

Spectral Clustering, Graph Neural Networks, and Transformers

CS 145: Introduction to Data Mining (Spring 2024)

Yanqiao Zhu

UCLA

May 17, 2024

Announcements

- HW4 scores have been released
- HW5 solution will be released tomorrow
- No more assignments!

Midterm Exam

Date • May 20, 12:00PM–1:45PM

- In-person only; no online or make-up exams offered

Format • Close-book but double-sided letter-sized cheat sheets allowed, **up to 10 pages**

- Simple calculators allowed
- Internet access strictly prohibited

Scope • Supervised Learning: Linear Regression, Logistic Regression, MLP, Gradient Descent, Backpropagation, Regularization, Batch/Layer Norm, Decision Trees, Random Forests, Mixture-of-Expert, Ensemble (Bagging, Boosting, Adaboost), K-Nearest Neighbor

- Unsupervised Learning: K-Means, Gaussian Mixture Models, EM Algorithm, (Variational) Auto Encoders
- Graphs and Networks: Random Walks, Spectral Clustering, Graph Representation Learning

Midterm Exam

Structure

- Part A. True/False Questions ($8 * 2 = 16$ points)
- Part B. Multiple-Choice Questions ($5 * 2 = 10$ points)
- Part C. Fill-in-the-Blank Questions ($21 * 1 = 21$ points)
- Part D. Open-Answer Questions (63 points)
 - Problem 17. Linear Regression with Regularization (16 points)
 - Problem 18. Multilayer Perceptron and Backpropagation (15 points)
 - Problem 19. Decision Trees and Bagging (16 points)
 - Problem 20. K-Means and Gaussian Mixture Models (16 points)

Scores

- 110 points in total, with the additional **10** points serving as bonus points

Outline

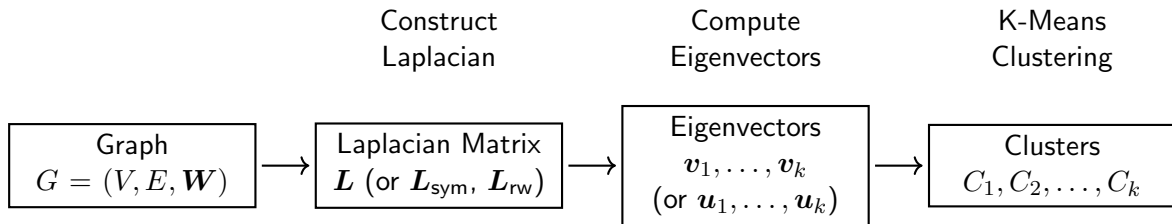
- 1 Spectral Clustering
- 2 Graph Neural Networks
- 3 Transformers

Graph Representations and Notations

- A graph is denoted as $G = (V, E)$, where V is the set of vertices (nodes) and E is the set of edges
- $|V| = n$ is the number of vertices, and $|E| = m$ is the number of edges
- Adjacency matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$: $\mathbf{A}_{ij} = 1$ if $(i, j) \in E$, and 0 otherwise
- Degree matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$: diagonal matrix with $\mathbf{D}_{ii} = \sum_{j=1}^n \mathbf{A}_{ij}$

Spectral Clustering

- Spectral clustering: K-Means in a transformed space
 - Preprocessing: Construct a graph and compute its Laplacian matrix
 - Transformation: Compute eigenvectors of the Laplacian matrix
 - Clustering: Run K-Means on the eigenvectors
- Motivation: Solve graph partitioning problems by optimizing graph cut objectives

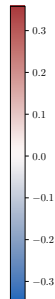
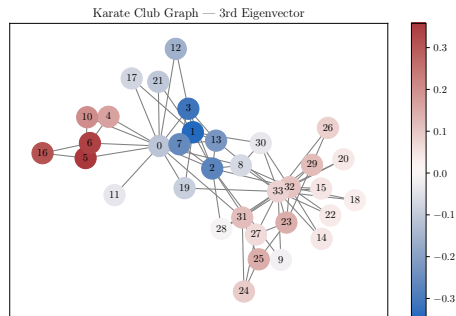
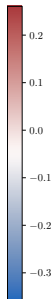
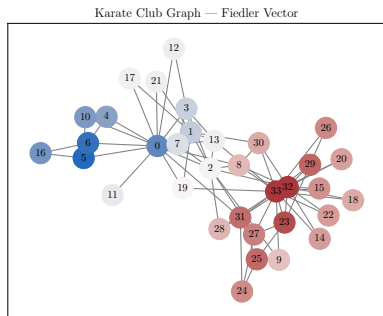


Graph Laplacian and Spectrum

- Undirected weighted graph $G = (V, E, \mathbf{W})$
- Degree matrix \mathbf{D} : diagonal matrix with $D_{ii} = \sum_{j=1}^n W_{ij}$
- Unnormalized Laplacian: $\mathbf{L} = \mathbf{D} - \mathbf{W}$
- Normalized Laplacians:
 - Symmetric: $\mathbf{L}_{\text{sym}} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2}$
 - Random walk: $\mathbf{L}_{\text{rw}} = \mathbf{D}^{-1} \mathbf{L}$
- Quadratic form: $\mathbf{x}^\top \mathbf{L} \mathbf{x} = \sum_{(i,j) \in E} W_{ij} (x_i - x_j)^2$
- Spectrum: eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$
- Key properties:
 - $\mathbf{L}, \mathbf{L}_{\text{sym}}, \mathbf{L}_{\text{rw}}$: symmetric, positive semi-definite
 - $\lambda_1 = 0, \mathbf{v}_1 = \mathbf{1}$ (constant vector)
 - Number of connected components = multiplicity of $\lambda_1 = 0$
 - Fiedler vector (eigenvector \mathbf{v}_2) for bi-partitioning

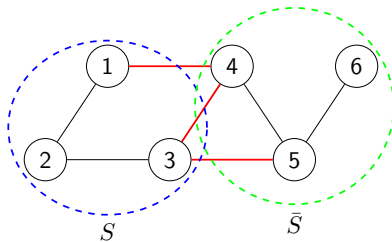
Graph Laplacian and Spectrum

- Eigenvalues and eigenvectors of these matrices encode important structural properties of the graph
- The Fiedler vector shows a gradual change across the graph, indicating a smooth transition of values



Graph Cuts and Optimization

- Graph cut: $\text{cut}(S, \bar{S}) = \sum_{i \in S, j \in \bar{S}} W_{ij}$
- Ratio cut objective: $\sum_{i=1}^k \frac{\text{cut}(S_i, \bar{S}_i)}{|S_i|}$
 - Relaxation: $\min_{\mathbf{X}} \text{tr}(\mathbf{X}^T \mathbf{L} \mathbf{X})$ s.t. $\mathbf{X}^T \mathbf{X} = \mathbf{I}$
 - Solution: $\mathbf{X} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k]$
- Normalized cut objective: $\sum_{i=1}^k \frac{\text{cut}(S_i, \bar{S}_i)}{\text{vol}(S_i)}$, where $\text{vol}(S_i) = \sum_{j \in S_i} D_{jj}$
 - Relaxation: $\min_{\mathbf{X}} \text{tr}(\mathbf{X}^T \mathbf{L}_{\text{rw}} \mathbf{X})$ s.t. $\mathbf{X}^T \mathbf{D} \mathbf{X} = \mathbf{I}$
 - Solution: $\mathbf{X} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k]$, where $\mathbf{L}_{\text{rw}} \mathbf{u}_i = \lambda_i \mathbf{D} \mathbf{u}_i$

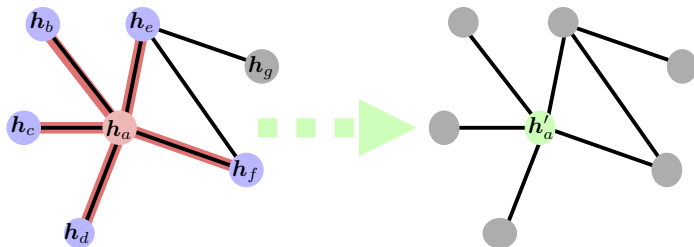


Spectral Clustering Algorithms

- Unnormalized spectral clustering (Ratio cut):
 - ① Construct unnormalized Laplacian \mathbf{L}
 - ② Compute first k eigenvectors of \mathbf{L} : $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$
 - ③ Run K-Means on the rows of $\mathbf{X} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k]$
- Normalized spectral clustering (Normalized cut):
 - ① Construct normalized Laplacian \mathbf{L}_{rw} (or \mathbf{L}_{sym})
 - ② Compute first k generalized eigenvectors: $\mathbf{L}_{rw}\mathbf{u}_i = \lambda_i \mathbf{D}\mathbf{u}_i$
 - ③ Run K-Means on the rows of $\mathbf{X} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k]$

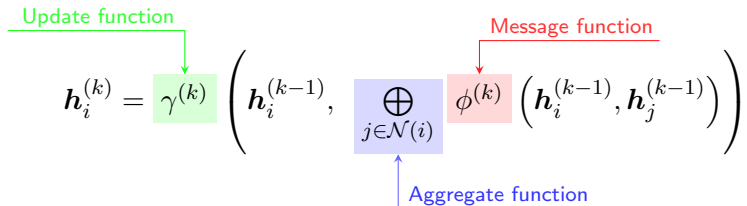
Graph Neural Networks (GNNs)

- GNNs generalize convolutional operations to the graph domain
- Key idea: neighborhood aggregation
 - Update node representations by aggregating information from neighboring nodes
 - Repeated for multiple layers to capture higher-order dependencies



Graph Neural Networks (GNNs)

- Formally, a GNN layer can be decomposed into three functions:



The diagram illustrates the decomposition of a GNN layer into three functions: Update, Message, and Aggregate. The equation is
$$\mathbf{h}_i^{(k)} = \gamma^{(k)} \left(\mathbf{h}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)} \left(\mathbf{h}_i^{(k-1)}, \mathbf{h}_j^{(k-1)} \right) \right)$$
. The Update function $\gamma^{(k)}$ is highlighted in a green box with a green arrow pointing to it from the label 'Update function'. The Message function $\phi^{(k)}$ is highlighted in a red box with a red arrow pointing to it from the label 'Message function'. The Aggregate function, represented by the summation symbol $\bigoplus_{j \in \mathcal{N}(i)}$, is highlighted in a blue box with a blue arrow pointing to it from the label 'Aggregate function'.

$$\mathbf{h}_i^{(k)} = \gamma^{(k)} \left(\mathbf{h}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)} \left(\mathbf{h}_i^{(k-1)}, \mathbf{h}_j^{(k-1)} \right) \right)$$

Three “Flavors” of GNNs

- Graph Convolutional Network (GCN):

- Aggregate neighborhood information using convolutional filters
- Simplifies spectral graph convolutions using first-order approximation
- $\mathbf{H}^{(l+1)} = \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{(l)} \mathbf{W}^{(l)})$

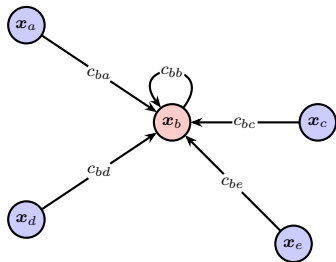
- Graph Attention Network (GAT):

- Introduces attention mechanism to assign different weights to neighbors
- $\mathbf{h}_i^{(l+1)} = \sigma(\sum_{j \in \mathcal{N}(i)} \alpha_{ij}^{(l)} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)})$
- $\alpha_{ij}^{(l)} = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^{(l)\top} [\mathbf{W}^{(l)} \mathbf{h}_i^{(l)} \parallel \mathbf{W}^{(l)} \mathbf{h}_j^{(l)}]))}{\sum_{k \in \mathcal{N}(i)} \exp(\text{LeakyReLU}(\mathbf{a}^{(l)\top} [\mathbf{W}^{(l)} \mathbf{h}_i^{(l)} \parallel \mathbf{W}^{(l)} \mathbf{h}_k^{(l)}]))}$

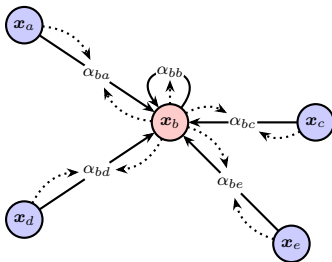
- Message Passing Neural Network (MPNN):

- Unifies various GNN architectures under a general message passing framework
- $\mathbf{m}_i^{(l+1)} = \sum_{j \in \mathcal{N}(i)} M^{(l)}(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)}, \mathbf{e}_{ij})$
- $\mathbf{h}_i^{(l+1)} = U^{(l)}(\mathbf{h}_i^{(l)}, \mathbf{m}_i^{(l+1)})$

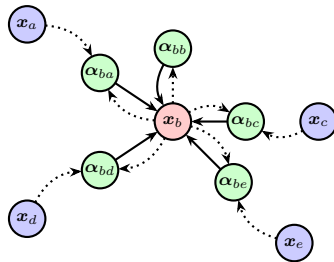
Three “Flavors” of GNNs



Convolutional



Attentional



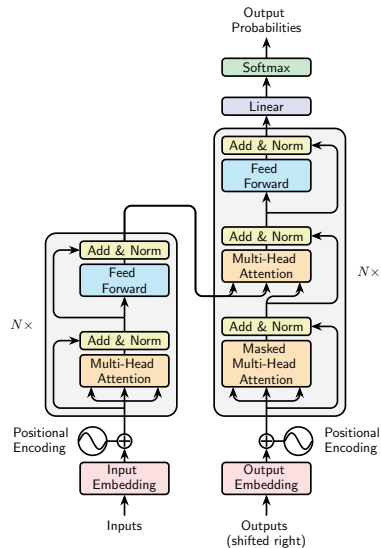
Message-Passing

Graph Tasks with GNN Embeddings

- Node-level tasks:
 - Node classification: $\hat{y}_i = f_c(\mathbf{h}_i)$, where f_c is a classifier
 - Node regression: $\hat{y}_i = f_r(\mathbf{h}_i)$, where f_r is a regression model
- Link-level tasks:
 - Link prediction: $\hat{y}_{ij} = f_l(\mathbf{h}_i \parallel \mathbf{h}_j)$, where f_l is a binary classifier and \parallel denotes concatenation
 - Edge classification: $\hat{y}_{ij} = f_e(\mathbf{h}_i \parallel \mathbf{h}_j)$, where f_e is a classifier
- Community-level tasks:
 - Community detection: $\{C_1, \dots, C_K\} = \text{Clustering}(\{\mathbf{h}_i \mid i \in V\})$
 - Subgraph classification: $\hat{y}_{S_k} = f_s(\mathbf{h}_{S_k})$, where $\mathbf{h}_{S_k} = \text{Aggregate}(\{\mathbf{h}_i \mid i \in S_k\})$ and f_s is a classifier
- Graph-level tasks:
 - Graph classification: $\hat{y}_{G_i} = f_g(\mathbf{h}_{G_i})$, where $\mathbf{h}_{G_i} = \text{Aggregate}(\{\mathbf{h}_i \mid i \in V_i\})$ and f_g is a classifier
 - Graph regression: $\hat{y}_{G_i} = f_r(\mathbf{h}_{G_i})$, where $\mathbf{h}_{G_i} = \text{Aggregate}(\{\mathbf{h}_i \mid i \in V_i\})$ and f_r is a regression model

Transformers

- Transformers are a type of deep learning model that rely heavily on the attention mechanism
- Introduced by Vaswani et al. (2017) for sequence-to-sequence tasks, particularly in natural language processing



Tokenization

- Tokenization is the process of breaking down a sequence (e.g., a sentence) into smaller units called tokens
- Tokens can be:
 - Words
 - Subwords (e.g., WordPiece, BPE)
 - Characters
- Each token is mapped to a unique integer ID, which is then used to look up the corresponding embedding vector

The quick brown fox jumps over the lazy dog



[The] [quick] [brown] [fox] [jumps] [over] [the] [lazy] [dog]



[12] [542] [1201] [783] [3120] [4500] [12] [7891] [1345]

Attention Mechanism

- Attention is a mechanism that allows the model to focus on different parts of the input sequence when generating each output token
- Mathematically, attention can be described as a weighted sum of values, where the weights are computed based on the query and keys:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right) \mathbf{V}$$

- \mathbf{Q} (Query): The current hidden state
- \mathbf{K} (Keys): The hidden states of the input sequence
- \mathbf{V} (Values): The hidden states of the input sequence (same as keys)
- The attention weights are computed by taking the dot product between the query and keys, scaled by $\sqrt{d_k}$ (the dimensionality of the keys), and then passed through a softmax function

Positional Encoding

- Positional encoding is a way to inject position information into the input embeddings
- It allows the model to learn relative positions of tokens in a sequence
- Positional encodings are typically sinusoidal functions of different frequencies

$$\text{PE}_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d_{\text{model}}}}\right)$$
$$\text{PE}_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{\text{model}}}}\right)$$

- pos : The position in the sequence
- i : The dimension index
- d_{model} : The dimensionality of the embeddings

Next Token Prediction

- Transformers are often used in a language modeling setup, where the goal is to predict the next token in a sequence given the previous tokens
- To perform next token prediction, the model:
 - ① Passes the input tokens through the Transformer layers
 - ② Takes the final hidden state corresponding to the last input token
 - ③ Passes this hidden state through a linear layer followed by a softmax to produce a probability distribution over the vocabulary
 - ④ Chooses the token with the highest probability as the predicted next token
- During training, the model is typically trained using a cross-entropy loss between the predicted probabilities and the actual next token
- During inference, the predicted next token can be sampled from the output distribution or chosen deterministically (e.g., by taking the argmax)

Transformer Architectures

- Encoder-only models:
 - Stack of encoder layers
 - Self-attention and feed-forward networks
 - Used for tasks like classification and extraction
- Decoder-only models:
 - Stack of decoder layers
 - Masked self-attention and feed-forward networks
 - Used for tasks like language modeling and generation
- Encoder-decoder models:
 - Encoder stack processes input sequence
 - Decoder stack attends to encoder output and generates output
 - Used for tasks like translation and summarization

Spectral Clustering

- Clustering based on graph Laplacian eigenvectors
- Partitions data using graph cut objectives
- Relies on the spectrum of the graph Laplacian matrix
- Algorithms: unnormalized, normalized spectral clustering

Graph Neural Networks

- Neural networks operating on graph-structured data
- Aggregate neighborhood information to update node representations
- Variants: convolutional, attentional, message-passing
- Applications: node classification, link prediction, graph classification

Transformers

- Sequence-to-sequence models based on attention mechanisms
- Architectures: encoder-only, decoder-only, encoder-decoder
- Key components: tokenization, attention mechanisms, positional encoding
- Applications: language modeling, machine translation, text generation