

Homework 3

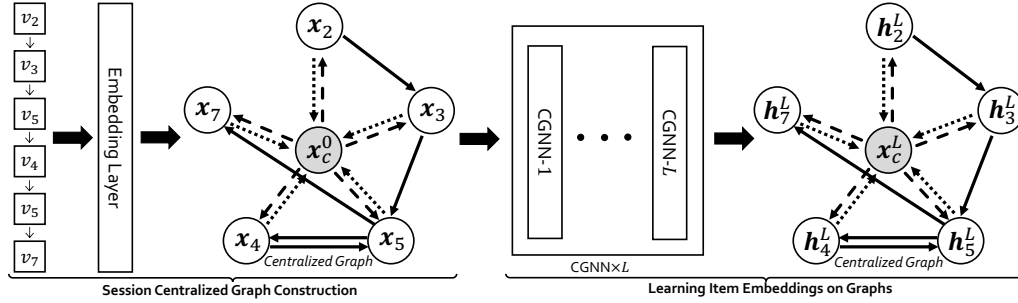
Release Date: Nov 19, 2024

Due Date: Nov 30, 2024

Max mark: 50

25 points

- I. Session-based Recommendation is a task that predicts users' next actions based on their sequential behaviors in the current session. Formally, let V denote all unique items in all sessions. Given a session as $S = \{v_1, v_2, \dots, v_n\}$ consisting of n sequential items, where $v_i \in V$ is the i -th item in the session, the goal of Session-based Recommendation is to predict v_{n+1} to click. In this question, we introduce the *Centralized Graph Neural Network (CGNN)* model to solve this problem.



Session Centralized Graph Construction For each session $S = \{v_1, v_2, \dots, v_n\}$, we construct a centralized graph to represent the transition relationship among the items in the session. Specifically, each session is represented as $G_s = \{\mathcal{V}_s, \mathcal{E}_s\}$, where $\mathcal{V}_s = \{\{x_1, x_2, \dots, x_m\}, x_c\}$ denotes the node set of the centralized graph. Here, the first part $\{x_1, x_2, \dots, x_m\}$ indicates all unique items in the session, which we call *satellite nodes*, and the latter part x_c is the newly added *central node*.

\mathcal{E}_s is the edge set in the centralized graph, which consists of two types of edges, i.e., central connections and satellite connections. For the central connections, we add a bidirectional edge between the central node and each satellite node. For the satellite connections, the edge $(x_i, x_j) \in \mathcal{E}_s$, i.e., the solid lines in the centralized graphs mean that the user clicks item x_j after clicking x_i . In this way, the adjacent relationship between two items in the session can be represented by an incoming matrix \mathbf{A}^I and an outgoing matrix \mathbf{A}^O :

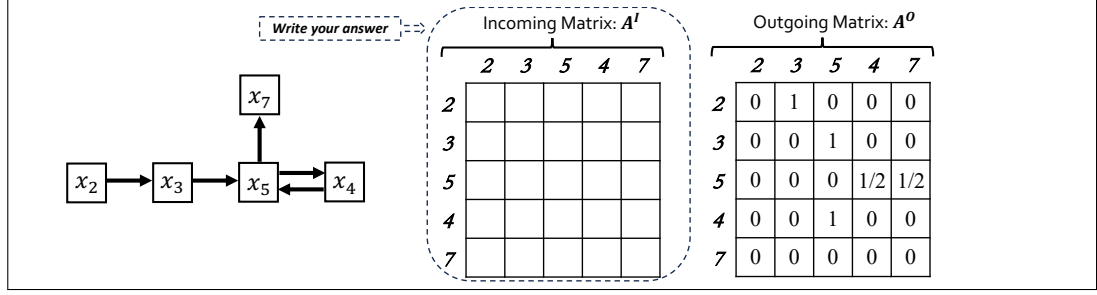
$$\mathbf{A}_{(i,j)}^I = \begin{cases} 1/in_i, & (x_j, x_i) \in \mathcal{E}_s \\ 0, & (x_j, x_i) \notin \mathcal{E}_s \end{cases}, \mathbf{A}_{(i,j)}^O = \begin{cases} 1/out_i, & (x_i, x_j) \in \mathcal{E}_s \\ 0, & (x_i, x_j) \notin \mathcal{E}_s \end{cases} \quad (1)$$

where out_i represents the number of outbound edges for x_i and in_i represents the number of inbound edges for x_i . Please note that **only satellite connections** need to be considered in the adjacency matrices here.

Given an example session $S = \{v_2, v_3, v_5, v_4, v_5, v_7\}$, we can construct the incoming matrix and the outgoing matrix in the following figure.

3 points

- (1) Construct the matrix \mathbf{A}^I for the example above. Please write the answer directly into the diagram above.



Learning Item Embeddings on Graphs We initialize the representation of satellite nodes and the central node. For the satellite nodes, we directly take the trainable embeddings of unique items in the session as the satellite nodes' representation:

$$\mathbf{h}^0 = [\mathbf{x}_1^0, \mathbf{x}_2^0, \dots, \mathbf{x}_m^0], \quad (2)$$

where $\mathbf{x}_i^0 \in \mathbb{R}^d$. As for the central node \mathbf{x}_c^0 , we apply an average pooling on the satellite nodes to get the initialization of the central node, i.e.,

$$\mathbf{x}_c^0 = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i^0 \quad (3)$$

To learn a node representation according to the graph, we update the satellite nodes and the central node as follows. First, we consider the information from adjacent nodes. For each satellite node x_i in the centralized graph at layer l , we utilize the incoming and outgoing matrices to obtain the propagation information as follows:

$$\begin{aligned} \mathbf{a}_i^l &= \text{CONCAT} \left(\mathbf{A}_i^I \left(\mathbf{h}^{l-1\top} \mathbf{W}^I + \mathbf{b}^I \right), \mathbf{A}_i^O \left(\mathbf{h}^{l-1\top} \mathbf{W}^O + \mathbf{b}^O \right) \right)^\top \\ &= \text{CONCAT} \left(\mathbf{A}_i^I \left([\mathbf{x}_1^{l-1}, \mathbf{x}_2^{l-1}, \dots, \mathbf{x}_m^{l-1}]^\top \mathbf{W}^I + \mathbf{b}^I \right), \mathbf{A}_i^O \left([\mathbf{x}_1^{l-1}, \mathbf{x}_2^{l-1}, \dots, \mathbf{x}_m^{l-1}]^\top \mathbf{W}^O + \mathbf{b}^O \right) \right)^\top \end{aligned} \quad (4)$$

where $\mathbf{A}_i^I, \mathbf{A}_i^O \in \mathbb{R}^{1 \times m}$ are the incoming the outgoing weights for node x_i , i.e., the i -th row in the incoming matrix and the outgoing matrix. $\mathbf{W}^I, \mathbf{W}^O$ are the learnable parameters for the incoming edges and the outgoing edges, respectively, while $\mathbf{b}^I, \mathbf{b}^O \in \mathbb{R}^d$ are the bias vectors. Hence we obtain $\mathbf{a}_i^l \in \mathbb{R}^{2d}$ to represent the propagation information for node x_i .

2 points

- (2) How many parameters are there in matrices \mathbf{W}^I and \mathbf{W}^O respectively? (Note: The answer is a polynomial with respect to d .)

Then we feed \mathbf{a}_i^l and node x_i 's previous state \mathbf{h}_i^{l-1} , i.e., the i -th column in \mathbf{h}^{l-1} , into the

gated graph neural networks:

$$\begin{aligned}
z_i^l &= \sigma \left(\mathbf{W}_z \mathbf{a}_i^l + \mathbf{U}_z \mathbf{h}_i^{l-1} \right), \\
r_i^l &= \sigma \left(\mathbf{W}_r \mathbf{a}_i^l + \mathbf{U}_r \mathbf{h}_i^{l-1} \right), \\
\tilde{\mathbf{h}}_i^l &= \tanh \left(\mathbf{W}_h \mathbf{a}_i^l + \mathbf{U}_h \left(r_i^l \odot \mathbf{h}_i^{l-1} \right) \right), \\
\hat{\mathbf{h}}_i^l &= \left(1 - z_i^l \right) \odot \mathbf{h}_i^{l-1} + z_i^l \odot \tilde{\mathbf{h}}_i^l,
\end{aligned} \tag{5}$$

where $\mathbf{W}_z, \mathbf{W}_r, \mathbf{W}_h$ and $\mathbf{U}_z, \mathbf{U}_r, \mathbf{U}_h$ are trainable parameters in the network. σ represents the *sigmoid* function and \odot is the element-wise multiplication. Note that $z_i^l, r_i^l, \tilde{\mathbf{h}}_i^l, \hat{\mathbf{h}}_i^l \in \mathbb{R}^d$.

6 points

- (3) How many parameters are there in $\mathbf{W}_z, \mathbf{W}_r, \mathbf{W}_h, \mathbf{U}_z, \mathbf{U}_r$ and \mathbf{U}_h , respectively? (Note: The answer is a polynomial with respect to d .)

Next, we consider information from the central node. In CGNN, the central node can represent the overall information from all items in a session. For each satellite node, we apply a gating network to decide how much information should be propagated from the central node and the adjacent nodes, respectively. Specifically, we calculate the similarity $\alpha_i^l \in \mathbb{R}$ of each satellite node x_i and the central node x_c with a self-attention mechanism as:

$$\alpha_i^l = \frac{\left(\mathbf{W}_{q1} \hat{\mathbf{h}}_i^l \right)^\top \mathbf{W}_{k1} \mathbf{x}_c^{l-1}}{\sqrt{d}} \tag{6}$$

where $\mathbf{W}_{q1}, \mathbf{W}_{k1}$ are trainable parameters. Then we apply a gating network to selectively integrate the information from the adjacent node $\hat{\mathbf{h}}_i^l$ and from the central node \mathbf{x}_c^{l-1} as follows:

$$\mathbf{h}_i^l = \left(1 - \alpha_i^l \right) \hat{\mathbf{h}}_i^l + \alpha_i^l \mathbf{x}_c^{l-1} \tag{7}$$

For updating the central node, we introduce a self-attention mechanism to assign different degrees of importance to the satellite nodes $\beta \in \mathbb{R}^m$ as follows:

$$\beta = \text{softmax} \left(\frac{\left(\mathbf{W}_{k2} \mathbf{h}^l \right)^\top \mathbf{W}_{q2} \mathbf{x}_c^{l-1}}{\sqrt{d}} \right) \tag{8}$$

where $\mathbf{W}_{q2}, \mathbf{W}_{k2}$ are the trainable parameters. Then, we combine the satellite nodes using a linear combination as a new representation of the central node:

$$\mathbf{x}_c^l = \mathbf{h}^l \beta \tag{9}$$

By updating the satellite nodes and the central node iteratively, we can obtain the item rep-

representations. Moreover, to accurately represent the transition relationship between items, multi-layer CGNNs can be stacked, where the l -layer of CGNN can be denoted as:

$$\mathbf{h}^l, \mathbf{x}_c^l = \text{CGNN}(\mathbf{h}^{l-1}, \mathbf{x}_c^{l-1}, \mathbf{A}^I, \mathbf{A}^O) \quad (10)$$

Finally, we obtain the new representations of the satellite nodes \mathbf{h}^L and the central node \mathbf{x}_c^L , as we employ an L -layer CGNN.

(4) Please briefly answer the following two questions.

2 points

- (i) If you want to demonstrate the effectiveness of CGNN, what experiments would you conduct to illustrate it?

4 points

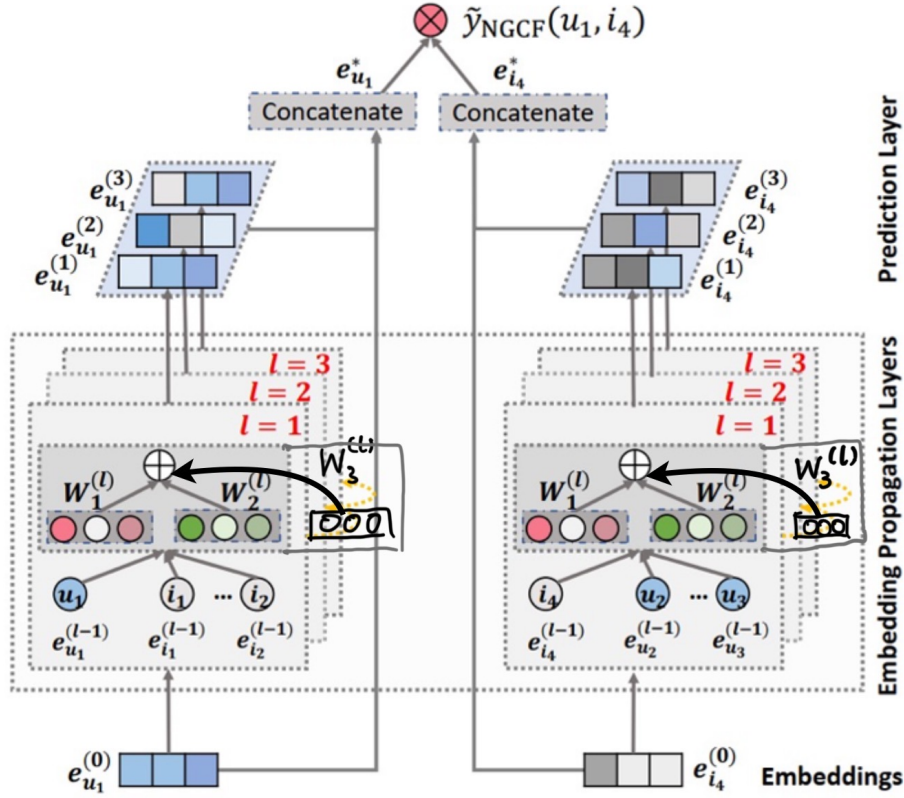
- (ii) When L is large, multiple layers of graph layers propagate a substantial amount of information between nodes, which can easily lead to overfitting. What can be done to improve the model? Please describe two potential approaches.

8 points

- (5) We define $\mathbf{W}_{k1} \in \mathbb{R}^{d \times d}$, $\mathbf{W}_{k2} \in \mathbb{R}^{2d \times d}$. Given $n = 27$, $d = 100$, $L = 1$, and V consists of 2,520 unique items, please compute the number of all trainable parameters of this network.

25 points

II. Consider a variant of the neural graph collaborative filtering (NGCF) model.



Here are some structural explanations of the NGCF-variant:

1. Embedding layer: we describe the user and item with a d -dimension vector. Both embeddings are trainable.
2. Embedding propagation layer: we build upon the message-passing architecture of GNNs in order to capture CF signal along the graph structure and refine the embeddings of users and items.

First-order propagation. We build upon preference basis to perform embedding propagation between the connected users and items, formulating the process with two major operations: message construction and message aggregation.

Message construction. For a connected user-item pair (u, i) , we define the message from i to u as:

$$m_{u \leftarrow i} = \frac{1}{\sqrt{|\mathcal{N}_u| |\mathcal{N}_i|}} (W_1 e_i + W_2 (e_i \odot e_u) + W_3 (e_i || e_u)),$$

where W_1, W_2, W_3 are the trainable weight matrices to distill useful information for propagation, \odot denotes the element-wise product, $||$ denotes the concatenation operation. \mathcal{N}_u and \mathcal{N}_i denote the first-hop neighbors of user u and item i .

Message Aggregation. In this stage, we aggregate the messages propagated from u 's neighborhood to refine u 's representation. Specifically, we define the aggregation function as:

$$e_u^{(1)} = \text{ReLU} \left(W_1 e_u + \sum_{i \in \mathcal{N}_u} m_{u \leftarrow i} \right)$$

where $\mathbf{e}_u^{(1)}$ denotes the representation of user u obtained after the first embedding propagation layer. \mathbf{W}_1 is the weight matrix shared with the one used in message construction. Analogously, we can obtain the representation $\mathbf{e}_i^{(1)}$ for the item i by propagating information from its connected users using a shared weight matrix.

High-order propagation. With the representations augmented by first-order connectivity modeling, we can stack more embedding propagation layers to explore the high-order connectivity information. In the l -th step, the representation of user u is recursively formulated as:

$$\mathbf{e}_u^l = \text{ReLU} \left(\mathbf{m}_{u \leftarrow u}^{(l)} + \sum_{i \in N_u} \mathbf{m}_{u \leftarrow i}^{(l)} \right),$$

wherein the messages being propagated are defined as:

$$\mathbf{m}_{u \leftarrow i}^{(l)} = \frac{1}{\sqrt{|\mathcal{N}_u| |\mathcal{N}_i|}} \left(\mathbf{W}_1^{(l)} \mathbf{e}_i^{(l-1)} + \mathbf{W}_2^{(l)} \left(\mathbf{e}_i^{(l-1)} \odot \mathbf{e}_u^{(l-1)} \right) + \mathbf{W}_3^{(l)} \left(\mathbf{e}_i^{(l-1)} \parallel \mathbf{e}_u^{(l-1)} \right) \right),$$

$$\mathbf{m}_{u \leftarrow u}^{(l)} = \mathbf{W}_1^{(l)} \mathbf{e}_u^{(l-1)},$$

where subscript l denotes the variable in layer l .

Model prediction: After propagating with L layers, we obtain multiple representations for the user and item. We concatenate the representations learned by different layers to get the final embeddings: $\mathbf{e}_u^* = \mathbf{e}_u^{(0)} \parallel \dots \parallel \mathbf{e}_u^{(L)}$, $\mathbf{e}_i^* = \mathbf{e}_i^{(0)} \parallel \dots \parallel \mathbf{e}_i^{(L)}$ where \parallel is the concatenation operation. Finally, we conduct the inner product to estimate the user's preference towards the target item: $\hat{y}_{NGCF(u,i)} = \mathbf{e}_u^{*\top} \mathbf{e}_i^*$.

Here are some assumptions:

1. There are 1,000 unique users and 520 unique items in the training set.
2. The dimensions of user and item embeddings are both 100.
3. The number of embedding propagation layer is 3.
4. The dimension of intermediate output of each embedding propagation layer is [80, 40, 20] respectively.

Compute the number of all trainable parameters of this model.