{vol(C_1)} + \frac{cut(C_1,C_2)}{vol(C_2)}$$





Generalization to $k \ge 2$ clusters

Partition V into disjoint clusters C_1, \ldots, C_k such that

- Cut: $\min_{C_1,...,C_k} \sum_{i=1}^k cut(C_i, V \setminus C_i)$
- Ratio Cut: $\min_{C_1,...,C_k} \sum_{i=1}^k \frac{cut(C_i,V\setminus C_i)}{|C_i|}$
- Normalized Cut: $\min_{C_1,...,C_k} \sum_{i=1}^k \frac{cut(C_i,V\setminus C_i)}{\operatorname{vol}(C_i)}$



Finding the optimal solution is NP-hard

How to compute an approximate solution efficiently?

Laplacian matrix L = D - A

• *A* = (weighted) adjacency matrix, *D* = degree matrix

Observation: For any vector f we have $f^T \cdot L \cdot f = \frac{1}{2} \cdot \sum_{(u,v) \in E} W_{uv} (f_u - f_v)^2$

Normalized Laplacian $L_{sym} = D^{-\frac{1}{2}}LD^{-\frac{1}{2}} = I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$

Physical interpretation of the Laplacian (I)

Let f be a heat distribution over a graph with f_i = the heat at node v_i The heat transferred between v_i and v_j is prop. to $(f_i - f_j)$ if $(i, j) \in E$



https://en.wikipedia.org/wiki/Laplacian_matrix#/media/ File:Graph_Laplacian_Diffusion_Example.gif

Physical interpretation of the Laplacian (I)

Graph is viewed as an electrical circuit with edges as wires (resistors)

Apply voltage at some nodes and measure induced voltage at other nodes

Induced voltages minimizes $\sum_{(u,v)\in E} (x_u - x_v)$

We can find the voltage by minimizing $x^T L x = 0 \sqrt{10}$

$$0.5V$$

$$0.5V$$

$$0.5V$$

$$0.5V$$

$$0.5V$$

$$0.5V$$

$$0.625V$$

Properties of the Graph Laplacian

L is symmetric and positive semi-definite

The number of eigenvectors of *L* with eigenvalue 0 corresponds to the number of connected components

Algebraic connectivity of a graph is $\lambda_2(L)$

• The magnitude reflects how well connected the graph overall is

The spectrum of L encodes useful information about the graph

• Unfortunately, there exist co-spectral graphs

 $\underbrace{\text{Minimum cut and the graph Laplacian}}_{\text{Define indicator vector: }: h_{C_k}[i] = \begin{cases} \frac{1}{\sqrt{|C_i|}} & \text{if } v_i \in C_k \\ 0 & else \end{cases}$

Let $H = [h_{C_1}; h_{C_2}; ...; h_{C_k}]$

Observations:

 $H^{T}H = Id \text{ is orthonormal}$ $h_{C_{i}}^{T} \cdot L \cdot h_{c_{i}} = \frac{cut(C_{i}, V \setminus C_{i})}{|C_{i}|} \text{ and } h_{C_{i}}^{T} \cdot L \cdot h_{c_{i}} = (H^{T}LH)_{ii}$ $Ratio(ut(C_{i}, C_{i})) = \sum_{i=1}^{k} \frac{cut(C_{i}, V \setminus C_{i})}{|C_{i}|} = \sum_{i=1}^{k} \frac{cut(C_{i}, V \setminus C_{i})}{$

$$RatioCut(C_1, \dots, C_k) = \sum_{i=1}^k \frac{cut(C_i, V \setminus C_i)}{|C_i|} = \sum_{i=1}^k (H^T L H)_{ii} = trace(H^T L H)_{ii}$$



, Minimum cut and the graph Laplacian

Minimizing ratio-cut (normalized cut with L_{sym}) is equivalent to $\min_{C_1,\dots,C_k} trace(H^T L H) \text{ subject to } H^T H = Id$

Constraint relaxation: allow arbitrary values for H $\min_{H \in \mathbb{R}^{V \times K}} trace(H^T L H) \text{ subject to } H^T H = Id$

Standard trace minimization problem Optimal *H* = First *K* smallest eigenvectors of *L*

Spectral embedding: random walk view

 $L_{rw} = D^{-1}L = I - D^{-1}A = I - P$ is the the random walk Laplacian

• λ is an eigenvalue of L_{rw} with eigenvector u if and only if λ is an eigenvalue of L_{sym} with eigenvector $w = D^{1/2}u$

Let $P(B|A) = P(X_1 \in B | X_0 \in A)$ be the probability of a random walker currently at any node in A to transition to any node in B, for $A \cap B = \emptyset$ and $A, B \subset V$. Sample $X_0 \sim \pi$ from the stationary distribution $Ncut(A, \overline{A}) = P(\overline{A}|A) + P(A|\overline{A})$



Finding the spectral embedding = Solving an optimization task $H^* = k$ -first eigenvectors of L(A)

Problem: sensitive to noisy data Clean



Robustness via Latent Decomposition



 $H^* = \arg \min_{H \in \mathbb{R}^{n \times d}} Trace(H^T \cdot L(A^g) \cdot H)$ subject to $H^T \cdot D(A^g) \cdot H = Id$

 $A^*, H^* = \arg \min_{\substack{H \in \mathbb{R}^{n \times d} \\ A^g \in (\mathbb{R}_{\geq 0})^{n \times n}}} Trace(H^T \cdot L(A^g) \cdot H)$ subject to $H^T \cdot D(A^g) \cdot H = Id$ $A = A^g + A^c$ $\|A^c\|_0 \leq 2\theta \quad \text{global}$ $\forall_i : \|a_i^c\|_0 \leq \omega_i \quad \text{local}$

Jointly learn decomposition & embedding

Decomposition steered by the underlying embedding / clustering



Update *H*, Given $A^g/A^c \rightarrow Easy$

- Trace minimization problem
- Solution for H are the k first generalized eigenvectors of $L(A^g)$

Solution: Alternating optimization

$$A^*, H^* = \arg \min_{\substack{H \in \mathbb{R}^{n \times d} \\ A^g \in (\mathbb{R}_{\geq 0})^{n \times n}}} Trace(H^T \cdot L(A^g) \cdot H)$$

$$f([a^c_{uv}]_{u,v \in E}) = \sum_{u,v \in E} a^c_{uv} \left(\underbrace{\|\mathbf{h}_u - \mathbf{h}_v\|_2^2}_{\substack{nodes far away in \\ the embedding space}} - \underbrace{\|\sqrt{\lambda \circ \mathbf{h}_u}\|_2 - \|\sqrt{\lambda \circ \mathbf{h}_v}\|_2}_{prefers edges close} \right)$$

subject to $\|\cdot\|_{o}$ constraints

Update A^g/A^c , Given $H \rightarrow (NP)$ Hard

- Express eigenvalues of A_{new}^g in closed form
- A_{new}^g that minimizes the trace equivalent to maximizing f

Solution: Alternating optimization

Equivalent to Multidimensional Knapsack problem

- Greedy approximation
- Best possible approximation ratio of $\frac{1}{\sqrt{N+1}}$

Efficient solution in O(#edges)







Spectral embedding is sensitive to noisy data

Robustness via latent decomposition $A = A^g + A^c$ A^c ariginalgraphgraphcorruptions

Removed corrupted edges \Rightarrow increased discrimination



Neglected aspects of graph embeddings

Capturing uncertainty Robustness to noise

Robustness to adversarial attacks

Uncertainty and Robustness of Graph Embeddings - Bojchevski

Adversarial attacks on graph embeddings

(Spectral) Embeddings are not robust to noise / but we can remedy that

Are graph embeddings robust to adversarial attacks?

In domains where graph embeddings are used (e.g. the Web) adversaries are common and false data is easy to inject

Adversarial attacks in the image domain

Image of a tabby cat correctly classified



Adversarial attacks in the image domain

Image of a tabby cat correctly classified

Add imperceptible perturbation

Model classifies the cat as guacamole



The relational nature of the data might

Improve Robustness

embeddings are computed jointly rather than in isolation

Cause Cascading Failures

perturbations in one part of the graph can propagate to the rest

Attack possibilities

General attack

Targeted attack

Goal: decrease the overall quality of the embeddings

Goal: attack a specific node or a specific downstream task

Actions:

•

- Add/remove (flip) an edge
- Add/remove a node

Examples:

- Misclassify a target node t
- Increase/decrease the similarity of a set of node pairs $\mathcal{T} \subset V \times V$

Adjacency matrix of the graph after the attacker modified some entries

$$\hat{A}^* = \arg \max_{\hat{A} \in \{0,1\}^{N \times N}} \mathcal{L}(\hat{A}, Z^*)$$

$$Z^* = \min_{Z} \mathcal{L}(\hat{A}, Z) \quad subj. to \|\hat{A} - A\|_0 = 2f$$

Optimal embedding from the to be optimized graph \hat{A}

The attacker's budget

Attack model formally
Adjacency matrix of the graph after
the attacker modified some entries

$$\hat{A}^* = \arg \max_{\hat{A} \in \{0,1\}^{N \times N}} \mathbf{\Sigma}(\hat{A}, Z^*)$$
General attack

$$Z^* = \min_{Z} \mathbf{\Sigma}(\hat{A}, Z) \quad subj. \ to \ \|\hat{A} - A\|_{0} = 2f$$
Optimal embedding from the
to be optimized graph \hat{A}
The attacker's budget

Attack model formally
Adjacency matrix of the graph after
the attacker modified some entries

$$\hat{A}^* = \arg \max_{\hat{A} \in \{0,1\}^{N \times N}} \mathbf{L}_{atck}(\hat{A}, Z^*)$$
Targeted attack

$$Z^* = \min_{Z} \mathbf{L}(\hat{A}, Z) \quad subj. \ to \ \|\hat{A} - A\|_{0} = 2f$$
Optimal embedding from the
to be optimized graph \hat{A}
The attacker's budget

49



Discrete and Combinatorial Bi-level optimization problem Transductive learning ⇒ network poisoning setting

Evasion





Random-walk based embeddings

RW-based embeddings solve:

$$Z^* = \min_{Z} \mathcal{L}(\{r_1, r_2, ...\}, Z)$$
 with $r_i = RW_l(A)$

- $Z^* \in \mathbb{R}^{N \times K}$: learned embedding
- RW_l : e stochastic procedure that generates RWs of length l
- \mathcal{L} : model-specific loss e.g. skip-gram with negative sampling (SGSN)

Challenge: RW sampling precludes gradient based optimization

Example: DeepWalk

DeepWalk is equivalent* to factorizing

 $\widetilde{M} = \log \max(M, 1)$ Shifted Positive Pointwise Mutual Information (PPMI) Matrix

$$M = \frac{vol(A)}{T \cdot h} S \qquad S = (\sum_{r=1}^{T} P^r) D^{-1} \qquad P = D^{-1} A$$

b negative samples

window size T

transition matrix

with Z^* obtained by the SVD of $\widetilde{M} = U\Sigma V^T$ using the top K largest singular values/vectors i.e. $Z^* = U_K \Sigma_K^{1/2}$

Equivalent to $\min_{\widetilde{M}_K} ||\widetilde{M} - \widetilde{M}_K||_F^2$

The loss using the optimal embedding is $\mathcal{L}_{DW_1}(A, Z^*) = \sqrt{\sum_{p=K+1}^{|V|} \sigma_p^2}$, where $\sigma_1 \ge \sigma_2 \ge \cdots \ge \sigma_{|V|}$ are the singular values of $\widetilde{M}(A)$ ordered decreasingly

Idea: Given a perturbation ΔA , find the change in the singular values of $\widetilde{M}(A + \Delta A)$

Example: DeepWalk

$$\widetilde{M} = \log \max(M, 1)$$
 $M = \frac{vol(A)}{T \cdot b}S$ $S = (\sum_{r=1}^{T} P^r) D^{-1}$

Linearization: ignore the $log(\cdot)$ and $max(\cdot, 1)$ Scalars vol(A), T, b can be also ignore

Rewrite
$$\mathcal{L}_{DW_1}(A, Z^*) = \sqrt{\sum_{p=K+1}^{|V|} |\lambda_p|^2}$$

Thus, find a change in the spectrum of S after the attacker perturbed the graph ΔA

Compute the generalized spectrum (generalized eigenvalues/vectors) of A i.e. compute and U, Λ that solve $Au = \lambda Du$

Rewrite
$$S = (\sum_{r=1}^{T} P^r) D^{-1}$$
 as $S = U (\sum_{r=1}^{T} \Lambda^r) U^T$

simple function of the generalized eigenvalues λ_i of the graph

The task is now to find the change in generalized eigenvalues λ_p of the adjacency matrix A given a perturbation ΔA

Second Eigenvalue perturbation theory

Given U, Λ that solve $Au = \lambda Du'$ and a small perturbation ΔA , ΔD

Find U', Λ' that solve $(A + \Delta A)u' = \lambda'(D + \Delta D)u'$

First order approximation:

$$\lambda'_p = \lambda_p + u_p^T (\Delta A + \lambda_p \Delta D) u_p$$

for small ΔA and ΔD higher order terms become negligible

 ΔA is a matrix with only 2 non-zero elements for a single edge flip (i, j)namely $\Delta A_{ij} = \Delta A_{ji} = 1 - 2A_{ij} \coloneqq \Delta w_{ij}$

Similarly, ΔD has only two non-zero elements on the diagonal

Then we can approximate the generalized eigenvalues of A + ΔA in closed-form computable in O(1) time:

$$\lambda'_{p} = \lambda_{p} + \Delta w_{ij} \left(2u_{pi} \cdot u_{pj} - \lambda_{p} (u_{pi}^{2} + u_{pj}^{2}) \right)$$

- 1. DeepWalk is equivalent to a SVD of $\widetilde{M} = \log \max \left(\frac{vol(A)}{T \cdot b} S, 1 \right)$
- 2. The loss can be computed from the singular values / the spectrum of S
- 3. The spectrum of *S* can be easily computed from the generalized spectrum of A
- 4. For any given edge flip (i, j) we can compute in O(1) the spectrum of $A + \Delta A$

However,

$$\hat{A}^* = \arg \max_{\hat{A} \in \{0,1\}^{N \times N}} \mathcal{L}_{closed-form} \quad subj. to \|\hat{A} - A\|_0 = 2f$$
is still hard to optimize $-\binom{N^2}{f}$ ways to choose the flips

Greedy solution:

- 1. For each edge (i, j) calculate its impact on the loss if flipped
- 2. Pick the top f edges





To target node v we need the change in its embedding Z_v^* . That is we need the change in eigenvectors

Apply eigenvalue perturbation again to approximate the top K eigenvectors

For a given edge flip (i, j) we get:

$$u'_{p} = u_{p} - \Delta w_{ij} (A - \lambda D)^{+} \left(-\Delta \lambda_{p} u_{p} \circ d + E_{i} (u_{pj} - \lambda_{p} u_{pi}) + E_{j} (u_{pi} - \lambda_{p} u_{pj}) \right)$$





ρ Targeted attack: Node classification





Degree attack



Analysis of adversarial edges



<u>م</u> Transferability l D

	DW (SVD)	DW (SGNS)	n2v	Spect. Embd.	Label Prop.	GCN
$f = 250 \ (0.8\%)$	-3.59	-2.37	-2.04	-2.11	-5.78	-3.34
f = 500 (1.6%)	-4.62	-3.97	-3.48	-4.57	-8.95	-2.33
f = 250 (1.7%)	-7.59	-5.73	-6.45	-3.58	-4.99	-2.21
f = 500 (3.4%)	-9.68	-11.47	-10.24	-4.57	-6.27	-8.61



Node embeddings are vulnerable to adversarial attacks

Poisoning has negative effect on the embeddings quality and the downstream tasks

Attacks are transferable – they generalize to many models



_____ Important aspects of graph embeddings

Capturing uncertainty Robustness to noise Robustness to adversarial attacks

