

Graph Neural Networks (GNN)

Makoto Yamada

`myamada@i.kyoto-u.ac.jp`

Kyoto University

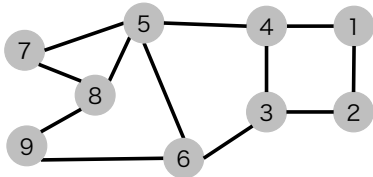
2020/7/20

Graph

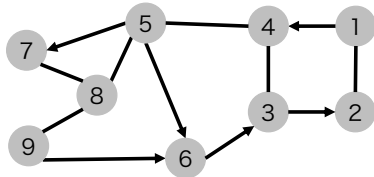
Graph

A graph consists of Nodes V and edges E . $G = (V, E)$

Undirected graph



Directed graph



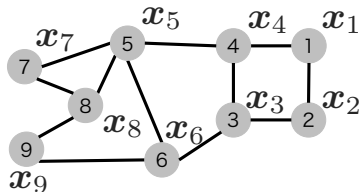
Attributed graph

Graph + features

In each node, a feature vector is given. $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$

We mainly use **Attributed graphs** for GNN.

Attributed graph



Adjacency matrix & Degree matrix

Adjacency matrix A

$n \times n$ matrix with **edge** information

$$a_{ij} = \begin{cases} 1 & (\text{i-th node is connected to j-th node}) \\ 0 & (\text{Otherwise}) \end{cases}$$

Degree matrix D

$n \times n$ matrix with **degree** information.

$$d_{ii} = \sum_{k=1}^n a_{ik}$$

$$D = \text{diag}(d_{11}, d_{22}, \dots, d_{nn}) \in \mathbb{R}^{n \times n}$$

Graph Laplacian

Graph Laplacian

$$L = D - A,$$

- $\mathbf{u}^\top L \mathbf{u} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} (u_i - u_j)^2 \geq 0$ (PSD)
- $L \mathbf{1}_n = \mathbf{0}_n$ (minimum eigenvalue is zero)

Normalized graph Laplacian

$$\bar{L} = D^{-\frac{1}{2}} L D^{-\frac{1}{2}} = I_n - D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$$

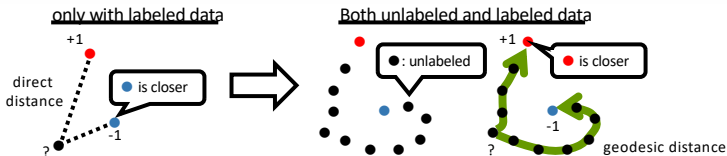
- $\mathbf{u}^\top \bar{L} \mathbf{u} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} \left(\frac{u_i}{\sqrt{d_{ii}}} - \frac{u_j}{\sqrt{d_{jj}}} \right)^2 \geq 0$ (PSD)
($\tilde{\mathbf{u}} = D^{-1/2} \mathbf{u}$ and $\mathbf{x}^\top \bar{L} \mathbf{u} = \tilde{\mathbf{u}}^\top L \tilde{\mathbf{u}}$)
- $\bar{L} D^{1/2} \mathbf{1}_n = \bar{L} \mathbf{d} = \mathbf{0}_n$ (minimum and maximum eigenvalues are 0 and 2, respectively)

Graph based Machine Learning

We construct graphs and use it for prediction.

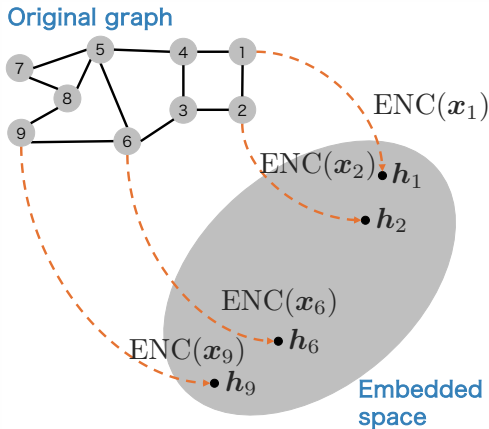
Graph stores the adjacent information between samples.

- K nearest neighbor graph ($A_{ij} = 0, 1$)
- Weighted graph ($A_{ij} = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2)$)



What is GNN?

It is an **embedding method** for graph.



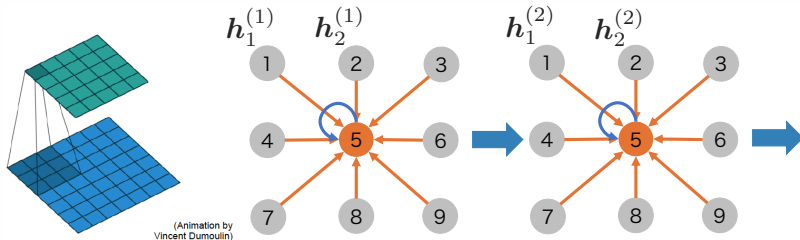
Convolution

Convolutional Neural Networks (CNN)

- Aggregate neighbor pixels.
- Apply $\sigma(\cdot)$ and propagate the information to higher layers.

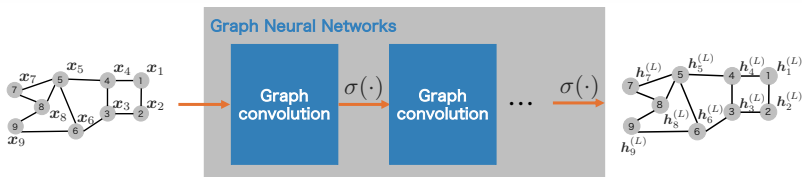
e.g., Convolution for 5th pixel.

$$h_5^{(\ell+1)} = \sigma \left(W_1^{(\ell)} h_1^{(\ell)} + W_2^{(\ell)} h_2^{(\ell)} + \dots + W_9^{(\ell)} h_9^{(\ell)} \right)$$



Graph Neural Networks (GNN)

GNN: **Neural network with graph convolution.**



Task $z_i = Wh_i^{(L)}$, $W \in \mathbb{R}^{K \times d_L}$

- Node classification : $\text{softmax}(z_i)$
- Graph classification : $\text{softmax}(\sum_i z_i)$
- Link prediction : $p(a_{ij}) = \sigma(z_i^\top z_j)$

GNN (Node classification)

We compute probability by using **an additional layer** ($\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_K)^\top \in \mathbb{R}^{K \times d_L}$) to the graph embedding $\mathbf{h}_i^{(L)}$

$$\mathbf{z}_i = \text{softmax}(\mathbf{W} \mathbf{h}_i^{(L)}) = \begin{pmatrix} \frac{\exp(\mathbf{w}_1^\top \mathbf{h}_i^{(L)})}{\sum_{k'=1}^K \exp(\mathbf{w}_{k'}^\top \mathbf{h}_i^{(L)})} \\ \vdots \\ \frac{\exp(\mathbf{w}_K^\top \mathbf{h}_i^{(L)})}{\sum_{k'=1}^K \exp(\mathbf{w}_{k'}^\top \mathbf{h}_i^{(L)})} \end{pmatrix}$$

Loss function (Training data: $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, $\mathbf{A} \in \mathbb{R}^{n \times n}$)

$$\mathcal{L} = \sum_{i=1}^n \sum_{k=1}^K y_{k,i} \log(z_{k,i})$$

$\mathbf{z}_i = (z_{1,i}, \dots, z_{K,i})^\top$, \mathbf{y}_i is an one-hot vector.

GNN (Graph classification)

We transform the output of graph convolution $\mathbf{h}_i^{(L)}$ by an additional layer and **aggregate them** to compute probability.

$$\mathbf{z}_i = \text{softmax}(\mathbf{W} \sum_i \mathbf{h}_i^{(L)}) = \begin{pmatrix} \frac{\exp(\mathbf{w}_1^\top \sum_i \mathbf{h}_i^{(L)})}{\sum_{k'=1}^K \exp(\mathbf{w}_{k'}^\top \sum_i \mathbf{h}_i^{(L)})} \\ \vdots \\ \frac{\exp(\mathbf{w}_K^\top \sum_i \mathbf{h}_i^{(L)})}{\sum_{k'=1}^K \exp(\mathbf{w}_{k'}^\top \sum_i \mathbf{h}_i^{(L)})} \end{pmatrix}$$

Loss function (Training data: $\{(\mathbf{X}_i, \mathbf{y}_i, \mathbf{A}_i)\}_{i=1}^n$)

$$\mathcal{L} = \sum_{i=1}^n \sum_{k=1}^K y_{k,i} \log(z_{k,i})$$

$\mathbf{z}_i = (z_{1,i}, \dots, z_{K,i})^\top$, \mathbf{y}_i is a one-hot vector.

GNN (Link prediction)

The output of an GNN: $\mathbf{z}_i = \mathbf{h}_i^{(L)}$

Generative model

$$p(\mathbf{A}|\mathbf{Z}) = \prod_{i=1}^n \prod_{j=1}^n p(a_{ij}|\mathbf{z}_i, \mathbf{z}_j),$$
$$p(a_{ij} = 1|\mathbf{z}_i, \mathbf{z}_j) = \sigma(\mathbf{z}_i^\top \mathbf{z}_j)$$

$\sigma(\cdot)$ is the sigmoid function.

Loss function ($\mathbf{X} \in \mathbb{R}^{d \times n}$, $\mathbf{A} \in \mathbb{R}^{n \times n}$)

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{Z}|\mathbf{X}, \mathbf{A})}[\log p(\mathbf{A}|\mathbf{Z})] - \text{KL}[q(\mathbf{Z}|\mathbf{X}, \mathbf{A})||p(\mathbf{Z})]$$

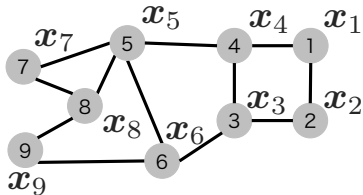
$\text{KL}[q(\cdot)||p(\cdot)]$ is the KL-divergence.

Graph Convolutional Networks (GCN)

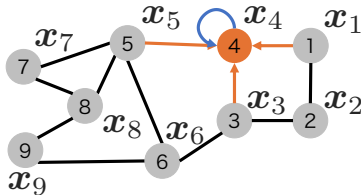
Graph convolution

$$h_i^{(\ell+1)} = \sigma \left(W_0^{(\ell)} h_i^{(\ell)} + \sum_{j: a_{i,j} \neq 0} \frac{1}{c_{i,j}} W_1^{(\ell)} h_j^{(\ell)} \right)$$

Attributed graph



Graph convolution



Graph Convolutional Networks (GCN)

Graph convolution

$$\mathbf{h}_i^{(\ell+1)} = \sigma \left(\mathbf{W}_0^{(\ell)} \mathbf{h}_i^{(\ell)} + \sum_{j: a_{i,j} \neq 0} \frac{1}{c_{i,j}} \mathbf{W}_1^{(\ell)} \mathbf{h}_j^{(\ell)} \right)$$

Definition of variables

- $\mathbf{A} \in \mathbb{R}^{n \times n}$: Adjacency matrix
- $\mathbf{W}_1^{(\ell)} \in \mathbb{R}^{m_\ell \times m_{\ell-1}}$: Transformation matrix
- $\mathbf{h}_i^{(\ell)} \in \mathbb{R}^{m_\ell}$: i -th sample's intermediate representation of the ℓ -th layer
- $\sigma(\cdot)$: Activation function

Explain about graph convolution through **graph Fourier transform**.

Fourier transform

Transform real signal to frequencies.

Discrete Fourier transform

$$X(k) = \sum_{i=0}^{n-1} x(k) e^{\frac{-j2\pi ki}{n}} = \sum_{i=0}^{n-1} x(k) w_n^{(ki)}$$

Matrix representation

$$\begin{pmatrix} X(0) \\ X(1) \\ \vdots \\ X(n-1) \end{pmatrix} = \begin{pmatrix} w_n^{(0)} & w_n^{(0)} & \cdots & w_n^{(0)} \\ w_n^{(1)} & w_n^{(2)} & \cdots & w_n^{(n-1)} \\ \vdots & \vdots & \ddots & \vdots \\ w_n^{(0)} & w_n^{(n-1)} & \cdots & w_n^{(1)} \end{pmatrix} \begin{pmatrix} x(0) \\ x(1) \\ \vdots \\ x(n-1) \end{pmatrix}$$

Low frequency component tends to be constant.

Graph Fourier transform

Definition Eigenvalue decomposition of the (normalized) graph Laplacian.

$$\bar{L} = U\Lambda U^\top$$

Normalized graph Laplacian

$$\bar{L} = D^{-\frac{1}{2}}LD^{-\frac{1}{2}} = I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$$

$$L = D - A,$$

where

$$U = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n)$$

$$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n), 0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

Graph Fourier transform

Eigenvalue decomposition (1st eigenvalue)

Since we have $L\mathbf{1}_n = 0\mathbf{1}_n$ and L is a PSD matrix, $\frac{1}{\sqrt{n}}\mathbf{1}_n$ is the corresponding eigenvector.

Moreover, we can write the eigenvalue decomposition as

$$\begin{aligned} \min_{\mathbf{u}} \quad & \mathbf{u}^\top L \mathbf{u} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} (u_i - u_j)^2 \\ \text{s.t.} \quad & \mathbf{u}^\top \mathbf{u} = 1. \end{aligned}$$

Thus, the minimum is achieved when

$u_1 = u_2 = \dots = u_n = \text{Const.}$ with its eigenvector $\frac{1}{\sqrt{n}}\mathbf{1}_n$. (i.e.,

Low pass filter)

Graph Fourier transform

Eigenvalue decomposition (2nd eigenvalue)

$$\begin{aligned} \min_{\mathbf{u}} \quad & \mathbf{u}^\top \mathbf{L} \mathbf{u} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} (u_i - u_j)^2 \\ \text{s.t.} \quad & \mathbf{u}^\top \mathbf{u} = 1, \mathbf{u}^\top \mathbf{u}_1 = 0. \end{aligned}$$

$\mathbf{u}_1 = \frac{1}{n} \mathbf{1}_n$ is the first eigenvector.

Interpretation

We want to have an eigenvector that is orthogonal to the 1st eigenvector \mathbf{u}_1 and minimizes $\mathbf{u}^\top \mathbf{L} \mathbf{u}$.

That is, eigenvectors with small eigenvalues are **low pass filters** and eigenvectors with high eigenvalues are **high pass filters**.

Graph convolutions (Graph Fourier)

We apply (graph) Fourier transform to the input and then apply filters to the transformed vector. Then, we apply inverse (graph) Fourier transform to the filtered signal.

Graph Fourier transform

$$\tilde{x} = U^\top x$$

Convolution (multiplication in frequency domain.)

$$\Theta \tilde{x} = \Theta U^\top x$$

Inverse Graph Fourier transform

$$x_{\text{filtered}} = U \Theta U^\top x$$

Θ is estimated from training data set.

Graph convolution (ChebyNet)

In graph Fourier transform, it needs to solve eigenvalue decomposition $O(n^3)$.

ChebyNet Compute filter **without** eigenvalue decomposition

$$\Theta(\Lambda) \approx \sum_{k=1}^K \theta_k T_k \left(\frac{2}{\lambda_{\max}} \Lambda - I_n \right)$$

θ_k are the model parameters, $T_k(\Lambda)$ is the k -th Chebychev polynomial.

$$T_k(\Lambda) = 2\Lambda T_{k-1}(\Lambda) - T_{k-2}(\Lambda) (k > 2), T_0(\Lambda) = I_n, T_1(\Lambda) = \Lambda$$

$$\mathbf{x}_{\text{filtered}} = \mathbf{U} \Theta \mathbf{U}^\top \mathbf{x} \approx \left[\sum_{k=1}^K \theta_k T_k \left(\frac{2}{\lambda_{\max}} \mathbf{L} - \mathbf{I}_n \right) \right] \mathbf{x}$$

Graph convolution (GCN)

We use the normalized graph Laplacian ($\lambda_{\max} = 2$) and use the 1st order ChebyNet.

$$\begin{aligned}\mathbf{x}_{\text{filtered}} &= \mathbf{U} [\theta_0 \mathbf{I}_n + \theta_1 (\mathbf{\Lambda} - \mathbf{I}_n)] \mathbf{U}^\top \mathbf{x} \\ &= [\theta_0 \mathbf{U} \mathbf{U}^\top + \theta_1 (\mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top - \mathbf{I}_n)] \mathbf{x} \\ &= [\theta_0 \mathbf{I}_n + \theta_1 (\mathbf{I}_n - \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{\frac{1}{2}} - \mathbf{I}_n)] \mathbf{x} \\ &= \theta_0 \mathbf{x} - \theta_1 \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{\frac{1}{2}} \mathbf{x}\end{aligned}$$

Thus, this can be written more generally as

$$\mathbf{h}_i^{(\ell+1)} = \sigma \left(\mathbf{W}_0^{(\ell)} \mathbf{h}_i^{(\ell)} + \sum_{j: a_{i,j} \neq 0} \frac{1}{c_{i,j}} \mathbf{W}_1^{(\ell)} \mathbf{h}_j^{(\ell)} \right)$$

Graph convolution (matrix form)

We can write the graph convolution by matrix multiplications.

Without normalization

$$\mathbf{H}^{(\ell+1)} = \sigma \left(\mathbf{W}_1^{(\ell)} \mathbf{H}^{(\ell)} \mathbf{A} \right)$$

With normalization

$$\mathbf{H}^{(\ell+1)} = \sigma \left(\mathbf{W}_1^{(\ell)} \mathbf{H}^{(\ell)} \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \right)$$

There exist many variants.

GNN (Node classification)

We compute probability by using **an additional layer** ($\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_K)^\top \in \mathbb{R}^{K \times d_L}$) to the graph embedding $\mathbf{h}_i^{(L)}$

$$\mathbf{z}_i = \text{softmax}(\mathbf{W} \mathbf{h}_i^{(L)}) = \begin{pmatrix} \frac{\exp(\mathbf{w}_1^\top \mathbf{h}_i^{(L)})}{\sum_{k'=1}^K \exp(\mathbf{w}_{k'}^\top \mathbf{h}_i^{(L)})} \\ \vdots \\ \frac{\exp(\mathbf{w}_K^\top \mathbf{h}_i^{(L)})}{\sum_{k'=1}^K \exp(\mathbf{w}_{k'}^\top \mathbf{h}_i^{(L)})} \end{pmatrix}$$

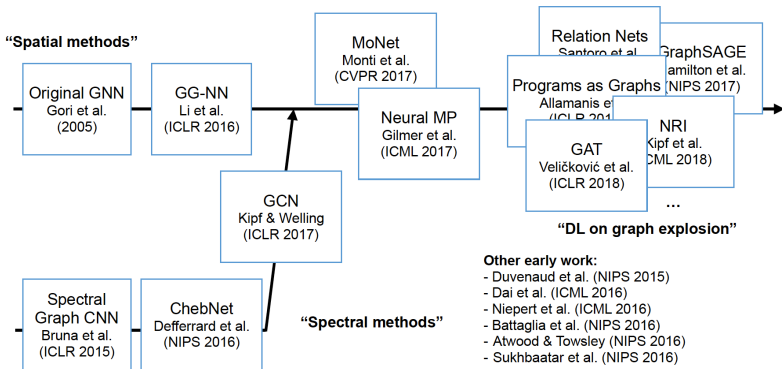
Loss function (Training data: $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, $\mathbf{A} \in \mathbb{R}^{n \times n}$)

$$\mathcal{L} = \sum_{i=1}^n \sum_{k=1}^K y_{k,i} \log(z_{k,i})$$

$\mathbf{z}_i = (z_{1,i}, \dots, z_{K,i})^\top$, \mathbf{y}_i is an one-hot vector.

Various types of GNNs

A brief history of graph neural nets



(slide inspired by Alexander Gaunt's talk on GNNs)

<http://tkipf.github.io/misc/SlidesCambridge.pdf>

Summary

- Graph data
- Introduction to graph neural networks
- Graph Fourier transform
- Graph convolution
- Graph Convolutional Networks (GCN)