# Graph Neural Network

Deep Learning (ENGG\*6600\*01)

School of Engineering,
University of Guelph, ON, Canada
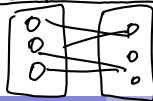
Course Instructor: Benyamin Ghojogh
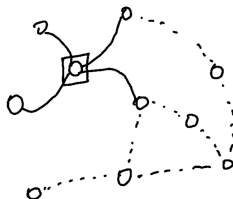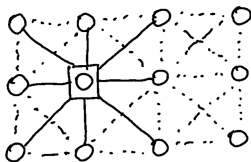Summer 2023

**Introduction**

# Introduction

- Many real-world datasets are in the form of graphs. Some examples are social networks, protein interaction networks, the internet (World Wide Web), molecules, etc.
- Image data can be considered as graph. Every image is a graph where each pixel represents a node (vertex) connected by edges to its adjacent pixels.
- Text data can also be considered as graph. Every token (word) can be a node connected by an edge to its next token (word).
- Tasks in graph processing:
  - Graph-level task: predict the property of an entire graph. Example: predict whether an antibody protein binds to an antigen protein or not.
  - Node-level task: predict the identity or role of every node in the graph. Example: Every node has some features and there is a label for every node. For instance, if the nodes correspond to people, the label can be whether the person lives in a specific city or not.
  - Edge-level task: predict the identity or role of every edge in the graph. Example: In recommender systems for movie suggestion to users, some nodes are the users and some nodes are the movies. An edge between a user and a movie exists if the user has rated that movie and the label of the edge is the rating score. It is possible to predict the label (score) of non-existing edges between a user and a movie.
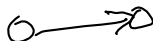
bipartite graph

# Introduction

- As was said, images are special cases of graphs. The graph of an image is called the Euclidean graph or a grid graph. But what if there is a graph with some arbitrary structure or irregular shape.
- In Convolutional Neural Network (CNN) [1], there is convolution of a filter kernel with the image. The question is how to define convolution of a filter kernel with the arbitrary graph.

**Graph Fourier
Transform**

# Laplacian of Graph

- Consider a graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$ with nodes (vertices) $\mathcal{V}$ and edges $\mathcal{E}$.
- Let the number of nodes be $n$. The adjacency matrix $\boldsymbol{A} \in \mathbb{R}^{n \times n}$ is a matrix whose $(i, j)$-th element is one if the node $i$ is connected to the node $j$ and is zero otherwise.
- The degree matrix of the matrix $\boldsymbol{A}$ is a diagonal matrix whose $(i, i)$-th element is the summation of the $i$-th row of the matrix $\boldsymbol{A}$, i.e.:

$$D(i, i) := \sum_{j=1}^{n} \boldsymbol{A}(i, j), \tag{1}$$

  where $\boldsymbol{A}(i, j)$ denotes the $(i, j)$-th element of $\boldsymbol{A}$.

- The Laplacian matrix of the graph $\mathcal{G}$ is defined as:

$$\mathbb{R}^{n \times n} \ni \boldsymbol{L} := \boldsymbol{D} - \boldsymbol{A}. \tag{2}$$

- It is noteworthy that there exist some other variants of Laplacian matrix such as [2, 3]:

$$\boldsymbol{L} \leftarrow \boldsymbol{D}^{-\alpha} \boldsymbol{A} \boldsymbol{D}^{-\alpha}, \tag{3}$$

  where $\alpha \geq 0$ is a parameter. A common value for this parameter is $\alpha = 0.5$:

$$\boldsymbol{L} = \boldsymbol{D}^{-1/2} \boldsymbol{A} \boldsymbol{D}^{-1/2}. \tag{4}$$

  This matrix is also referred to as the normalized Laplacian matrix.

- Here, the normalized Laplacian is used.

# Fourier Functions

- Consider the eigenvalue decomposition of the normalized Laplacian matrix [4]:

$$L = U \Lambda U^\top, \tag{5}$$

  where $U = [u_1, \ldots, u_n] \in \mathbb{R}^{n \times n}$ and $\Lambda = \mathbf{diag}([\lambda_1, \ldots, \lambda_n]^\top) \in \mathbb{R}^{n \times n}$ contain the eigenvectors and eigenvalues of the normalized Laplacian matrix, respectively.

- The eigenvectors of the (normalized) Laplacian, i.e., $u_1, \ldots, u_n$, are called the **Fourier functions**.

- The **Fourier transform** is projecting a signal $x$ on the Fourier functions.

- The result is the coefficients of the **Fourier series**.

$$u^\top x$$

$$U^\top x$$

# Graph Fourier Transform

$$U = [u_1, \ldots, u_n]$$

$$U^\top x$$

- **Graph Fourier transform** projects the input graph signal to a space whose orthonormal bases are the eigenvectors of the normalized Laplacian of the graph.
- For now, assume that every node of graph has a scalar feature value. Let $\mathbb{R}^n \ni x = [x_1, \ldots, x_n]^\top$ be the vector of features of all nodes in the graph, where $x_i \in \mathbb{R}$ is the feature vector of the $i$-th node.
- The graph Fourier transform of $x$ is its projection onto the column space of the matrix $U$:

$$f(x) = \hat{x} = U^\top x. \tag{6}$$

- The inverse graph Fourier transform reconstructs the signal back from projection:

$$f^{-1}(\hat{x}) = U f(x) = U U^\top x. \tag{7}$$

# Graph Convolution

- The **graph convolution** of the input signal $x$ with the filter $g \in \mathbb{R}^n$ is defined as:

$$x * g = f^{-1}(\hat{f(x)}f(g)) \stackrel{(6)}{=} f^{-1}(U^\top x U^\top g) \stackrel{(7)}{=} U(U^\top x U^\top g). \tag{8}$$

- We define:

$$\mathbb{R}^{n \times n} \ni G := \text{diag}(U^\top g) = \begin{bmatrix} u_1^\top g & 0 & \cdots & 0 \\ 0 & u_2^\top g & \cdots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \cdots & u_n^\top g \end{bmatrix}. \tag{9}$$

- Hence, the graph convolution can be stated as:

$$\mathbb{R}^n \ni x * g = UGU^\top x. \tag{10}$$

- If every node has a feature vector rather than a feature value, the features become a matrix $X \in \mathbb{R}^{n \times d}$ where every row is the $d$-dimensional feature vector of a node. Then, the graph convolution becomes:

$$\mathbb{R}^{n \times d} \ni X * g = UGU^\top X \tag{11}$$

**ChebNet**

# ChebNet

- Convolutional graph neural networks have been built upon two main approaches:
    - **spectral** methods which have a graph signal processing perspective.
    - **spatial** methods which define graph convolution by information propagation.
- **Graph Convolutional Network (GCN)** [5] bridged the gap between spectral and spatial approaches.
- Recall Eq. (11). If the input of the $\ell$-th layer is denoted by $\boldsymbol{H}^{(\ell-1)}$ and the output of the $\ell$-th layer be $\boldsymbol{H}^{(\ell)}$, then Eq. (11) becomes:

$$\boldsymbol{H}^{(\ell)} = \sigma(\boldsymbol{U}\boldsymbol{G}\boldsymbol{U}^{\top}\boldsymbol{H}^{(\ell-1)}), \tag{12}$$

where the activation function $\sigma(.)$ has been applied on the result of graph convolution. The first layer accepts the data features as input:

$$\boldsymbol{H}^{(0)} = \boldsymbol{X}. \tag{13}$$

# ChebNet

- A big limitation with Eq. (12) is that $U$ in that equation is the matrix of eigenvectors of the Laplacian of its input graph. The computational complexity of the eigenvalue decomposition of the $n \times n$ Laplacian matrix is $\mathcal{O}(n^3)$.
- **ChebNet** (2016) [6] improves the computational complexity of the convolutional neural network. It approximates the filter $g$ by Chebyshev polynomials of the diagonal matrix of eigenvalues $\Lambda$.
- The **Chebyshev polynomials** are:

$$\begin{cases} T_0(x) = 1, \\ T_1(x) = x, \\ T_i(x) = 2xT_{i-1}(x) - T_{i-2}(x). \end{cases} \tag{14}$$

The domain of input $x$ for Chebyshev polynomials is $[-1, 1]$. for example, the Chebyshev polynomials are widely used for cosine expressions:

$$\cos(i\alpha) = T_i(\cos(\alpha)).$$

$$\cos(2\alpha) = T_2(\cos\alpha) =$$
$$2\cos\alpha\cos\alpha - 1$$
$$= 2\cos^2\alpha - 1$$

$$\cos(0\alpha) = T_0(\cos\alpha) = 1$$
$$\cos(\alpha) = T_1(\cos\alpha) = \cos\alpha$$

# ChebNet

- ChebNet approximates the filter **G** by a linear combination of Chebyshev polynomials of the eigenvalues **Λ**:

$$G = \sum_{i=0}^{k} \theta_i T_i(\mathbf{\Lambda}),$$

where $k$ is the order of Chebyshev polynomials.

- However, there is a problem with the domain of the Chebyshev polynomials in this equation. The eigenvalues, i.e., the diagonal elements of **Λ** are between zero and the largest eigenvalue $\lambda_{max}$. Therefore, the eigenvalues need to be normalized as:

$$\mathbb{R}^{n \times n} \ni \widetilde{\mathbf{\Lambda}} := \frac{2}{\lambda_{max}} \mathbf{\Lambda} - \mathbf{I}_n, \tag{15}$$

where $\mathbf{I}_n$ is the $n \times n$ identity matrix. The values in the normalized eigenvalue matrix are in range $[-1, 1]$ as required by the domain of Chebyshev polynomials.

- Hence, the approximation of the filter $g$ is:

$$G = \sum_{i=0}^{k} \theta_i T_i(\widetilde{\mathbf{\Lambda}}). \tag{16}$$

# ChebNet

- Recall Eq. (10):

$$\boldsymbol{x} * \boldsymbol{g} = \underbrace{\boldsymbol{U} \boldsymbol{G} \boldsymbol{U}^\top \boldsymbol{x}}_{} \overset{(16)}{=} \boldsymbol{U} \left( \sum_{i=0}^{k} \theta_i T_i(\widetilde{\boldsymbol{\Lambda}}) \right) \boldsymbol{U}^\top \boldsymbol{x} = \sum_{i=0}^{k} \theta_i \boldsymbol{U} T_i(\widetilde{\boldsymbol{\Lambda}}) \boldsymbol{U}^\top \boldsymbol{x}. \tag{17}$$

- The matrix $\boldsymbol{U}$ is orthogonal, i.e., its columns are orthonormal, because it is the matrix of eigenvectors. For an orthonormal transformation, the following holds:

$$\boldsymbol{U} T_i(\widetilde{\boldsymbol{\Lambda}}) \boldsymbol{U}^\top = T_i(\boldsymbol{U}\widetilde{\boldsymbol{\Lambda}}\boldsymbol{U}^\top). \tag{18}$$

- Similar to Eq. (15), we define:

$$\widetilde{\boldsymbol{L}} := \frac{2}{\lambda_{\max}} \boldsymbol{L} - \boldsymbol{I}_n, \tag{19}$$

where $\lambda_{\max}$ is largest eigenvalue of the normalized Laplacian $\boldsymbol{L}$. Then, according to Eq. (15) and similar to Eq. (5), the eigenvalue decomposition of $\widetilde{\boldsymbol{L}}$ becomes:

$$\boldsymbol{L} = \boldsymbol{U}\boldsymbol{\Lambda}\boldsymbol{U}^\top \implies \widetilde{\boldsymbol{L}} = \boldsymbol{U}\widetilde{\boldsymbol{\Lambda}}\boldsymbol{U}^\top. \tag{20}$$

- Combining Eqs. (18) and (20) gives:

$$\boldsymbol{U} T_i(\widetilde{\boldsymbol{\Lambda}}) \boldsymbol{U}^\top = T_i(\widetilde{\boldsymbol{L}}). \tag{21}$$

- Putting Eq. (21) in Eq. (17) gives:

$$\boldsymbol{x} * \boldsymbol{g} = \sum_{i=0}^{k} \theta_i T_i(\widetilde{\boldsymbol{L}})\boldsymbol{x}. \tag{22}$$

# ChebNet

- Comparing Eqs. (12) and (22):

$$H^{(\ell)} = \sigma(\boldsymbol{U}\boldsymbol{G}\boldsymbol{U}^\top \boldsymbol{H}^{(\ell-1)}), \quad \rightarrow \text{regular GC}$$

$$\boldsymbol{x} * \boldsymbol{g} = \sum_{i=0}^{k} \theta_i T_i(\tilde{\boldsymbol{L}})\boldsymbol{x}, \quad \rightarrow \text{ChebNet}$$

shows that ChebNet resolves the limitation of eigenvalue decomposition of the Laplacian. In fact, it uses the approximation of Chebyshev polynomials and does not perform eigenvalue decomposition.

**Graph Convolutional Network**

# Graph Convolutional Network

- **Graph Convolutional Network (GCN)** (2017) [5] is the first-order approximation of the ChebNet. In Eq. (22), it approximates the Chebyshev polynomials to its first order ($k = 1$):

$$T_i(\widetilde{L}) \approx T_0(\widetilde{L}) + T_1(\widetilde{L}). \tag{23}$$

In other words:

$$x * g \approx \sum_{i=0}^{1} \theta_i T_i(\widetilde{L})x = \theta_0 T_0(\widetilde{L})x + \theta_1 T_1(\widetilde{L})x \overset{(14)}{=} \theta_0 x + \theta_1 \widetilde{L}x.$$

- More number of learnable parameters may result in overfitting [7]. To reduce the number of parameters and to avoid overfitting, it is assumed that $\theta_0 = \theta_1 = \theta$, so:

$$x * g = \theta x + \theta \widetilde{L}x = \theta(I + \widetilde{L})x \overset{(19)}{=} \theta(I + \frac{2}{\lambda_{max}} L - I)x = \theta\left(\frac{2}{\lambda_{max}}\right)Lx.$$

- It is possible to absorb the constant $2/\lambda_{max}$ into the learnable parameters and simply the graph convolution as:
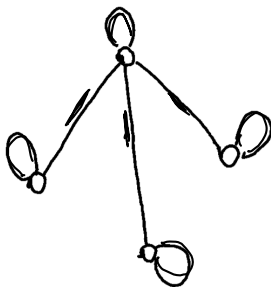
$$x * g = \theta Lx \overset{(4)}{=} \theta D^{-1/2} A D^{-1/2} x. \tag{24}$$

# Graph Convolutional Network

- We found:

$$\boldsymbol{x} * \boldsymbol{g} = \theta \boldsymbol{D}^{-1/2} \boldsymbol{A} \boldsymbol{D}^{-1/2} \boldsymbol{x}.$$

- It has been empirically observed that this results in some instability in training of GCN. Therefore, an additional assumption is added to have self-loops on the nodes meaning that every node has an edge from it to itself.

# Graph Convolutional Network

- Mathematically, it means that the main diagonal of the adjacency matrix should become one by adding the identity matrix to it. Therefore, we define:

$$\widetilde{A} := A + I,$$
$$\widetilde{D}(i,j) := \sum_{j=1}^{n} \widetilde{A}(i,j), \qquad (25)$$
$$\overline{L} := \widetilde{D}^{-1/2} \widetilde{A} \widetilde{D}^{-1/2}.$$

- As a result, the Eq. (24) is replaced by:

$$x * g = \theta \widetilde{D}^{-1/2} \widetilde{A} \widetilde{D}^{-1/2} x = \theta \overline{L} x.$$

In matrix form, if every row of $X \in \mathbb{R}^{n \times d}$ is the $d$-dimensional feature vector of a node, this equation becomes $x * g = \overline{L} X \theta$ where $\theta \in \mathbb{R}^d$. If there is a need to have $f$ feature maps after the convolution, then this equation can become $x * g = \overline{L} X \Theta$ where $\Theta \in \mathbb{R}^{d \times f}$.

- As a result, if the input of the $\ell$-th layer is denoted by $H^{(\ell-1)}$ and the output of the $\ell$-th layer be $H^{(\ell)}$, then Eq. (11) becomes:

$$H^{(\ell)} = \sigma(\overline{L} H^{(\ell-1)} \Theta), \qquad (26)$$

where the activation function $\sigma(.)$ has been applied on the result of graph convolution. The first layer accepts the data features as input, as stated in Eq. (13).

# Graph Convolutional Network

- Eq. (26) is the graph convolution performed in every layer of GCN where $\boldsymbol{\Theta}$ is the matrix of learnable weights in the layer.
- Comparing Eqs. (12) and (26):

$$\boldsymbol{H}^{(\ell)} = \sigma(\boldsymbol{U}\boldsymbol{G}\boldsymbol{U}^{\top}\boldsymbol{H}^{(\ell-1)}),$$
$$\boldsymbol{H}^{(\ell)} = \sigma(\tilde{\boldsymbol{L}}\boldsymbol{H}^{(\ell-1)}\boldsymbol{\Theta}),$$

shows that GCN resolves the limitation of eigenvalue decomposition of the Laplacian. It makes use of the approximation of Chebyshev polynomials and does not perform eigenvalue decomposition.

# Graph Convolutional Network vs. Feedforward Network

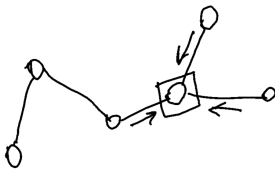- In the fully connected layer of a feedforward neural network, the operation of the layer is:

$$H^{(\ell)} = \sigma(H^{(\ell-1)}\Theta). \qquad (27)$$

However, according to Eqs. (25) and (26), the operation of convolution in a layer of GCN is:

$$H^{(\ell)} = \sigma(\widetilde{D}^{-1/2}\widetilde{A}\widetilde{D}^{-1/2}H^{(\ell-1)}\Theta). \qquad (28)$$

- Comparing Eqs. (27) and (28) shows the relation of GCN and feedforward network. In a fully connected layer of feedforward network, all the features of previous layer $H^{(\ell-1)}$ are fed to the next layer through a linear transformation by the weight matrix $\Theta$ followed by a nonlinear activation function. However, in graph neural network, firstly the adjacency matrix defines which nodes (or features) are connected to each other, and then the linear transformation by the weight matrix $\Theta$ is performed followed by a nonlinear activation function. In other words, the adjacency matrix determines which nodes should impact the features of every node (see this this figure).



when complete graph
(all nodes are connected)
reduces to fully conneed layer

when not connected
to a node
$\widetilde{A} = 0 \longrightarrow H^{(\ell)} = 0$
from not have impact
that node

**More General
Frameworks of Graph
Convolutional Network**

# More General Frameworks of GCN

- The update rule of every layer, i.e., Eq. (28), can be restated as:

$$h_i^{(\ell)} = \sigma\Big(\sum_{j \in \mathcal{N}_i} \Theta h_j^{(\ell-1)}\Big), \tag{29}$$

  for the $i$-th neuron in the $\ell$-th layer, where $\mathcal{N}_i$ denotes the neighbors of the $i$-th node (or neuron) in the input of the layer. This update rule is called **sum pooling** because it sums over the neighbors.

- There is a problem with sum pooling. Summing the contents of the neighboring nodes (or neurons) increases the scale of the output feature gradually over multiple layers.

- To resolve this issue, it is possible to normalize the input of the activation function by $\widetilde{D}^{-1}$:

$$H^{(\ell)} = \sigma\Big(\widetilde{D}^{-1}\widetilde{A}H^{(\ell-1)}\Theta\Big), \tag{30}$$

  where $\widetilde{D}$ is defined in Eq. (25). Eq. (30) can be stated for every node $i$:

$$\widetilde{D}(i,i) = |\mathcal{N}_i|$$
$$\widetilde{D}(i,i)^{-1} = \frac{1}{|\mathcal{N}_i|}$$

$$h_i^{(\ell)} = \sigma\Big(\sum_{j \in \mathcal{N}_i} \frac{1}{|\mathcal{N}_i|}\Theta h_j^{(\ell-1)}\Big), \tag{31}$$

  where $|\mathcal{N}_i|$ denotes the number of neighbors for the $i$-th node. This is because the degree matrix counts the number of neighbors for nodes.

- The update rule in Eq. (30) or (31) is called the **mean pooling**.

# More General Frameworks of GCN

- Rather than Eq. (30), it is possible to use symmetric normalization in mean pooling:
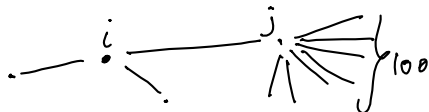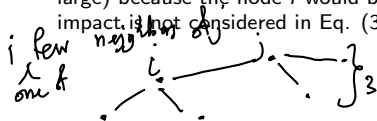
$$H^{(\ell)} = \sigma\left(\tilde{D}^{-1/2}\tilde{A}\tilde{D}^{-1/2}H^{(\ell-1)}\Theta\right). \tag{32}$$

- Eq. (32) can be stated for every node $i$:

$$h_i^{(\ell)} = \sigma\left(\sum_{j\in\mathcal{N}_i} \frac{1}{\sqrt{|\mathcal{N}_i||\mathcal{N}_j|}}\Theta h_j^{(\ell-1)}\right), \tag{33}$$

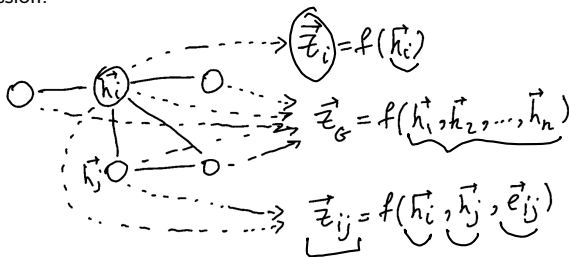which is called **mean pooling** with symmetric normalization.

- Comparing Eqs. (28) and (32) shows that the original GCN uses mean pooling with symmetric normalization.

- Eq. (33) means that for every node $i$, if the node $j$ is a neighbor, its impact on the node $i$ should be more if the node $j$ has few number of neighbors ($|\mathcal{N}_j|$ is small). However, its impact on the node $i$ should be less if the node $j$ has large number of neighbors ($|\mathcal{N}_j|$ is large) because the node $i$ would be one of the many neighbors of node $j$. Note that this impact is not considered in Eq. (31).

# More General Frameworks of GCN

- Different tasks:
  - ▶ <u>Node classification/regression</u>: after the multiple layers of convolution, the $h_i$'s of the last layer are used in the loss function for the node classification or regression.
  - ▶ <u>Graph classification/regression</u>: after the multiple layers of convolution, all the $h_i$'s of the last layer are aggregated and used in the loss function for the graph classification or regression.
  - ▶ <u>Link classification/regression</u>: after the multiple layers of convolution, the $h_i$'s and the edges of the last layer are used in the loss function for the link classification or regression.

$$\vec{z}_i = f(\vec{h_i})$$

$$\vec{z}_G = f(\vec{h_1}, \vec{h_2}, \dots, \vec{h_n})$$

$$\vec{z}_{ij} = f(\vec{h_i}, \vec{h_j}, \vec{e}_{ij})$$

**Graph Attention Network**

# Graph Attention Network

- As was seen in Eqs. (29), (31), and (33), the linear combination in pooling can have weights.
- **Graph Attention Network (GAT)** (2017) [8] adopts attention mechanisms to learn the relative weights between two connected nodes. In the pooling operation, the weights of attention are added:

$$h_i^{(\ell)} = \sigma\Big(\sum_{j \in \mathcal{N}_i} \alpha_{ij} h_j^{(\ell-1)}\Big), \tag{34}$$

where the attention weight $\alpha_{ij}$ measures the influence of node $j$ to node $i$.

$$\alpha_{ij} = \text{attention}(h_i^{(\ell-1)}, h_j^{(\ell-1)}). \tag{35}$$

- The attention weight can be computed by a attention function $a(.)$ between $h_i^{(\ell-1)}$ and $h_j^{(\ell-1)}$:

$$a_{ij} = a(h_i^{(\ell-1)}, h_j^{(\ell-1)}). \tag{36}$$

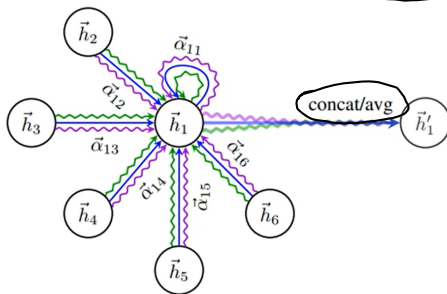This attention function may also consider the edge between the nodes $i$ and $j$:

$$a_{ij} = a(h_i^{(\ell-1)}, h_j^{(\ell-1)}, e_{ij}). \tag{37}$$

# Graph Attention Network

- This attention function $a(.)$ can be a transformer autoencoder [9].
- However, GAT models the attention function $a(.)$ as a single-layer feedforward neural network. This single-layer neural network calculates the attention between nodes.
- Finally, the attention values of every node are normalized in a softmax form to obtain the attention weights:

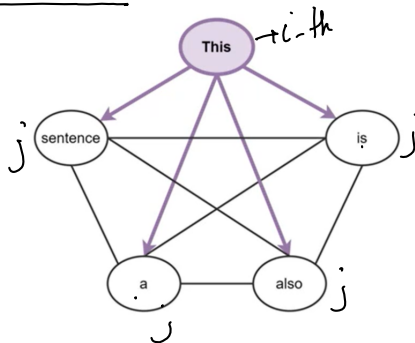$$\alpha_{ij} = \frac{e^{a_{ij}}}{\sum_{k \in \mathcal{N}_i} e^{a_{ik}}}, \tag{38}$$

where the summation in the denominator is over the neighbors of the $i$-th node.



credit of image: [9]

# Graph Attention Network

- Transformers [9] are special cases of graph neural networks.
- In fact, every sentence or sequence can be considered as a graph where GAT can calculate the attention between the tokens in the sequence. For example, the graph for the sentence "This is also a sentence" is depicted below.

# Graph Attention Network

- In the following, GAT and transformer are compared.
- In GAT, the attention is $a_{ij} = a(\mathbf{h}_i^{(\ell-1)}, \mathbf{h}_j^{(\ell-1)})$ where the $\mathbf{h}_i^{(\ell-1)}$ and $\mathbf{h}_j^{(\ell-1)}$ are passed through a single-layer network with some weight $\mathbf{W}$. Therefore, the attention is calculated between $\mathbf{W}^\top \mathbf{h}_i^{(\ell-1)}$ and $\mathbf{W}^\top \mathbf{h}_j^{(\ell-1)}$ after feeding to the single layer of network. In transformer, on the other hand, the attention is $a_{ij} = a(\mathbf{q}_i, \mathbf{k}_j)$ where the query $\mathbf{q}_i$ and key $\mathbf{k}_j$ are different linear transformations of the same tokens, i.e., $\mathbf{q}_i = \mathbf{W}_Q^\top \mathbf{x}$ and $\mathbf{k}_i = \mathbf{W}_K^\top \mathbf{x}$ [10]. Therefore, the difference of GAT and transformer is that GAT uses the shared learnable wight for the query and key but transformer uses different learnable weights for them.
- Another difference between GAT and transformer is that GAT uses a single-layer feedforward neural network as the attention function $a(.)$. However, in transformer, this function is [10]:

$$a(\mathbf{q}_i, \mathbf{k}_j) = \frac{1}{\sqrt{p}} \mathbf{q}_i^\top \mathbf{k}_j,$$

where $p$ is the dimensionality of query and key.
- The last difference of GAT and transformer is the softmax form. GAT sums over the neighbors in the denominator of the softmax form (see Eq. (38)). However, transformer sums over all tokens in the sequence:

$$\alpha_{ij} = \frac{e^{a_{ij}}}{\sum_{k=1}^n e^{a_{ik}}}.$$

**Graph Autoencoder**

# Graph Autoencoder

- Consider the following autoencoder where the encoder has two layers. This autoencoder accepts a graph as its input; hence, its name is **Graph Autoencoder (GAE)** (2016) [11].
- According to Eq. (26), the first layer of the encoder is:

$$H^{(1)} = \sigma(\bar{L}X\Theta_1), \qquad (39)$$

where $\Theta_1$ is the learnable weight matrix of the first layer, $X \in \mathbb{R}^{n \times d}$ is the feature vectors of nodes stacked row-wise, $\bar{L}$ is defined in Eq. (25) based on the adjacency matrix of the graph, $\sigma(.)$ is usually the ReLU activation function [12], and $H^{(1)}$ is the output of the first layer.

- Again, according to Eq. (26), the second layer of the encoder is:

$$H^{(2)} = \bar{L}H^{(1)}\Theta_2, \qquad (40)$$

where $\Theta_2$ is the learnable weight matrix of the second layer, $H^{(2)}$ is the output of the second layer, and the second layer is assumed not to have an actvation function.

- Putting Eq. (39) in Eq. (40) gives the following function which we denote by $GCN(X, A; \Theta_1, \Theta_2)$:

$$GCN(X, A; \Theta_1, \Theta_2) := \bar{L}\,\sigma(\bar{L}X\Theta_1)\Theta_2. \qquad (41)$$

- There are two types of GAE, i.e., graph reconstruction autoencoder and graph variational autoencoder [11]. These autoencoders are introduced in the following.

# Graph Reconstruction Autoencoder

- In the graph reconstruction autoencoder, also called the non-probabilistic GAE, the encoder is Eq. (41) with two layers. The $p$-dimensional latent embeddings of nodes, denoted by $Z \in \mathbb{R}^{n \times p}$, are obtained as:

$$Z = \text{GCN}(X, A; \Theta_1, \Theta_2) := \bar{L}\,\sigma(\bar{L}X\Theta_1)\Theta_2.$$

- The decoder of graph reconstruction autoencoder does not contain any layers but models measuring similarity between the embedding vectors of the nodes (see this figure). It is the sigmoid function of $z_i^\top z_j$ to show the score of similarity (inner product) of the latent variables $z_i$ and $z_j$. In other words, it reconstructs the adjacency matrix but with the latent embeddings of nodes rather than the nodes directly:

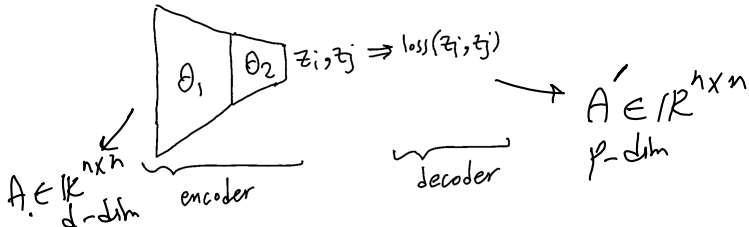$$\widehat{A} = \text{sigmoid}(ZZ^\top), \text{ or} \tag{42}$$

$$\widehat{A}(i,j) = \frac{1}{1 + e^{-z_i^\top z_j}}. \tag{43}$$

- The graph reconstruction autoencoder is depicted in this figure.

# Graph Reconstruction Autoencoder

- The loss is the mean squared error between the adjacency matrix and the reconstructed adjacency matrix:

$$\underset{\theta}{\text{minimize}} \quad \|\widehat{A} - A\|_F^2, \tag{44}$$

where $\|.\|_F$ denotes the Frobenius norm and $\theta := \{\Theta_1, \Theta_2\}$ is the learnable parameters. This loss function is minimized by backpropagation [13].

no decoder
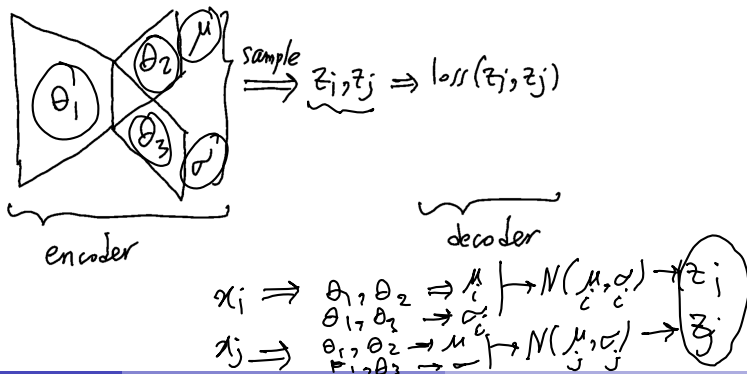
decoder : calculating $\widehat{A}$

# Graph Variational Autoencoder

- Graph variational autoencoder uses these two GCN modules for estimating the mean and variance in the latent space by the encoder:

$$\text{GCN}_\mu(\boldsymbol{X}, \boldsymbol{A}; \boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2) := \bar{\boldsymbol{L}}\,\sigma(\bar{\boldsymbol{L}}\boldsymbol{X}\boldsymbol{\Theta}_1)\boldsymbol{\Theta}_2, \tag{45}$$

$$\text{GCN}_\sigma(\boldsymbol{X}, \boldsymbol{A}; \boldsymbol{\Theta}_1, \boldsymbol{\Theta}_3) := \bar{\boldsymbol{L}}\,\sigma(\bar{\boldsymbol{L}}\boldsymbol{X}\boldsymbol{\Theta}_1)\boldsymbol{\Theta}_3, \tag{46}$$

where the first layer is shared between them as shown in this figure.

# Graph Variational Autoencoder

- As the latent variables of the nodes are independent, the encoder of graph variational autoencoder models the following conditional distribution:

$$q(\mathbf{Z}|\mathbf{X}, \mathbf{A}) = \prod_{i=1}^{n} q(\mathbf{z}_i|\mathbf{X}, \mathbf{A}), \tag{47}$$

where $\mathbf{Z} \in \mathbb{R}^{n \times p}$ contains the $p$-dimensional latent variables and $\mathbf{z}_i \in \mathbb{R}^p$ is the latent variable of the $i$-th node whose conditional distribution is a Gaussian distribution:

$$q(\mathbf{z}_i|\mathbf{X}, \mathbf{A}) = \mathcal{N}(\mathbf{z}_i \mid \boldsymbol{\mu}_i, \mathbf{diag}(\boldsymbol{\sigma}_i^2)), \tag{48}$$

where $\mathbf{diag}(.)$ makes a diagonal matrix with its input as as the diagonal of matrix.

- The latent variables $\{\mathbf{z}_i\}_{i=1}^n$ are sampled from the multivariate joint distribution in Eq. (47).

$$z_i \sim q(z_i \mid X, A)$$
$$z_j \sim g(z_j \mid X, A)$$

# Graph Variational Autoencoder

- The decoder of the autoencoder models the following conditional distribution:

$$q(\boldsymbol{A}|\boldsymbol{Z}) = \prod_{i=1}^{n} \prod_{j=1}^{n} p(\boldsymbol{A}(i,j)|\boldsymbol{z}_i, \boldsymbol{z}_j), \tag{49}$$

where $p(\boldsymbol{A}(i,j)|\boldsymbol{z}_i, \boldsymbol{z}_j)$ is the sigmoid function of $\boldsymbol{z}_i^\top \boldsymbol{z}_j$ to show the probability of similarity (inner product) of the latent variables $\boldsymbol{z}_i$ and $\boldsymbol{z}_j$:

$$p(\boldsymbol{A}(i,j) = 1|\boldsymbol{z}_i, \boldsymbol{z}_j) = \frac{1}{1 + e^{-\boldsymbol{z}_i^\top \boldsymbol{z}_j}}. \tag{50}$$

As a result, the decoder of graph variational autoencoder does not contain any layers but models measuring similarity between the sampled latent variables in the latent space (see this figure).

# Graph Variational Autoencoder

- The graph variational autoencoder maximizes the Evidence Lower Bound (ELBO) in variational inference [14]:

$$\underset{\theta}{\text{maximize}} \; \underbrace{\mathbb{E}_{q(\mathbf{Z}|\mathbf{X},\mathbf{A})}\big[\log(p(\mathbf{A}\,|\,\mathbf{Z}))\big]} - \underbrace{\text{KL}\big(q(\mathbf{Z}|\mathbf{X},\mathbf{A})\|\,p(\mathbf{Z})\big)}. \qquad (51)$$

where $\text{KL}(.\|.)$ denotes the Kullback-Leibler (KL) divergence [15], $p(\mathbf{Z})$ is the desired prior distribution such as some Gaussian distribution, and $\theta := \{\Theta_1, \Theta_2, \Theta_3\}$ is the learnable parameters.

- The graph variational autoencoder is trained by backpropagation [13]. In backpropagation, the loss function should be minimized; therefore, the loss is the ELBO times $-1$:

$$\underset{\theta}{\text{minimize}} \quad -\mathbb{E}_{q(\mathbf{Z}|\mathbf{X},\mathbf{A})}\big[\log(p(\mathbf{A}\,|\,\mathbf{Z}))\big] + \text{KL}\big(q(\mathbf{Z}|\mathbf{X},\mathbf{A})\|\,p(\mathbf{Z})\big). \qquad (52)$$

- minimizing this loss function tries to learn generation of the adjacency matrix $\mathbf{A}$ given the sampled latent variables $\mathbf{Z}$ while the conditional distribution of the latent variable given the graph and its adjacency matrix becomes similar to the desired prior distribution of the latent space.

# Acknowledgment

- Some slides of this slide deck are inspired by teachings of Prof. Ali Ghodsi at University of Waterloo, Department of Statistics.
- Graph neural network in PyTorch Geometric: `https://pytorch-geometric.readthedocs.io/en/latest/get_started/introduction.html`
- Good tutorial on PyTorch Geometric by Antonio Longa: `https://www.youtube.com/playlist?list=PLGMXrbDNfqTzqxB1IGgimuhtfAhGd8lHF`

# References

[1] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[2] Y. Weiss, "Segmentation using eigenvectors: a unifying view," in *Proceedings of the seventh IEEE international conference on computer vision*, vol. 2, pp. 975–982, IEEE, 1999.

[3] A. Ng, M. Jordan, and Y. Weiss, "On spectral clustering: Analysis and an algorithm," *Advances in neural information processing systems*, vol. 14, pp. 849–856, 2001.

[4] B. Ghojogh, F. Karray, and M. Crowley, "Eigenvalue and generalized eigenvalue problems: Tutorial," *arXiv preprint arXiv:1903.11240*, 2019.

[5] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations*, 2017.

[6] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," *Advances in neural information processing systems*, vol. 29, pp. 3844–3852, 2016.

[7] B. Ghojogh and M. Crowley, "The theory behind overfitting, cross validation, regularization, bagging, and boosting: tutorial," *arXiv preprint arXiv:1905.12787*, 2019.

[8] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," in *International Conference on Learning Representations*, 2017.

# References (cont.)

[9] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, vol. 30, 2017.

[10] B. Ghojogh and A. Ghodsi, "Attention mechanism, transformers, BERT, and GPT: tutorial and survey," 2020.

[11] T. N. Kipf and M. Welling, "Variational graph auto-encoders," *arXiv preprint arXiv:1611.07308*, 2016.

[12] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 807–814, 2010.

[13] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, no. 6088, pp. 533–536, 1986.

[14] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "Factor analysis, probabilistic principal component analysis, variational inference, and variational autoencoder: Tutorial and survey," *arXiv preprint arXiv:2101.00734*, 2021.

[15] S. Kullback and R. A. Leibler, "On information and sufficiency," *The annals of mathematical statistics*, vol. 22, no. 1, pp. 79–86, 1951.