

A Restricted Black-box Adversarial Framework Towards Attacking Graph Embedding Models

Members: 朱偉綸、黃秉茂、許齊巖



Abstract

- Graph Embedding Model can effectively transform Node information, but is vulnerable to specific adversarial attacks
- Many researches now focus on White-box Attack, which can learn information such as model, model parameters, gradient direction, etc.
- This research proposes GF-Attack to attack Graph Embedding Model from the perspective of Graph Signal Processing
- In order to examine the generalization of the attack, this article will focus on model verification of GCN, SGN, DeepWalk, LINE, etc.

Introduction



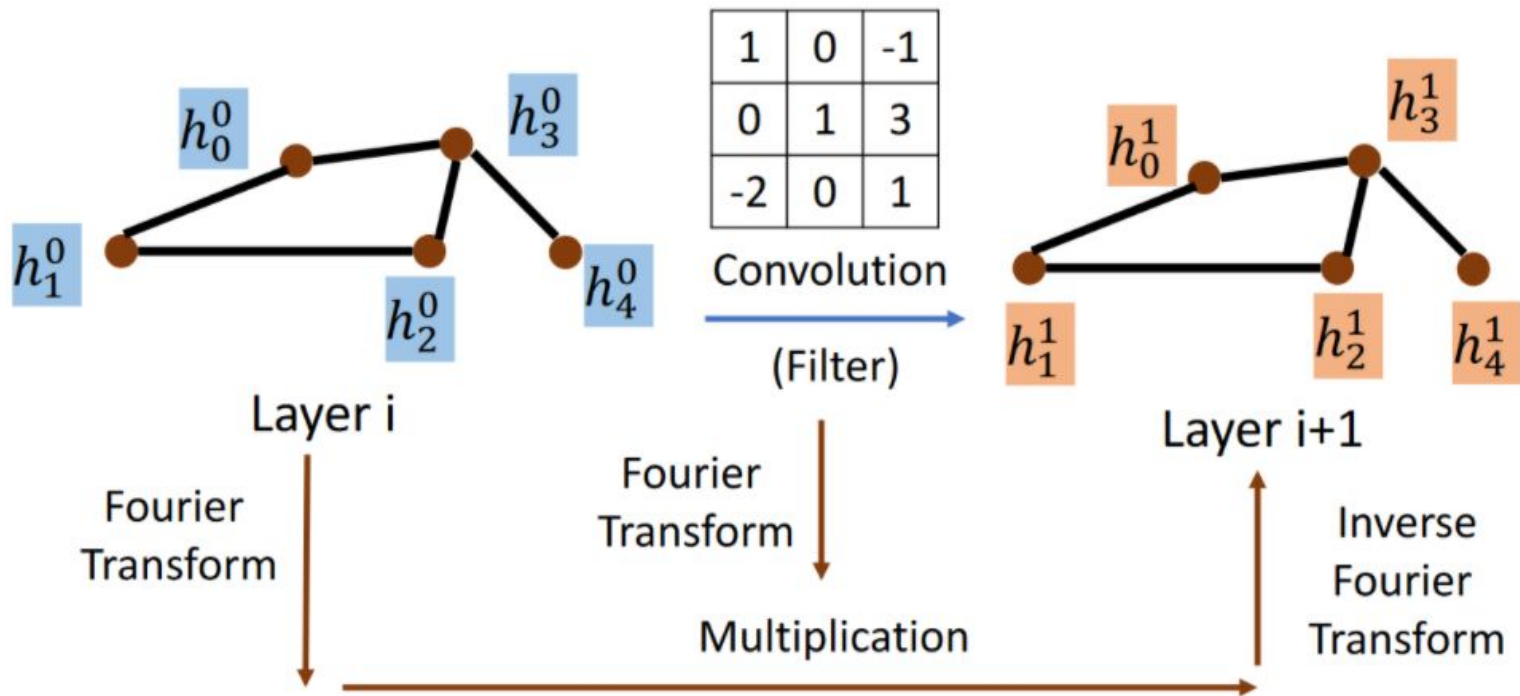
- We try to understand the graph embedding model from a new perspective and propose an attack framework: GF-Attack, which can perform adversarial attack on various kinds of graph embedding models.
- We formulate the graph embedding model as a general graph signal processing with corresponding graph filter which can be computed by the input adjacency matrix.
- GF-Attack is capable to perform the adversarial attack on any graph embedding models which can be formulate to a general graph signal processing.



Preliminary

Spectral-Based Convolution

將Graph Data與Graph Filter都做Fourier Transform, 也就是將Graph Data與Graph Filter從Spatial Domain轉換至Frequency Domain, 如此一來對要Graph Data做Convolution直接相乘就好, 比較簡單, 做完後再將Graph Data做Inverse Fourier Transform, 使其從Frequency Domain變回Spatial Domain。



Spectral Graph Theory

- 考慮一無向圖, Graph Signal通常反應在Graph Node上(ie. Node Feature)

Graph: $G = (V, E)$, $N = |V|$

$A \in \mathbb{R}^{N \times N}$, adjacency matrix (weight matrix).

$$A_{i,j} = 0 \text{ if } e_{i,j} \notin E, \text{ else } A_{i,j} = w(i,j)$$

We only consider undirected graph

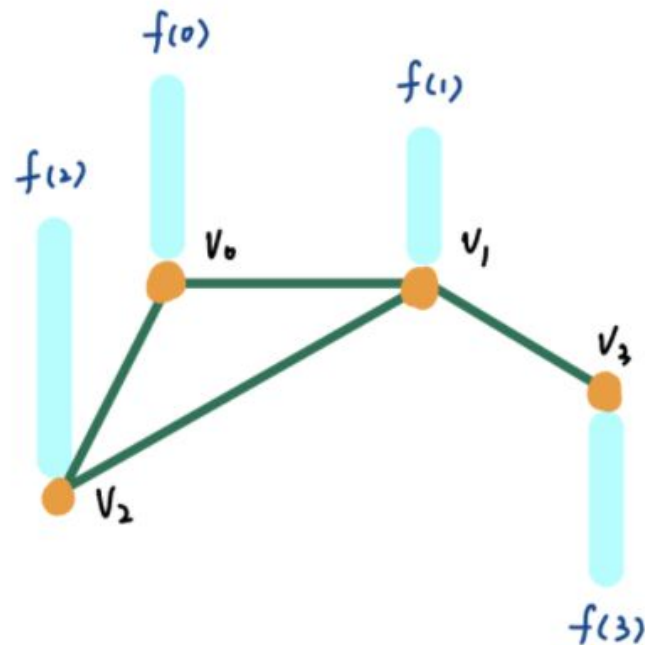
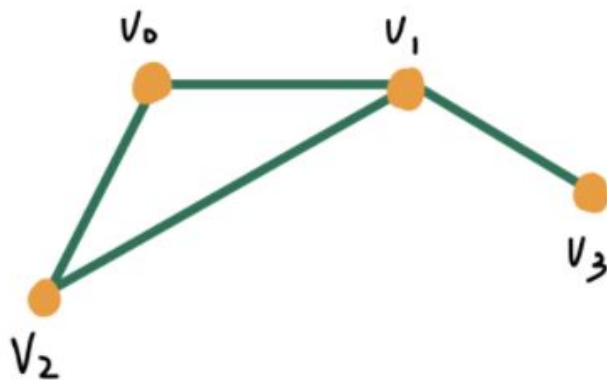
$D \in \mathbb{R}^{N \times N}$, degree matrix

$$D_{i,j} = \begin{cases} d(i) & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \quad (\text{Sum of row } i \text{ in } A)$$

$f: V \rightarrow \mathbb{R}^N$, signal on graph (vertices). $f(i)$ denotes the signal on vertex i

Spectral Graph Theory

- 假設一簡單圖 G 有 v_0 、 v_1 、 v_2 、 v_3 等節點， $f(0)$ 、 $f(1)$ 、 $f(2)$ 、 $f(3)$ 指節點訊號



Spectral Graph Theory

- 定義Laplacian Matrix, 屬於半正定矩陣, 對L特徵分解則可以得到Eigen-Vector U和Eigen-Value λ 。下圖將Eigenvalue定義為Frequency, Eigenvector定義為Basis, 也就是說L可以作為一運算子, 對Graph Signal進行處理

Graph Laplacian $L = D - A$, $L \succcurlyeq 0$ (Positive semidefinite)

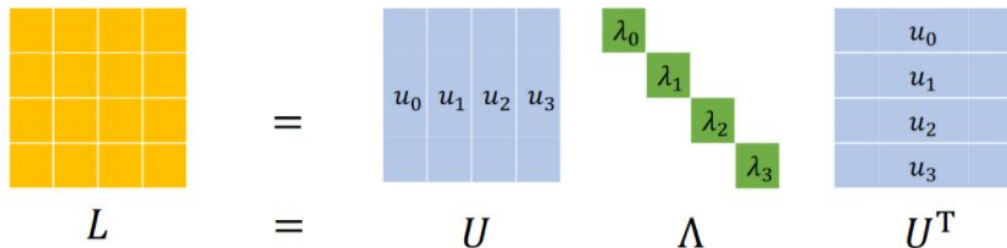
L is symmetric (for undirected graph)

$L = U\Lambda U^T$ (spectral decomposition)

$\Lambda = \text{diag}(\lambda_0, \dots, \lambda_{N-1}) \in \mathbb{R}^{N \times N}$

$U = [u_0, \dots, u_{N-1}] \in \mathbb{R}^{N \times N}$, orthonormal

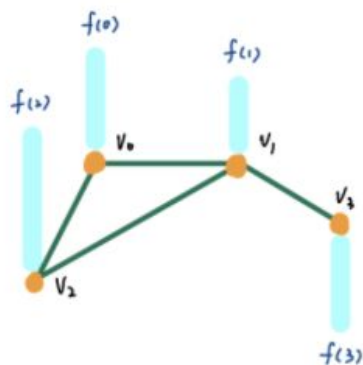
λ_l is the frequency, u_l is the basis corresponding to λ_l



Spectral Graph Theory

- 根據G, 可以得到下圖矩陣, f為事先定義之Signal, 對L特徵分解可以得到對應的 Eigenvalue和Eigenvector

✓ Vertex domain signal



$$f = \begin{bmatrix} 4 \\ 2 \\ 4 \\ -3 \end{bmatrix}$$

$$D = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

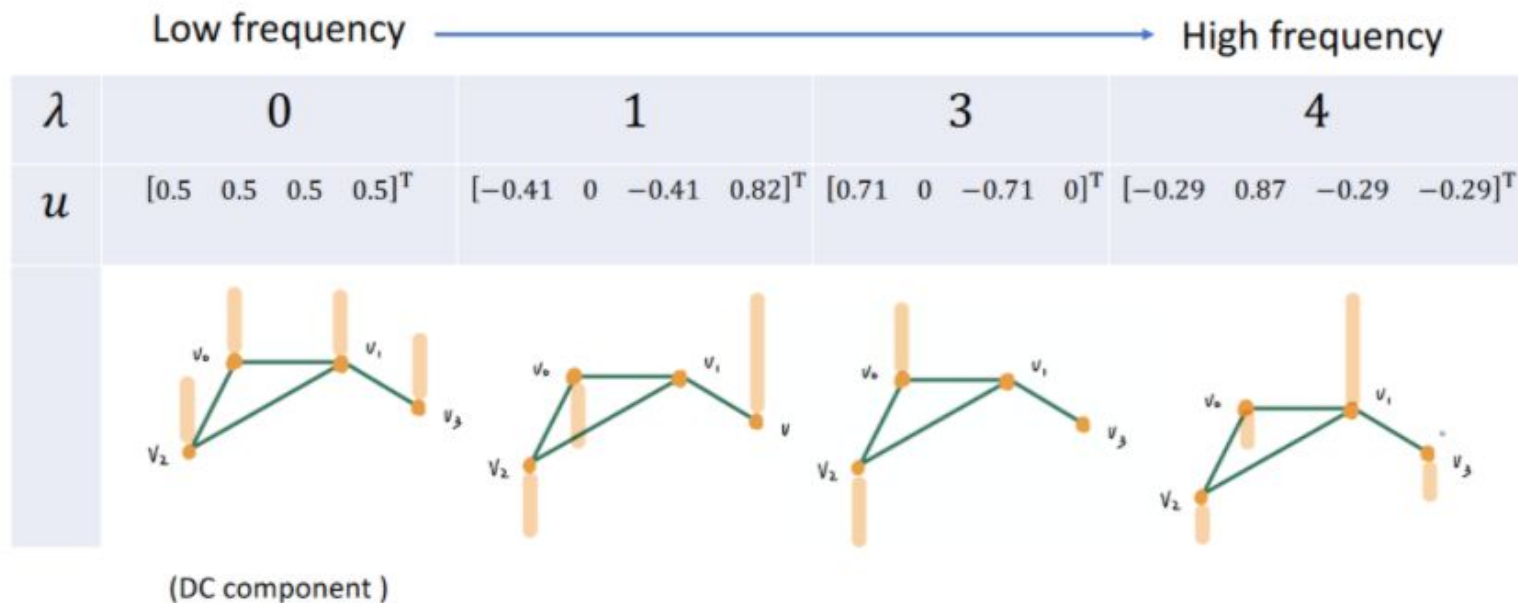
$$L = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 2 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix}$$

$$\Lambda = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 4 \end{bmatrix}$$

$$U = \begin{bmatrix} 0.5 & -0.41 & 0.71 & -0.29 \\ 0.5 & 0 & 0 & 0.87 \\ 0.5 & -0.41 & -0.71 & -0.29 \\ 0.5 & 0.82 & 0 & -0.29 \end{bmatrix}$$

Spectral Graph Theory

- 若將Eigenvalue視為Frequency，則不同Eigenvalue會對應不同Eigenvector，而它能被視為影響Signal大小的Transform



Spectral Graph Theory

- L 的物理意義在於作為一運算子，能考量不同鄰居間的Node Signal進行運算。以下圖 V_0 節點的Signal為例，它會需要考慮 V_1 和 V_2 的Signal影響

Spectral Graph Theory

✓ Interpreting vertex frequency

- L as an operator on graph
- Given a graph signal f , what does Lf mean?
- $Lf = (D - A)f = Df - Af$

$$D = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$A = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

$$f = \begin{bmatrix} 4 \\ 2 \\ 4 \\ -3 \end{bmatrix}$$

$$Lf = \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}$$

➤ Lets focus on the first row of Lf

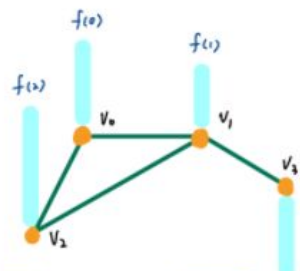
$$➤ a = [2 \ 0 \ 0 \ 0] \cdot [4 \ 2 \ 4 \ -3] - [0 \ 1 \ 1 \ 0] \cdot [4 \ 2 \ 4 \ -3]$$

of neighbors of v_0 Signal on v_0 's neighbors

$$= 2 \times 4 - 2 - 4 = (4 - 2) + (4 - 4) = 2$$

Signal on v_0

Sum of difference between v_0 and its neighbors



Spectral Graph Theory

- 受到其他Signal影響可以定義為下圖公式，也就是計算Signal之間的平滑程度

✓ $(Lf)(v_i) = \sum_{v_j \in V} w_{i,j} (f(v_i) - f(v_j))$, where $w_{i,j}$ is the $(i,j)^{th}$ entry of A

$$\begin{aligned} f^T L f &= \sum_{v_i \in V} f(v_i) \sum_{v_j \in V} w_{i,j} (f(v_i) - f(v_j)) \\ &= \sum_{v_i \in V} \sum_{v_j \in V} w_{i,j} (f^2(v_i) - f(v_i)f(v_j)) \\ &= \frac{1}{2} \sum_{v_i \in V} \sum_{v_j \in V} w_{i,j} (f^2(v_i) - f(v_i)f(v_j) + f^2(v_j) - f(v_j)f(v_i)) \\ &= \frac{1}{2} \sum_{v_i \in V} \sum_{v_j \in V} w_{i,j} (f(v_i) - f(v_j))^2 \end{aligned}$$

“Power” of signal variation
between nodes, i.e.,
smoothness of graph signal

Spectral Graph Theory

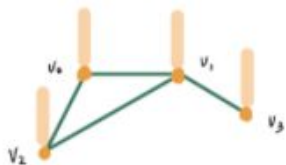
- 根據上述公式，可以將Eigenvalue視為Frequency

✓ $u_i^T L u_i = u_i^T \lambda_i u_i = \lambda u_i^T u_i = \lambda_i$

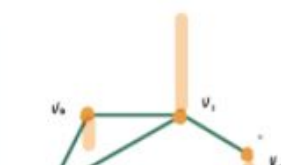
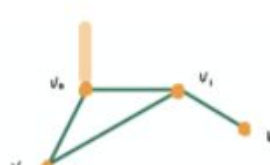
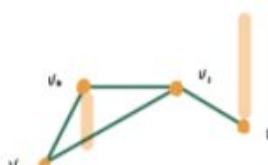
- ✓ The eigenvectors corresponding to small λ belong to the low-pass part of a graph signal

Low frequency \longrightarrow High frequency

λ	0	1	3	4
u	$[0.5 \ 0.5 \ 0.5 \ 0.5]^T$	$[-0.41 \ 0 \ -0.41 \ 0.82]^T$	$[0.71 \ 0 \ -0.71 \ 0]^T$	$[-0.29 \ 0.87 \ -0.29 \ -0.29]^T$



(DC component)



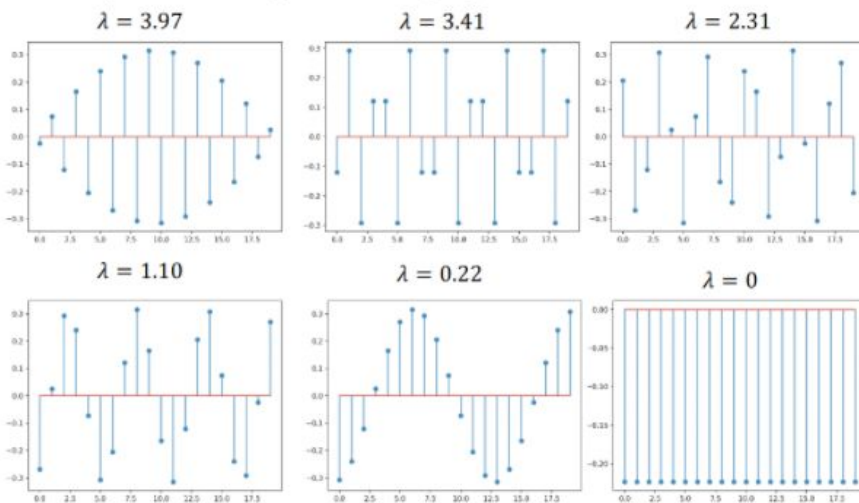
Spectral Graph Theory

- 以下圖證明Eigenvalue擁有Frequency特性

Spectral Graph Theory



✓ A special example: a line graph with 20 nodes



Graph Fourier Transform

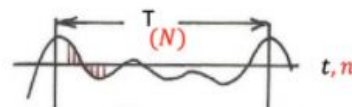
- 應用Fourier Transform的概念，將Vertex Domain(Spatial Domain)的訊號，透過Eigenvector Transpose投射到Spectral Domain上

✓ Graph Fourier Transform of signal x : $\hat{x} = U^T x$

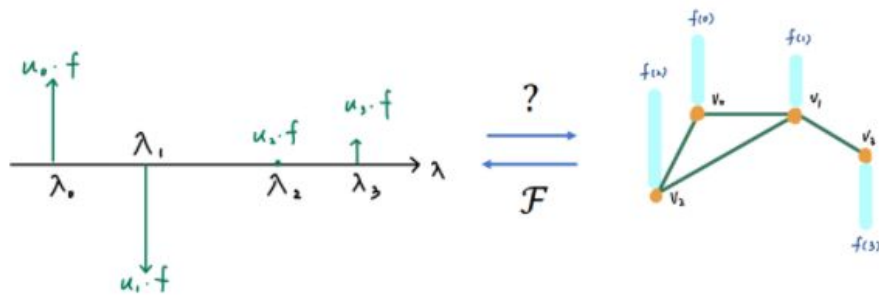
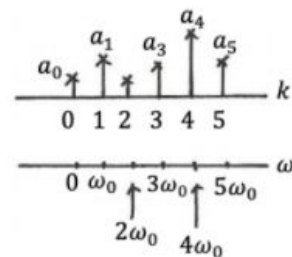
$$\begin{array}{c}
 \text{分析} \\
 \text{分析} \\
 \text{分析} \\
 \text{分析}
 \end{array}
 \begin{array}{c}
 u_0 \cdot x \\
 u_1 \cdot x \\
 u_2 \cdot x \\
 u_3 \cdot x
 \end{array}
 =
 \begin{array}{c}
 u_0 \\
 u_1 \\
 u_2 \\
 u_3
 \end{array}
 \begin{array}{c}
 x
 \end{array}
 \begin{array}{c}
 \text{合成}
 \end{array}$$

\hat{x} $=$ U^T x

Spectral domain Vertex domain



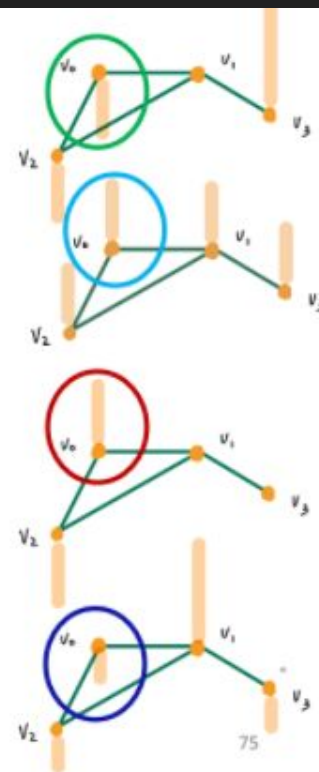
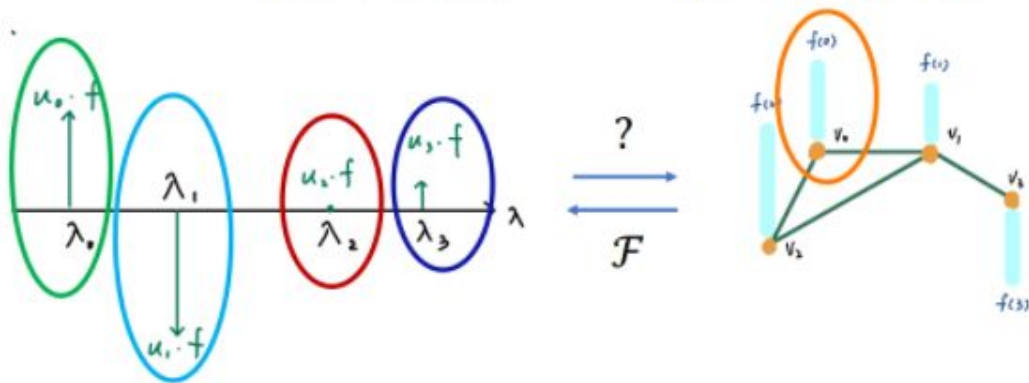
\updownarrow
 \mathcal{F}



Inverse Graph Fourier Transform

✓ Inverse Graph Fourier Transform of signal \hat{x} : $x = U\hat{x}$

$$\begin{array}{ccc}
 \begin{array}{|c|} \hline \text{orange} \\ \hline \end{array} & = & \begin{array}{|c|c|c|c|} \hline \text{green} & \text{cyan} & \text{red} & \text{blue} \\ \hline u_0 & u_1 & u_2 & u_3 \\ \hline \end{array} \begin{array}{|c|} \hline u_0 \cdot x \\ \hline u_1 \cdot x \\ \hline u_2 \cdot x \\ \hline u_3 \cdot x \\ \hline \end{array} \\
 x & = & U \hat{x} \quad \text{合成} \\
 \text{Vertex domain} & & \text{Spectral domain}
 \end{array}$$



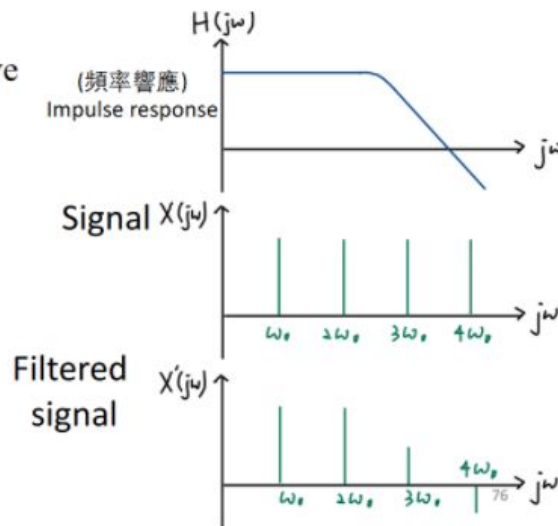
Filtering

- 針對不同頻率間的訊號進行過濾

Filtering

modifying the amplitude/ phase of the different frequency components in a signal, including eliminating some frequency components entirely

- frequency shaping, frequency selective



Filtering

- 定義一 $g(\theta)$ function, 作為 Convolution Function, 讓 Spectral Domain Signal 進行捲積

Spectral Graph Theory

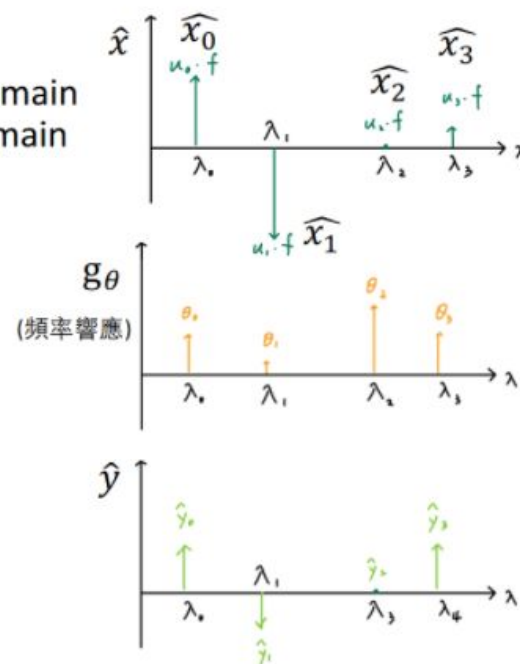
✓ Filtering: Convolution in time domain is multiplication in frequency domain

✓ $\hat{y} = g_{\theta}(\Lambda) \hat{x}$

$$\begin{bmatrix} \hat{y}_0 \\ \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} \theta_0 & & & \\ & \theta_1 & & \\ & & \theta_2 & \\ & & & \theta_3 \end{bmatrix} \begin{bmatrix} \hat{x}_0 \\ \hat{x}_1 \\ \hat{x}_2 \\ \hat{x}_3 \end{bmatrix}$$

$$\hat{y} = g_{\theta}(\Lambda) \hat{x}$$

$$g_{\theta}(\lambda_i) = \theta_i$$



Inverse Graph Fourier Transform

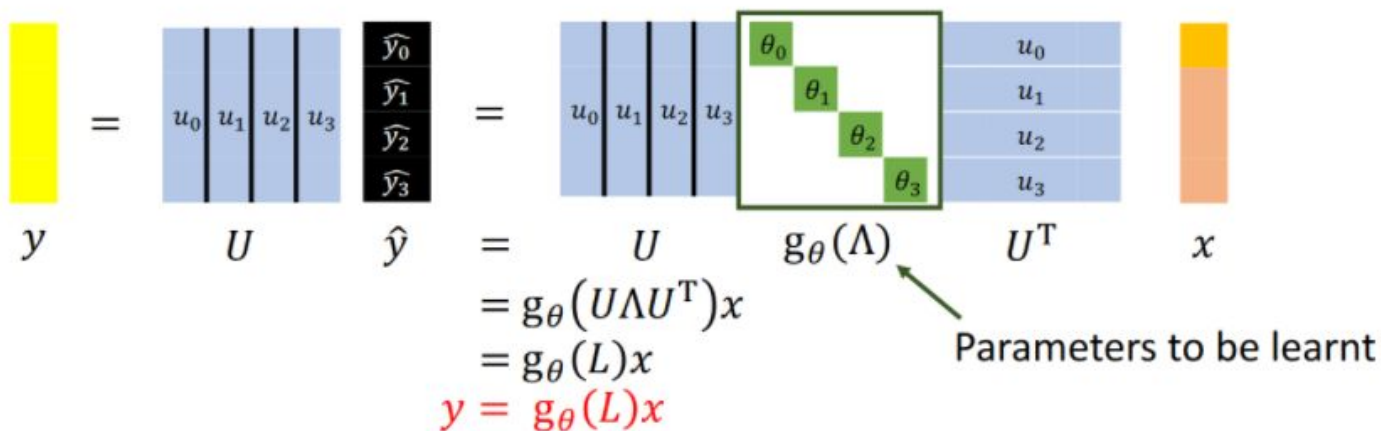
$$\begin{array}{c}
 \begin{array}{c} \hat{y}_0 \\ \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{array} = \begin{array}{c} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{array} \begin{array}{c} u_0 \\ u_1 \\ u_2 \\ u_3 \end{array} x \\
 \hat{y} = g_{\theta}(\Lambda) U^T x
 \end{array}$$

$$\begin{array}{c}
 \begin{array}{c} u_0 \\ u_1 \\ u_2 \\ u_3 \end{array} \begin{array}{c} \hat{y}_0 \\ \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{array} = \begin{array}{c} u_0 \\ u_1 \\ u_2 \\ u_3 \end{array} \begin{array}{c} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{array} \begin{array}{c} u_0 \\ u_1 \\ u_2 \\ u_3 \end{array} x \\
 U \hat{y} = U g_{\theta}(\Lambda) U^T x
 \end{array}$$

Inverse Fourier Transform

Spectral Graph Theory

- 最後得到下圖式子，但這樣會遇到第一個問題，整體會受限於Input Data的N維大小，效率為 $O(N)$



$$\begin{aligned}
 y &= U \hat{y} = U U^T x g_{\theta}(\Lambda) \\
 &= g_{\theta}(U \Lambda U^T) x \\
 &= g_{\theta}(L) x \\
 \text{Parameters to be learnt} &\rightarrow g_{\theta}(\Lambda)
 \end{aligned}$$

$y = g_{\theta}(L)x$

$g_{\theta}(\cdot)$ can be any function. For example,

$$g_{\theta}(L) = \log(I + L) = L - \frac{L^2}{2} + \frac{L^3}{3} \dots, \lambda_{\max} < 1$$

Problem 1: Learning complexity is $O(N)$!!!

Spectral Graph Theory

- 但若考慮 $g = L^N$ function, 則會考慮到 N 個 Node 的影響, 因此缺乏 Localize 特性
 - ✓ If we select $g_\theta(L) = \cos(L) = I - \frac{L^2}{2!} + \frac{L^4}{4!} \dots$
 - ✓ If a connected graph has N nodes, then L^N will make the all nodes be able to share their signals with each other

Problems 2: Not localize!!!

ChebNet

- 將上述g function引入ChebNet機制，但會產生第三個問題為高Time Complexity

✓ Solution to Problem 1 and 2:

➤ Use polynomial to parametrize $g_{\theta}(L)$

$$g_{\theta}(L) = \sum_{k=0}^K \theta_k L^k$$

Now it is K-localized

$$g_{\theta}(\Lambda) = \sum_{k=0}^K \theta_k \Lambda^k$$

Parameters to be learnt: $O(K)$

Problem 3:
Time complexity: $O(N^2)$

$$y = U g_{\theta}(\Lambda) U^T x = U \left(\sum_{k=0}^K \theta_k \Lambda^k \right) U^T x$$



GCN

GCN

$$\checkmark y = g_{\theta'}(L)x = \sum_{k=0}^K \theta'_k T_k(\tilde{L})x, K = 1$$

<https://openreview.net/pdf?id=SJU4ayYgl>

$$y = g_{\theta'}(L)x = \theta'_0 x + \theta'_1 \tilde{L}x \quad \because \tilde{L} = \frac{2L}{\lambda_{\max}} - I$$

$$= \theta'_0 x + \theta'_1 \left(\frac{2L}{\lambda_{\max}} - I \right) x \quad \because \lambda_{\max} \approx 2$$

$$= \theta'_0 x + \theta'_1 (L - I)x \quad \because L = I - D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$$

$$= \theta'_0 x - \theta'_1 (D^{-\frac{1}{2}} A D^{-\frac{1}{2}}) x \quad \because \theta = \theta'_0 = -\theta'_1$$

$$= \theta (I + D^{-\frac{1}{2}} A D^{-\frac{1}{2}}) x$$

$$\text{renormalization trick: } I_N + D^{-\frac{1}{2}} A D^{-\frac{1}{2}} \rightarrow \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}$$

$$H^{(l+1)} = \sigma \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)} \right)$$

GCN

$$H^{(l+1)} = \sigma\left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)}\right)$$

Can be rewritten as:

$$h_v = f\left(\frac{1}{|\mathcal{N}(v)|} \sum_{u \in \mathcal{N}(v)} W x_u + b\right), \quad \forall v \in \mathcal{V}.$$

Methodology

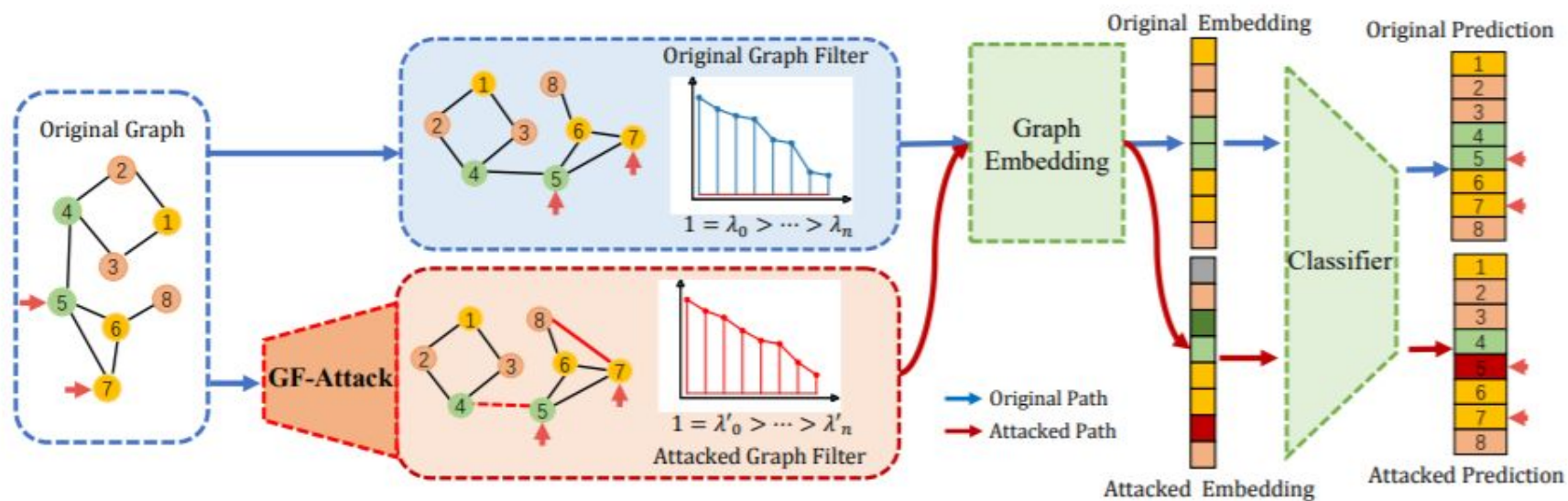


GF Attack

- Attack on \mathcal{V} : Add/delete vertices in graphs. This operation may change the dimension of the adjacency matrix A .
- Attack on A : Add/delete edges in graphs. This operation would lead to the changes of entries in the adjacency matrix A . This kind of attack is also known as *structural attack*.
- Attack on X : Modify the attributes attached on vertices.

Here, we mainly focus on adversarial attacks on graph structure A , since attacking A is more practical than others in real applications (Tong et al. 2012).

Architecture



Adversarial Attack Definition

- Adversarial Attack數學表示, β 指budget

$$\arg \max_{A'} \mathcal{L}(A', Z) \quad (1)$$

$$\text{s.t. } Z = \mathcal{M}_{\Theta}(A', X),$$

$$\Theta^* = \arg \min_{\Theta} \mathcal{L}(\Theta; A', X), \|A' - A\| = 2\beta,$$

Methodologies

- Graph Signals According to Graph Filter H together with feature transformation.

$$\tilde{X} = \mathcal{H}(X), X' = \sigma(\tilde{X}\Theta), \quad (2)$$

- 攻擊目標為quality of the output embedding Z, S'指干擾後的Graph Shift Filter

$$\mathcal{L}(A', Z) = \|h(S')X - h(S')_T X\|_F^2, \quad (3)$$

- 根據Low-rank Approximation可得下方式子

$$\mathcal{L}(A', Z) = \left\| \sum_{i=T+1}^n \lambda'_i \mathbf{u}_i \mathbf{u}_i^T X \right\|_F \leq \sum_{i=T+1}^n \lambda_i'^2 \cdot \sum_{i=T+1}^n \|\mathbf{u}_i^T X\|_2^2, \quad (4)$$

Methodologies

- 以Upper Bound作為最後 Loss function

$$\begin{aligned} \arg \max_{A'} \quad & \sum_{i=T+1}^n \lambda_i'^2 \cdot \sum_{i=T+1}^n \|\mathbf{u}_i^T X\|_2^2, \\ \text{s.t.} \quad & \|A' - A\| = 2\beta. \end{aligned} \tag{5}$$

Algorithm

Algorithm 1 Graph Filter Attack (GF-Attack) adversarial attack algorithm under RBA setting


Input:

Adjacent Matrix A ; feature matrix X ; target vertex t ;
number of top- T smallest singular values/vectors selected T ;
order of graph filter K ; fixed budget β .

Output:

Perturbed adjacent Matrix A' .

- 1: Initial the candidate flips set as $\mathcal{C} = \{(v, t) | v \neq t\}$, eigenvalue decomposition of $\hat{A} = U_{\hat{A}} \Lambda_{\hat{A}} U_{\hat{A}}^T$;
 - 2: **for** $(v, t) \in \mathcal{C}$ **do**
 - 3: Approximate $\Lambda'_{\hat{A}}$ resulting by removing/inserting edge (v, t) via Equation (9);
 - 4: Update $Score_{(v, t)}$ from loss Equation (8) or Equation (12);
 - 5: **end for**
 - 6: $\mathcal{C}_{sel} \leftarrow$ edge flips with top- β $Score$;
 - 7: $A' \leftarrow A \pm \mathcal{C}_{sel}$;
 - 8: **return** A'
-



Experiment
Result

Dataset

scientific publications

links

unique words

classes

Dataset

	Cora	Citeseer	Pubmed
Task	Transductive	Transductive	Transductive
# Nodes	2708 (1 graph)	3327 (1 graph)	19717 (1 graph)
# Edges	5429	4732	44338
# Features/Node	1433	3703	500
# Classes	7	6	3
# Training Nodes	140	120	60
# Validation Nodes	500	500	500
# Test Nodes	1000	1000	1000

Steps

train GCN

GF Attack

result

Our result

Cora

The test accuracy performance is: 0.8390

Final accuracy after attack is: 0.8109

Attack performance: -0.0281

Citeseer

The test accuracy performance is: 0.7465

Final accuracy after attack is: 0.73

Attack performance: -0.0165

Table 1: Summary of the change in classification accuracy (in percent) compared to the clean/original graph. Single edge perturbation under RBA setting. Lower is better.

Dataset	Cora				Citeseer				Pubmed			
Models (unattacked)	GCN	SGC	DeepWalk	LINE	GCN	SGC	DeepWalk	LINE	GCN	SGC	DeepWalk	LINE
<i>Random</i>	80.20	78.82	77.23	76.75	72.50	69.68	69.68	65.15	80.40	80.21	78.69	72.12
<i>Degree</i>	-1.90	-1.22	-1.76	-1.84	-2.86	-1.47	-6.62	-1.78	-1.75	-1.77	-1.25	-1.01
<i>RL-S2V</i>	-2.21	-4.42	-3.08	-12.40	-4.68	-5.21	-9.67	-12.55	-3.86	-4.44	-2.43	-13.05
\mathcal{A}_{class}	-5.20	-5.62	-5.24	-10.38	-6.50	-4.08	-12.13	-20.10	-6.40	-6.11	-6.10	-13.21
	-3.62	-2.96	-6.29	-7.55	-3.48	-2.83	-12.56	-10.28	-4.21	-2.25	-3.05	-6.75
<i>GF-Attack</i>	-7.60	-9.73	-5.31	-13.27	-7.78	-6.19	-12.50	-22.11	-7.96	-7.20	-7.43	-14.16

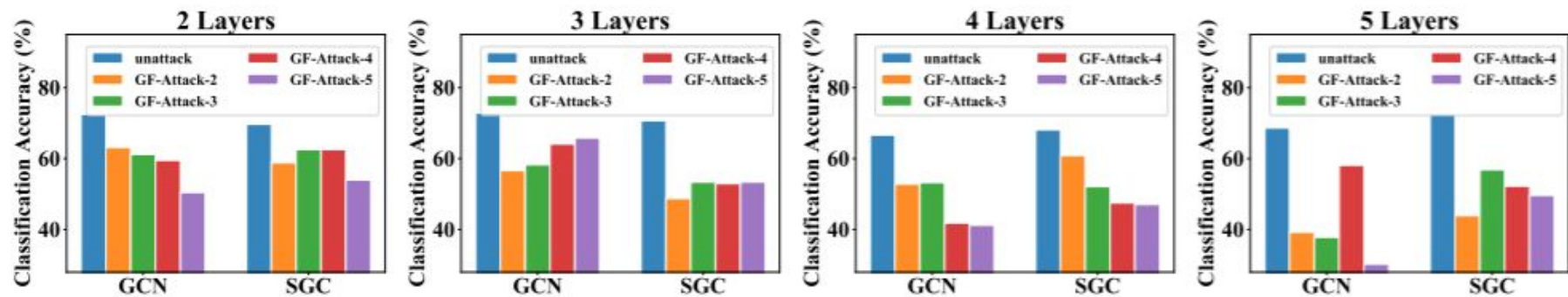
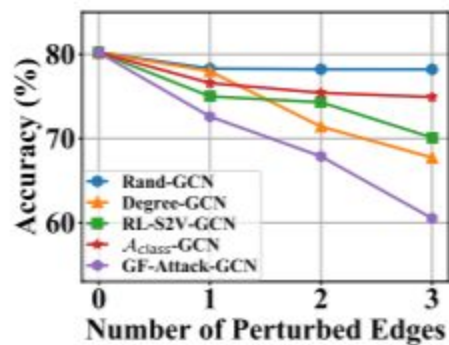


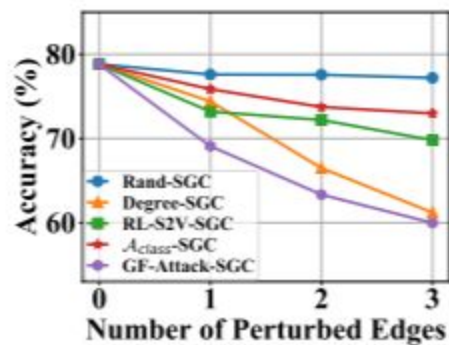
Figure 2: Comparison between order K of *GF-Attack* and number of layers in GCN/SGC on Citeseer.

Table 2: Running time (s) comparison over all baseline methods on Citeseer. We report the 10 times average running time of processing single node for each model.

Models	<i>Random</i>	<i>Degree</i>	<i>RL-S2V</i>	\mathcal{A}_{class}	<i>GF-Attack</i>
Citeseer	0.19	42.21	222.80	146.58	12.78



(a) GCN



(b) SGC

Figure 3: Multiple-edge attack results on Cora under RBA setting. Lower is better.