

# CS224W Homework 3 - Kyuhyeok Seo

April 4, 2025

## 1 GNNs as MLP of eigenvectors

### 1.1 Batch Node Update

Consider the update for Graph Isomorphism Network:

$$\mathbf{x}_v^{(l+1)} = \text{MLP} \left( (1 + \epsilon) \mathbf{x}_v^{(l)} + \sum_{u \in \mathcal{N}(v)} \mathbf{x}_u^{(l)} \right), \quad (1)$$

where  $\mathbf{x}_v^{(l)} \in \mathbb{R}^{d_l}$  is the embedding of node  $v$  at layer  $l$ . Let  $\mathbf{X}^{(l)} \in \mathbb{R}^{N \times d_l}$  be a matrix containing the embeddings of all the nodes in the graph, i.e.,  $\mathbf{X}^{(l)}[:, v] = \mathbf{x}_v^{(l)}$ . Also, let  $\mathbf{A} \in \{0, 1\}^{N \times N}$  represent the adjacency matrix of the graph. Write down the update of  $\mathbf{X}^{(l+1)}$  as a function of  $\mathbf{X}^{(l)}$  and  $\mathbf{A}$ .

$$\star \text{ Solution } \star \mathbf{X}^{(l+1)} = \text{MLP}((1 + \epsilon)\mathbf{X}^{(l)} + \mathbf{A}\mathbf{X}^{(l)})$$

### 1.2 Single Layer MLP

Assume that  $\text{MLP}()$  represents a single layer MLP with no bias term. Write down the update of  $\mathbf{X}^{(l+1)}$  as a function of  $\mathbf{X}^{(l)}$  and  $\mathbf{A}$ , and the trainable parameters  $\mathbf{W}^{(l)}$  of layer  $l$ .

$$\star \text{ Solution } \star \mathbf{X}^{(l+1)} = \sigma(\mathbf{W}^{(l)}((1 + \epsilon)\mathbf{X}^{(l)} + \mathbf{A}\mathbf{X}^{(l)}))$$

### 1.3 Eigenvector Extension

Let  $\{\lambda_n, \mathbf{v}_n\}_{n=1}^N$  represent the eigenvalues and eigenvectors of the graph adjacency. Then we can write  $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$ , where  $\mathbf{V} \in \mathbb{R}^{N \times N}$  is the matrix of eigenvectors with  $\mathbf{V}[:, n] = \mathbf{v}_n$  and  $\mathbf{\Lambda} \in \mathbb{R}^{N \times N}$  is the diagonal matrix of eigenvalues with  $\mathbf{\Lambda}[n, n] = \lambda_n$ . Show that

$$\mathbf{X}^{(l+1)} = \sigma \left( \mathbf{V}\hat{\mathbf{W}}^{(l)} \right), \quad \hat{\mathbf{W}}^{(l)}[n, j] = (\lambda_n + 1 + \epsilon) \sum_{i=1}^{d_l} \mathbf{W}^{(l)}[i, j] \langle \mathbf{v}_n, \mathbf{X}^{(l)}[:, i] \rangle,$$

where  $\langle \cdot \rangle$  denotes the dot product. Hint: Use the fact that the eigenvectors are orthonormal. Next, show that each feature across all nodes,  $\mathbf{X}^{(l+1)}[:, i]$ , can be expressed as a linear combination of eigenvectors, followed by a pointwise nonlinearity.

★ **Solution** ★  $(1 + \epsilon)\mathbf{X}^{(l)} + \mathbf{A}\mathbf{X}^{(l)} = \mathbf{V}(I + \epsilon I + \Lambda)\mathbf{V}^T\mathbf{X}^{(l)}$ , then Define  $\Lambda' = I + \epsilon I + \Lambda$ ,  $\lambda'_n = 1 + \epsilon + \lambda_n$ . Therefore,  $\mathbf{X}^{(l+1)} = \sigma(\mathbf{W}^{(l)}(\mathbf{V}\Lambda'\mathbf{V}^T\mathbf{X}^{(l)}))$ . Let's focus on  $i$ -th column of  $\mathbf{V}\mathbf{W}^{(l)}$ . It is  $\sum_{j=1}^N v_j(1 + \epsilon + \lambda_j) \sum_{k=1}^{d_l} W^{(l)}[k, i] v_j^T x_k^{(l)} = \sum_{k=1}^{d_l} W^{(l)}[k, i] \sum_{j=1}^N v_j(1 + \epsilon + \lambda_j) v_j^T x_k^{(l)} = \mathbf{W}^{(l)}(\mathbf{V}\Lambda'\mathbf{V}^T\mathbf{X}^{(l)})$ 's  $i$ -th column.

## 1.4 GraphSAGE

Perform the same analysis for the GraphSAGE update when the aggregation function is sum pooling. Recall that the GraphSAGE update function is

$$\begin{aligned} \mathbf{x}_v^{(l+1)} &= \sigma \left( \mathbf{W}^{(l)} \cdot \text{CONCAT} \left( \mathbf{x}_v^{(l)}, \mathbf{x}_{N(v)}^{(l)} \right) \right) \\ &= \sigma \left( \mathbf{W}_1^{(l)} \mathbf{x}_v^{(l)} + \mathbf{W}_2^{(l)} \text{AGG} \left( \mathbf{x}_u^{(l)}, \forall u \in N(v) \right) \right) \end{aligned}$$

★ **Solution** ★ The sum aggregation function in GraphSAGE updates node embeddings by summing over neighboring node features. Mathematically, this update can be rewritten as:  $X^{(l+1)} = \sigma \left( W_1^{(l)} X^{(l)} + W_2^{(l)} A X^{(l)} \right)$ .

Using the eigendecomposition of the adjacency matrix:  $A = V\Lambda V^T$ ,  $X^{(l+1)} = \sigma \left( W_1^{(l)} X^{(l)} + W_2^{(l)} V\Lambda V^T X^{(l)} \right)$ .

Multiplying both sides by  $V^T$ :  $V^T X^{(l+1)} = \sigma \left( V^T W_1^{(l)} X^{(l)} + V^T W_2^{(l)} V\Lambda V^T X^{(l)} \right)$ .

Define the transformed representation:  $\tilde{X}^{(l)} = V^T X^{(l)}$ , which gives:  $V^T X^{(l+1)} = \sigma \left( V^T W_1^{(l)} V \tilde{X}^{(l)} + V^T W_2^{(l)} V \Lambda \tilde{X}^{(l)} \right)$ .

Defining the transformed weight matrix:  $\hat{W}^{(l)}[n, j] = \left( W_1^{(l)} + \lambda_n W_2^{(l)} \right) \sum_{i=1}^{d_l} W^{(l)}[i, j] \langle v_n, X^{(l)}[:, i] \rangle$ ,

we obtain the final spectral form:  $X^{(l+1)} = \sigma \left( V \hat{W}^{(l)} \right)$ .

This result shows that GraphSAGE with sum pooling behaves like a spectral transformation where eigenvalues  $\lambda_n$  modulate the contribution of neighbor information via  $W_2^{(l)}$ , while self-information is preserved via  $W_1^{(l)}$ .

## 2 LightGCN

We learned in class about **LightGCN**, a GNN model for recommender systems. Given a bipartite user-item graph  $G = (V, E)$ , let  $\mathbf{A} \in \mathbb{R}^{|V| \times |V|}$  be its unnormalized adjacency matrix,  $\mathbf{D} \in \mathbb{R}^{|V| \times |V|}$  be its degree matrix and  $\mathbf{E}^{(k)} \in \mathbb{R}^{|V| \times d}$  be its node embedding matrix at layer  $k$  where  $d$  is the embedding dimension.

Let  $\tilde{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$  be the normalized adjacency matrix.

The original GCN updates node embeddings across layers according to  $\mathbf{E}^{(k+1)} = \text{ReLU}(\tilde{\mathbf{A}} \mathbf{E}^{(k)} \mathbf{W}^{(k)})$ , while LightGCN removes the non-linearity and uses the equation for each layer  $k \in \{0, 1, \dots, K-1\}$ :

$$\mathbf{E}^{(k+1)} = \tilde{\mathbf{A}} \mathbf{E}^{(k)} \quad (2)$$

Moreover, LightGCN adopts multi-scale diffusion to compute the final node embeddings for link prediction, averaging across layers:

$$\mathbf{E} = \sum_{i=0}^K \alpha_i \mathbf{E}^{(i)}, \quad (3)$$

where we have uniform coefficients  $\alpha_i = \frac{1}{K+1}$ .

## 2.1 Advantages of Average Embeddings

Why does LightGCN average over layer embeddings? What benefits does it bring, in a recommendation systems setting?

★ **Solution** ★ Averaging over layer embeddings in LightGCN enables the model to capture multi-hop neighborhood information without over-smoothing, which often occurs in deep GCNs. This helps balance the influence of both immediate and distant neighbors, improving representation quality for link prediction in recommendation systems. Additionally, it enhances robustness and generalization by aggregating signals from different levels of the graph hierarchy.

## 2.2 Self-connection

We denote the embedding of an item  $i$  at layer- $k$   $\mathbf{e}_i^{(k)}$  and that of a user  $u$   $\mathbf{e}_u^{(k)}$ . The graph convolution operation (a.k.a., propagation rule) in LightGCN is defined as:

$$\begin{aligned} \mathbf{e}_u^{(k+1)} &= \sum_{i \in \mathcal{N}_u} \frac{1}{\sqrt{|\mathcal{N}_u|} \sqrt{|\mathcal{N}_i|}} \mathbf{e}_i^{(k)} \\ \mathbf{e}_i^{(k+1)} &= \sum_{u \in \mathcal{N}_i} \frac{1}{\sqrt{|\mathcal{N}_i|} \sqrt{|\mathcal{N}_u|}} \mathbf{e}_u^{(k)} \end{aligned}$$

The symmetric normalization term  $\frac{1}{\sqrt{|\mathcal{N}_u|} \sqrt{|\mathcal{N}_i|}}$  follows the design of standard GCN, which can avoid the scale of embeddings increasing with graph convolution operations.

However, from the equations above, we can find that in LGCN, we only aggregate the connected neighbors and do not integrate the target node itself (i.e., there is no **self-connection**). This is different from most existing graph

convolution operations that typically aggregate extended neighbors and also specifically handle self-connection.

Does LightGCN contain implicit self-connection? If your answer is yes, which operation captures the same effect as self-connection? If no, what do you think is the reason why LightGCN doesn't need self-connection or similar effects?

★ **Solution** ★ Yes, LightGCN contains an implicit self-connection through the layer-wise embedding averaging in the final embedding computation. By including the initial embedding  $\mathbf{e}^{(0)}$  (which is unaffected by neighbors) in the final average  $\mathbf{E} = \frac{1}{K+1} \sum_{k=0}^K \mathbf{E}^{(k)}$ , the model retains each node's original identity, achieving a similar effect to self-connections without explicitly adding them during message passing.

## 2.3 Relation with APPNP

There is a work that connects GCN with Personalized PageRank, where the authors propose a GCN variant named APPNP that can propagate long range without the risk of oversmoothing. Inspired by the teleport design in Personalized PageRank, APPNP complements each propagation layer with the starting features (i.e., the 0-th layer embeddings), which can balance the need of preserving locality (i.e., staying close to the root node to alleviate oversmoothing) and leveraging the information from a large neighborhood. The propagation layer in APPNP is defined as:

$$\mathbf{E}^{(k+1)} = \beta \mathbf{E}^{(0)} + (1 - \beta) \tilde{\mathbf{A}} \mathbf{E}^{(k)}$$

where  $\beta$  is called the “teleport probability” to control the retention of starting features in the propagation, and  $\tilde{\mathbf{A}}$  denotes the normalized adjacency matrix.

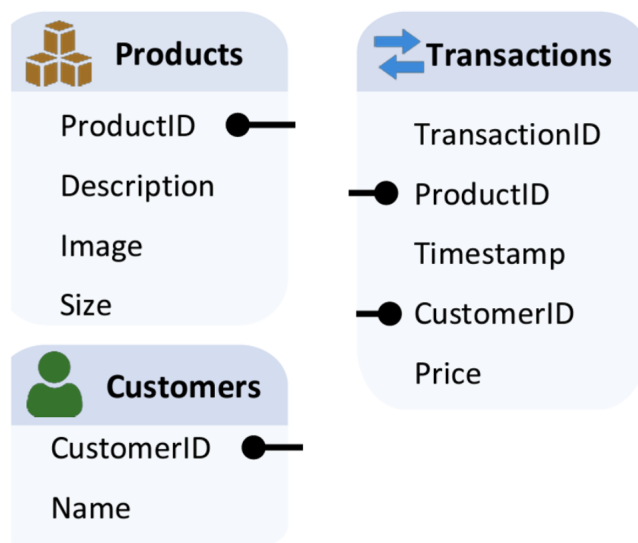
Aligning with Equation (3), we can see that by setting  $\alpha_k$  accordingly, LightGCN can fully recover the prediction embedding used by APPNP. As such, LightGCN shares the strength of APPNP in combating oversmoothing — by setting the  $\alpha$  properly, LightGCN allows using a large  $K$  for long-range modeling with controllable oversmoothing.

Express the layer- $K$  embeddings  $\mathbf{E}^{(K)}$  of APPNP as a function of the initial embeddings  $\mathbf{E}^{(0)}$  and the normalized adjacency matrix  $\tilde{\mathbf{A}}$ . Show all work.

**What to submit?** Multi-line mathematical derivation of the relationship between  $\mathbf{E}^{(K)}$  and  $\mathbf{E}^{(0)}$

$$\star \text{ **Solution** } \star \mathbf{E}^{(K)} = \beta \mathbf{E}^{(0)} \sum_{k=0}^K (1 - \beta)^k \tilde{\mathbf{A}}^k + (1 - \beta)^{K+1} \tilde{\mathbf{A}}^K \mathbf{E}^{(0)}$$

### 3 Relational Deep Learning



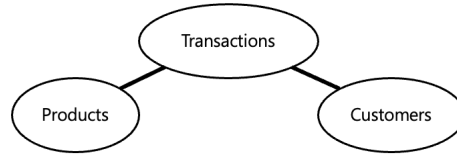
Assume we have the relational database as seen above, which consists of three tables. These tables contain information about products, customers, and transactions in which customers purchase products. Each table contains a unique identifier, known as a *primary key*, potentially along with other attributes. *Foreign keys* in a table create connections between tables by referencing primary keys in other tables. In the three tables shown above, “ProductID”, “TransactionID”, and “CustomerID” are the primary keys, while “ProductID” and “CustomerID” are also foreign keys for the “Transactions” table.

#### 3.1 Schema Graphs

A key component of a relational deep learning framework is the schema graph, which illustrates the relationships between tables in a database. In a schema graph, each table is represented as a node, and an edge is drawn between two nodes if a primary key from one table appears as a foreign key in another. This graph helps visualize how data is linked across the database.

Describe or draw what the schema graph of this database would look like (hint: it’s very simple).

★ **Solution** ★



### 3.2 Relational Entity Graphs

Table 1: Products





ProductID	Description	Image	Size
1	Smartphone		Small
2	Laptop		Medium
3	TV		Large
4	Headphones		Small

Table 2: Customers

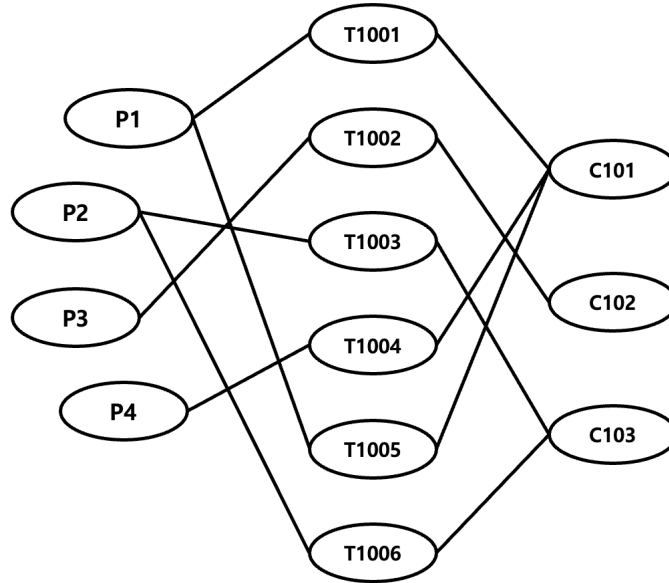
CustomerID	Name
101	Alice
102	Bob
103	Carol

Table 3: Transactions

TransactionID	ProductID	Timestamp	CustomerID	Price (\$)
1001	1	2024-10-15	101	600
1002	3	2024-10-20	102	500
1003	2	2024-10-26	103	1,300
1004	4	2024-11-01	101	100
1005	1	2024-11-02	101	600
1006	2	2024-11-12	103	1,300

Another component of this framework is the relational entity graph. The nodes of this graph are all the individual entities rather than tables. Links are again made by primary-foreign key connections - that is, two entities are linked if they appear together in the same entry of any table in the database. Given the list of transactions above, produce a relational entity graph describing this database.

★ **Solution** ★



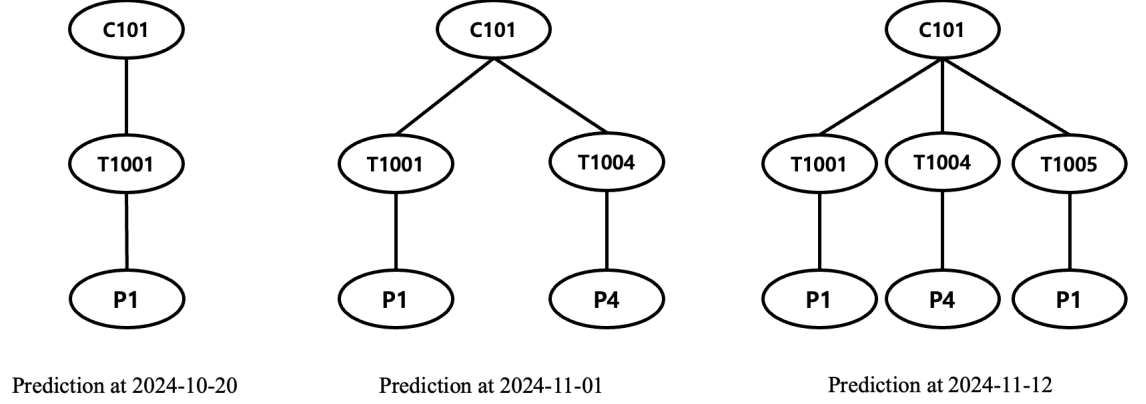
### 3.3 Computation Graphs [6 points]

The computational graphs used for training are dependent on the specific timestep used for prediction. For example, let's assume our training table (which defines the information we seek to predict) contains the following information:

- Target: How much total money a customer spends in the next 30 days
- ID: Customer ID
- Timestep: The time at which the 30 day period starts

When predicting, we can only use the information in the database that takes place before our prediction period. That means the computational graphs (the specific set of nodes and connections we send messages over) we use for predictions are directly dependent on the timestep in our training table. Let's say we want to make predictions for customer 101. Using the tables from the previous part, draw out the computation graphs if we wanted to make predictions at 2024-10-20, 2024-11-01, 2024-11-12.

★ **Solution** ★



### 3.4 Message Passing

A relational database will produce a heterogeneous graph. What are example message passing and update rules that can be used to make predictions like the one mentioned above?

★ **Solution** ★

#### 1. Message Passing

- Item to Customer: Customers aggregate information from the items they purchased (They are connected through an edge).
- Customer to Item: Items aggregate information from the customers who purchased them (They are connected through an edge).

**2. Update Rules** Each node updates its embedding using information from its neighbors. For a customer  $u$ , the update rule is:

$$h_u^{(k+1)} = \sigma \left( W_{u1}^{(k+1)} h_u^{(k)} + W_{u2}^{(k+1)} \text{AGG}(h_i^{(k)}, i \in N(u)) \right)$$

where  $W_{u1}^{(k+1)}$  and  $W_{u2}^{(k+1)}$  are learnable weight matrices.

For a product  $i$ , the update rule is:

$$h_i^{(k+1)} = \sigma \left( W_{i1}^{(k+1)} h_i^{(k)} + W_{i2}^{(k+1)} \text{AGG}(h_u^{(k)}, u \in N(i)) \right)$$

where  $W_{i1}^{(k+1)}$  and  $W_{i2}^{(k+1)}$  are learnable weight matrices.



**3. Prediction** After  $K$  message passing steps, the final customer embedding is used for prediction:

$$\hat{y}_u = f(h_u^{(K)})$$

where  $f$  is a learned function such as a multi-layer perceptron (MLP).

**Summary:** The model constructs a heterogeneous graph from customer-product interactions, updates embeddings through message passing, and predicts total spending using the final customer embeddings.