# LightGNN: Simple Graph Neural Network for Recommendation

Guoxuan Chen
University of Hong Kong
Hong Kong, China
guoxchen@foxmail.com

Lianghao Xia
University of Hong Kong
Hong Kong, China
aka_xia@foxmail.com

Chao Huang*
University of Hong Kong
Hong Kong, China
chaohuang75@gmail.com

## Abstract

Graph neural networks (GNNs) have demonstrated superior performance in collaborative recommendation through their ability to conduct high-order representation smoothing, effectively capturing structural information within users' interaction patterns. However, existing GNN paradigms face significant challenges in scalability and robustness when handling large-scale, noisy, and real-world datasets. To address these challenges, we present LightGNN, a lightweight and distillation-based GNN pruning framework designed to substantially reduce model complexity while preserving essential collaboration modeling capabilities. Our LightGNN framework introduces a computationally efficient pruning module that adaptively identifies and removes redundant edges and embedding entries for model compression. The framework is guided by a resource-friendly hierarchical knowledge distillation objective, whose intermediate layer augments the observed graph to maintain performance, particularly in high-rate compression scenarios. Extensive experiments on public datasets demonstrate LightGNN's effectiveness, significantly improving both computational efficiency and recommendation accuracy. Notably, LightGNN achieves an 80% reduction in edge count and 90% reduction in embedding entries while maintaining performance comparable to more complex state-of-the-art baselines. The implementation of our LightGNN model is available at the github repository: https://github.com/HKUDS/LightGNN.

## CCS Concepts

• **Information systems → Recommender systems**.

## Keywords

Graph Learning, Recommendation, Knowledge Distillation

---

*Chao Huang is the Corresponding Author.

---

## 1 Introduction

Recommender systems [7, 38] have become indispensable in modern online platforms, effectively addressing information overload and enhancing user engagement through personalized service delivery. At the core of these systems, Collaborative Filtering (CF) [14, 20] stands as a dominant paradigm, leveraging users' historical interactions to model latent preferences for behavior prediction.

The evolution of collaborative filtering has spawned diverse approaches, from classical matrix factorization methods (*e.g.* [13]) to sophisticated neural architectures (*e.g.* [9]). Among these developments, Graph Neural Networks (GNNs) have emerged as particularly powerful tools for CF-based recommendation, distinguished by their ability to capture complex, high-order interaction patterns through iterative embedding smoothing. Pioneering works include NGCF [25], which introduced graph convolutional networks (GCNs) to model user-item relationships, and LightGCN [8], which simplifies GCNs to their essential components for recommendation. To address the challenge of sparse interactions in GNN-based recommendation, researchers have developed innovative self-supervised learning (SSL) techniques, including SGL [27], NCL [15], and HCCF [30]. These approaches significantly enhance recommendation accuracy by leveraging self-augmented supervision signals.

Despite significant advancements in GNNs, we would like to emphasize two inherent limitations that continue to challenge GNN-based CF models. **i) Limited scalability of GNNs**: Online recommendation services typically handle vast amounts of relational data (e.g., millions of interactions). This causes the size of user-item graphs to increase dramatically, resulting in a considerable number of information propagation operations within GNNs. Such scalability issues present challenges concerning storage, computational time, and memory requirements. Furthermore, GNN-based CF relies heavily on id-corresponding embeddings for user and item representation [8], with the complexity of these embeddings directly linked to the growing number of users and items, incurring significant memory costs. **ii) Presence of pervasive noise in interaction graphs**: Collaborative recommenders mainly utilize users' implicit feedback, such as clicks and purchases, because of its abundance. However, these interaction records often contain substantial noise that diverges from users' true preferences, including misclicks and popularity biases [23]. Although some existing methods address scalability through techniques like random dropping (e.g., PinSage [33]) or knowledge distillation (KD) (e.g., SimRec [29]), they remain susceptible to misinformation, which can result in inaccurate predictions from their compressed recommenders.

To address these limitations, this paper proposes pruning redundant and noisy components in GNNs, specifically targeting graph edges and embedding entries. We aim to enhance model scalability while preserving essential user preference features. However,

achieving this objective presents non-trivial challenges, outlined as:

- How to identify the graph edges and embedding entries that are genuinely redundant or noisy in the user-item interaction graph?
- How to maintain the high performance of GNN-based CF when significant structural and node-specific information is removed?

As illustrated in Figure 1(a), a considerable proportion of items that users interact with fall into the same category, leading to redundant information about users' preferences. By identifying and removing this redundancy from both structures and parameters, we can significantly reduce the complexity of GNN-based CF. Additionally, many observed interactions represent noise linked to users' negative feedback, as revealed by the review text. This noise can disrupt the preference modeling of existing compressed CF methods, which often fail to explicitly identify such noisy information. Regarding the second challenge, depicted in Figure 1(b), traditional knowledge distillation approaches struggle to effectively maintain performance when compressing the GNN model at a high ratio due to the limited number of edges and parameters. In contrast, our innovative hierarchical KD offers enhanced preservation capabilities.
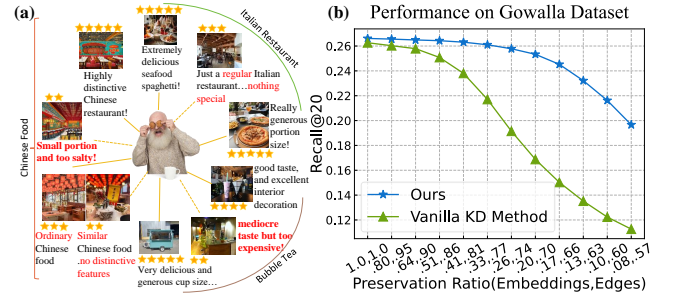
Fully aware of these challenges, we introduce a GNN pruning framework called LightGNN that facilitates efficient and denoised recommendations. LightGNN incorporates graph structure learning to explicitly assess the likelihood of redundancy or noise for each edge and embedding entry. This learning process is supervised in an end-to-end fashion, leveraging the downstream recommendation task alongside a hierarchical knowledge distillation paradigm. Inspired by the advantages of global relation learning in recommendation [30], our KD approach features an intermediate distillation layer that utilizes high-order relations to enhance candidate edges in the compressed model. This augmentation improves the model's capacity to maintain recommendation performance under high-rate compression. Through innovative importance distillation and prediction-level and embedding-level alignments, our hierarchical knowledge distillation enriches learnable pruning with abundant supervisory signals, boosting its compression capability.

The contributions of our LightGNN are summarized as follows:

- We introduce a novel GNN pruning framework for recommendation, explicitly identifying and eliminating redundancy and noise in GNNs to enable efficient and denoised recommendations.
- Our LightGNN framework integrates an innovative hierarchical knowledge distillation paradigm, seamlessly compressing GNNs at high ratios while preserving prediction accuracy.
- We conduct extensive experiments to demonstrate the superiority of LightGNN in terms of recommendation accuracy, inference efficiency, model robustness, and interpretability.

## 2 GNN-based Collaborative Filtering

Graph neural network (GNN) has been shown a most effective solution to collaborative filtering (CF) [4, 28]. The CF task typically involves a user set $\mathcal{U}$ ($|\mathcal{U}| = I$), an item set $\mathcal{V}$ ($|\mathcal{V}| = J$), and a user-item interaction matrix $\mathbf{A} \in \mathbb{R}^{I \times J}$. For a user $u_i \in \mathcal{U}$ and an item $v_j \in \mathcal{V}$, the entry $a_{i,j} \in \mathbf{A}$ equals 1 if user $u_i$ has interacted with item $v_j$, otherwise $a_{i,j} = 0$. Common interactions include users' rating, views, and purchases. GNN-based CF methods construct



Figure 1: Illustrations depicting (a) redundant and noisy user interactions, with red text indicating noisy feedback, and (b) the superior performance retention of LightGNN compared to vanilla KD, especially under high-rate pruning.

the user-item graph based on the interaction matrix $\mathbf{A}$. This graph can be denoted by $\mathcal{G} = (\mathcal{U}, \mathcal{V}, \mathcal{E})$, where $\mathcal{U}, \mathcal{V}$ serve as the graph vertices, and $\mathcal{E}$ denotes the edge set. For each $(u_i, v_j)$ that satisfies $a_{i,j} = 1$, there exists bidirectional edges $(u_i, v_j), (v_j, u_i) \in \mathcal{E}$.

Based on the user-item graph $\mathcal{G}$, GNNs conduct information propagation to smooth user/item embeddings for better reflecting the interaction data. Specifically, it firstly assigns initial embeddings $\mathbf{e}_i, \mathbf{e}_j \in \mathbb{R}^d$ to each user $u_i$ and item $v_j$, respectively. Here $d$ represents the hidden dimensionality. Then it iteratively propagates each node's embedding to its neighboring nodes for representation smoothing. Take the widely applied LightGCN [8] as an example, the embeddings for user $u_i$ and item $v_j$ in the $l$-th iteration are:

$$\mathbf{e}_{i,l} = \sum_{(v_j, u_i) \in \mathcal{E}} \frac{1}{\sqrt{d_i d_j}} \mathbf{e}_{j,l-1}, \quad \mathbf{e}_{j,l} = \sum_{(u_i, v_j) \in \mathcal{E}} \frac{1}{\sqrt{d_i d_j}} \mathbf{e}_{i,l-1} \quad (1)$$

where $\mathbf{e}_{i,l}, \mathbf{e}_{i,l-1} \in \mathbb{R}^d$ denote the embedding vectors for $u_i$ in the $l$-th and the $(l-1)$-th iterations, and analogous notations are used in $\mathbf{e}_{j,l}, \mathbf{e}_{j,l-1}$. The 0-th embedding vectors $\mathbf{e}_{i,0}, \mathbf{e}_{j,0}$ uses the initial embeddings $\mathbf{e}_i, \mathbf{e}_j$. And $d_i, d_j$ represent the degrees of nodes $u_i, v_j$, for Lapalacian normalization. After a total $L$ iterations, GNN-based CF aggregates the multi-order embeddings for final representations $\bar{\mathbf{e}}_i, \bar{\mathbf{e}}_j \in \mathbb{R}^d$ and user-item relation predictions $\hat{y}_{i,j}$, as follows:

$$\hat{y}_{i,j} = \bar{\mathbf{e}}_i^\top \bar{\mathbf{e}}_j, \quad \bar{\mathbf{e}}_i = \sum_{l=0}^{L} \mathbf{e}_{i,l}, \quad \bar{\mathbf{e}}_j = \sum_{l=0}^{L} \mathbf{e}_{j,l} \quad (2)$$

With the prediction scores $\hat{y}_{i,j}$, the GNN models are optimized by minimizing the BPR loss function [18] over all positive user-item pairs $(u_i, v_{j^+}) \in \mathcal{E}$, and sampled negative pairs $(u_i, v_{j^-})$, as follows:

$$\mathcal{L}_{bpr} = \sum_{(u_i, v_{j^+}, v_{j^-})} -\log \text{sigm}(\hat{y}_{i,j^+} - \hat{y}_{i,j^-}) \quad (3)$$

Though the above GNN framework achieves state-of-the-art performance in recommendation, its scalability is limited by the large-scale interaction graph and embedding table. In light of this, this paper proposes LightGNN aiming to effectively prune the GNN model for efficient graph neural collaborative filtering.

## 3 Methodology

This section goes through the proposed LightGNN to show the technical details. The overall framework is illustrated in Figure 2.
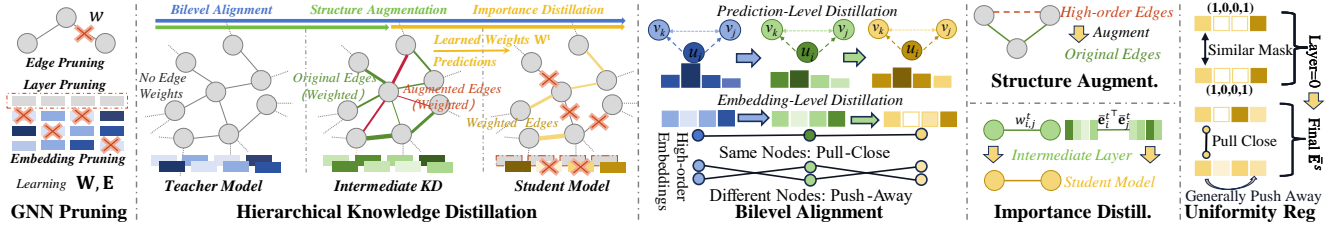
**Figure 2: Overall framework of the proposed LightGNN model.**

## 3.1 Graph Neural Network Pruning

Inspired by the lottery ticket hypothesis for GNNs [5, 6], we propose to use only a subset of GNN's parameters that maximally preserve the model functionality, to improve its efficiency. Specifically, the time complexity for a typical GNN model as aforementioned is $O(L \times |\mathcal{E}| \times d)$, and the space complexity is correspondingly $O(|\mathcal{E}| + (I + J) \times d)$. Therefore by reducing the number of edges $|\mathcal{E}|$, and the number of non-zero elements in the $d$ embedding dimensions, our LightGNN is able to optimize both the computational efficiency and memory efficiency. To achieve this, it is essential to identify the noisy and redundant parts in the edges $\mathcal{E}$ and the embedding table $\mathbf{E} = \{\mathbf{e}_i, \mathbf{e}_j | u_i \in \mathcal{U}, v_j \in \mathcal{V}\}$, to prevent performance degradation.

*3.1.1* **Edge Pruning.** To this end, LightGNN employs a sparse weight matrix $\mathbf{W} \in \mathbb{R}^{I \times J}$ for edge pruning. If an edge $(u_i, v_j)$ is a candidate for pruning, the corresponding weight $w_{i,j}$ in $\mathbf{W}$ is a learnable parameter. Otherwise $w_{i,j}$ is set as 0 and is not optimized. With the weight matrix $\mathbf{W}$, the graph information propagation process for the pruned GNN is conducted as follows:

$$\mathbf{E}_{\mathcal{U},l} = \mathbf{D}_{\mathcal{U}}^{-\frac{1}{2}} \cdot (\mathbf{A} \odot \mathbf{W}) \cdot \mathbf{D}_{\mathcal{V}}^{-\frac{1}{2}} \cdot \mathbf{E}_{\mathcal{V},l-1} + \mathbf{E}_{\mathcal{U},l-1} \quad (4)$$

where $\odot$ denotes the element-wise product operator which injects the learnable weights $\mathbf{W}$ into the information propagation process. Here $\mathbf{E}_{\mathcal{U},l}, \mathbf{E}_{\mathcal{U},l-1} \in \mathbb{R}^{I \times d}$ denote the user embedding table in the $l$-th and the $(l-1)$-th iteration, and $\mathbf{E}_{\mathcal{V},l-1} \in \mathbb{R}^{J \times d}$ denotes the embedding matrix for items in the $(l-1)$-th iteration. And $\mathbf{D}_{\mathcal{U}} \in \mathbb{R}^{I \times I}, \mathbf{D}_{\mathcal{V}} \in \mathbb{R}^{J \times J}$ denote the degree matrices for users and items, respectively. The information propagation to obtain higher-order item embeddings $\mathbf{E}_{\mathcal{V},l}$ is analogously using $(\mathbf{A} \odot \mathbf{W})^{\top}$.

Based on the parametric information propagation, the weights $\mathbf{W}$ participate in the calculation for final user/item embeddings, which are then used for predictions and loss calculations. Through the back propagation, $\mathbf{W}$ is tuned to reflect the importance of edges, wherein larger $|w_{i,j}|$ denotes the edge $(u_i, v_j)$ having a larger influence on producing better recommendation results. In light of this property, our LightGNN framework prunes the less important edges (noises or redundancies) after training, specifically by setting the $\rho\%$ candidate edges with the least importance to 0 (see 3.2.3), where $\rho \in (0, 100)$ denotes the proportion to drop. The pruning algorithm follows an iterative manner with multiple runs. In each run, LightGNN first conducts parameter optimization for model training and pruning weight tuning, and then prunes the GNN by dropping edges and other parameters.

*3.1.2* **Embedding and Layer Pruning.** As indicated by the complexity analysis for GNNs, the parameters for representing users and items (*i.e.* embeddings $\mathbf{E}$) also contribute significantly to the

running time and the memory costs of GNNs. Therefore LightGNN follows the similar pruning algorithm for edges to prune the entries in the embedding matrix $\mathbf{E}$. As the scalar parameters in $\mathbf{E}$ already reflect the importance of their corresponding entries, LightGNN does not employ extra pruning weights for embeddings. Analogously, LightGNN alternately conducts model training and parameter pruning with ratio $\rho'\%$ according to the absolute value $|e_{i,d'}|$, where $e_{i,d'}$ represents the $d'$-th dimension in $i$'s embedding vector.

In addition to the edges and embeddings, the time complexity of GNNs suggests that the number of graph propagation layers $L$ also greatly impacts the computation time of GNNs. Moreover, in practice, $L$ is also significant to influence the temporary memory costs for stacking the intermediate results. Thus our LightGNN further reduces the number of graph iterations $L$ for efficiency, which also alleviates the over-smoothing effect of GNNs [30].

## 3.2 Hierarchical Knowledge Distillation

*3.2.1* **Bilevel Alignment.** Motivated by the strength of knowledge distillation (KD) in compressing the learned knowledge of advanced models into light-weight architectures [29], the proposed LightGNN develops a hierarchical knowledge distillation framework to maximally retain the original high performance in the pruned GNN model. Taking a well-trained GNN model (*e.g.* Light-GCN [8]) as the teacher, LightGNN aligns the student model with pruned structures, embeddings, and GNN layers to the teacher model with respect to both hidden embeddings and final predictions. In the prediction level, the following loss function is applied:

$$\mathcal{L}_{p-kd} = \sum_{\mathbf{v}} -\left( \sigma(\epsilon_{\mathbf{v}}^t/\tau) \cdot \log \sigma(\epsilon_{\mathbf{v}}^s/\tau) + \overline{\sigma}(\epsilon_{\mathbf{v}}^t/\tau)) \cdot \log \overline{\sigma}(\epsilon_{\mathbf{v}}^s/\tau) \right)$$

$$\text{where } \mathbf{v} = (u_i, v_{j^1}, v_{j^2}), \ \overline{\sigma}(x) = 1 - \sigma(x), \ \epsilon_{\mathbf{v}}^* = \hat{y}_{i,j^1}^* - \hat{y}_{i,j^2}^* \quad (5)$$

Here $(u_i, v_{j^1}, v_{j^2})$ denotes the randomly sampled training tuples analogous to the BPR loss, while $v_{j^1}$ and $v_{j^2}$ are not fixed to be positive or negative samples. $\sigma(\cdot)$ denotes the sigmoid function to constrain the values to be within $(0, 1)$. And $\tau \in \mathbb{R}$ is known as the temperature coefficient [10]. We denote the predictions made by the student model using the superscript $s$, and denote the predictions made by the teacher model with the superscript $t$. With this training objective, our LightGNN framework encourages the pruned GNN model to mimic the predictions made by the complete GNN model with all the edges, embedding entries and propagation iterations, to obtain the teacher's prediction ability as much as possible.

Besides the prediction-level alignment, our LightGNN aligns the teacher model and the student model by treating their learned embeddings as paired data views for contrastive learning. In specific,

the following infoNCE loss function [16] is applied:

$$\mathcal{L}_{e-kd} = -\sum_{u_i \in \mathcal{U}} \log \operatorname{softmax}(\mathbf{S}_\mathcal{U}, u_i) - \sum_{v_j \in \mathcal{V}} \log \operatorname{softmax}(\mathbf{S}_\mathcal{V}, v_j)$$

$$\text{where } \operatorname{softmax}(\mathbf{S}_\mathcal{U}, u_i) = \frac{\exp s_{i,i}}{\sum_{u_{i'}} \exp s_{i',i}}, \quad s_{i',i} = \cos(\bar{\mathbf{e}}_{i'}^s, \bar{\mathbf{e}}_i^t) \quad (6)$$

Here $s_{i',i} \in \mathbf{S}_\mathcal{U}$ denotes the cosine similarity between the final embeddings $\bar{\mathbf{e}}_{i'}^s, \bar{\mathbf{e}}_i^t$ for the users $u_{i'}$ and $u_i$, given by the student model and the teacher model, respectively. The item-side embedding-level KD is calculated analogously. With this embedding-level KD objective, our LightGNN can better guide the pruned GNN to preserve the essential graph structures and parameters in a deeper level.

### 3.2.2 Intermediate KD Layer for Structure Augmentation.
Due to the sparsity nature of the user-item interaction data, some key preference patterns are not reflected by the direct neighboring relations but preserved by the high-order relations. To facilitate the capturing of these high-order connections during our edge pruning, we augment the knowledge distillation of LightGNN with an intermediate KD layer model for edge augmentation.

To be specific, LightGNN conducts a two-stage distillation, firstly from the original GNN to an augmented GNN, and then from the augmented GNN to the final pruned GNN. The augmented GNN does not prune any edges or embedding entries, but instead includes the high-order connections as augmented edges. Formally, the augmented GNN has the same model architecture (Eq. 4) as the student but works over the following augmented interaction graph:

$$\bar{\mathcal{G}} = (\mathcal{U}, \mathcal{V}, \bar{\mathcal{E}}), \quad \bar{\mathcal{E}} = \{(u_i, v_j), (v_j, u_i) | \bar{a}_{i,j}^{(h)} \neq 0\} \quad (7)$$

where $\bar{a}_{i,j}^{(h)}$ denotes the entry for $(u_i, v_j)$ in the $h$-th power of the symmetric adjacent matrix with self loop [25]. In other words, edge $(u_i, v_j)$ exists in the augmented graph $\bar{\mathcal{G}}$ if $u_i$ can be connected to $v_j$ via any path with its length shorter than or equal to $h$ hops in the original graph. With this structure augmentation, the augmented GNN directly includes the high-order connections in the model parameters, to prevent losing the key high-order patterns in radical edge pruning. During the intermediate KD, the augmented GNN is supervised by the original GNN (no weights), not only to mimic its accurate predictions, but also to learn proper weights $\mathbf{W}^t$ for all the edges. The intermediate KD layer prevents the augmented larger graph from introducing noises using the supervision of the bilevel distillation from original GNN and the adaptive edge weights.

### 3.2.3 Importance Distillation for Pruning.
After the first knowledge distillation from the original GNN to the augmented GNN model, our LightGNN then distills its learned knowledge with structure augmentation to the final pruned GNN model. Apart from the aforementioned bilevel alignment, LightGNN further enhances this second KD with the importance distillation, which explicitly leverages the learned importance weights in the intermediate model to increase the precision of pruning weights in the final model. Specifically, the pruning weight matrix in the final pruned GNN is a compound variable whose entries are calculated as follows:

$$\bar{w}_{i,j}^s = w_{i,j}^s + \beta_1 \cdot w_{i,j}^t + \beta_2 \cdot \sigma(\bar{\mathbf{e}}_i^{t\top} \bar{\mathbf{e}}_j^t) \quad \text{for } (u_i, v_j) \in \mathcal{E} \quad (8)$$

where $\bar{w}_{i,j}^s \in \mathbb{R}$ denotes the weight to decide if edge $(u_i, v_j)$ should be pruned, and it is acquired using the independent edge weight

$w_{i,j}^s \in \mathbf{W}^s$ of the final student model, the tuned edge weight $w_{i,j}^t \in \mathbf{W}^t$ of the intermediate GNN as the teacher model, and the edge prediction made by the intermediate GNN's final embeddings $\bar{\mathbf{e}}_i^t, \bar{\mathbf{e}}_j^t \in \mathbb{R}^d$. Here $\beta_1, \beta_2$ denote two hyperparameters for weighting and we define the sparse decision matrix $\bar{\mathbf{W}}^s = \{\bar{w}_{i,j}^s\}_{I \times J}$.

With this importance distillation in the edge pruning, the pruning weights $\bar{\mathbf{W}}^s$ in the final student model are not only trained in the end-to-end manner using the bilevel KD objectives, but also directly adjusted by the well-trained weights in the intermediate teacher model. Moreover, by utilizing the edge weights obtained in the augmented graph, the pruned GNN is injected with the high-order connectivity to facilitate edge dropping and global relation learning. It is worth noting that, apart from the edge pruning, the student's edge weights are also employed in the graph information propagation, to enrich the pruned GNN with less edges but compensatory, adaptive and informative edge importance.

## 3.3 Optimization with Uniformity Constraint
Inspired by the advantage of learning uniform embeddings in CF [22, 28], our LightGNN proposes to regularize the model optimization with an adaptive uniformity constraint based on contrastive learning. In specific, the constraint minimizes the pairwise inner-product between embeddings to enforce representation uniformity, while maximizing the embedding similarity between nodes with similar pruning masks. In this way, the positive relations are augmented by the learned pruning weights for enhancement. Formally, the adaptive uniformity constraint is as follows:

$$\mathcal{L}_{u-reg} = \sum_{u_i \in \mathcal{U}} \left( -\log \frac{\sum_{u_{i^1} \in \mathcal{S}_i} \exp\left(\bar{\mathbf{e}}_i^{s\top} \bar{\mathbf{e}}_{i^1}^s / \tau\right)}{\sum_{u_{i^2} \in \mathcal{U}} \exp\left(\bar{\mathbf{e}}_i^{s\top} \bar{\mathbf{e}}_{i^2}^s / \tau\right)} \right)$$
$$+ \sum_{v_j \in \mathcal{V}} \left( -\log \frac{\sum_{v_{j^1} \in \mathcal{S}_j} \exp\left(\bar{\mathbf{e}}_j^{s\top} \bar{\mathbf{e}}_{j^1}^s / \tau\right)}{\sum_{v_{j^2} \in \mathcal{V}} \exp\left(\bar{\mathbf{e}}_j^{s\top} \bar{\mathbf{e}}_{j^2}^s / \tau\right)} \right) \quad (9)$$

where $\mathcal{S}_i$ and $\mathcal{S}_j$ denote the positive sets of user $u_i$ and item $v_j$, respectively, which are determined by picking the users/items that share the highest similarity in embedding pruning. Take the user side as an example, the neighborhood set $\mathcal{S}_i$ is acquired by:

$$\mathcal{S}_i = \left\{ u_{i^1} \mid \|\mathbb{o}_i \odot \mathbb{o}_{i^1}\|_0 \geq \max\left(\|\mathbb{o}_i\|_0, \|\mathbb{o}_{i^1}\|_0\right) - \delta \right\} \quad (10)$$

where $\mathbb{o}_i, \mathbb{o}_{i^1} \in \{0, 1\}^d$ denote binary pruning masks for the 0-th embedding vectors $\mathbf{e}_i^s$ and $\mathbf{e}_{i^1}^s$, respectively. Operator $\odot$ denotes the element-wise multiplication, and $\| * \|_0$ denotes the $l_0$ norm of vectors. $\delta$ represents the threshold hyperparameter for similarity relaxation, which is selected according to the pruning ratio.

With the above contrastive loss using similarly-pruned embeddings as positive sets, LightGNN can learn uniformly-distributed embeddings while capturing the node-wise similarity during the pruning process. Combining it with the collaborative filtering loss $\mathcal{L}_{bpr}$, the bilevel KD losses $\mathcal{L}_{p-kd}$ and $\mathcal{L}_{e-kd}$, and a weight-decay regularization term over parameters $\Theta$, LightGNN applies the following multi-task training loss with hyperparameters $\lambda_*$:

$$\mathcal{L} = \lambda_0 \mathcal{L}_{bpr} + \lambda_1 \mathcal{L}_{p-kd} + \lambda_2 \mathcal{L}_{e-kd} + \lambda_3 \mathcal{L}_{u-reg} + \lambda_4 \|\Theta\|_F^2. \quad (11)$$

**Table 1: Statistical details of experimental datasets.**

| Dataset | # Users | # Items | # Interactions | Interaction Density |
|---------|---------|---------|----------------|---------------------|
| Gowalla | 25557 | 19747 | 294983 | $5.85 \times 10^{-4}$ |
| Yelp | 42712 | 26822 | 182357 | $1.59 \times 10^{-4}$ |
| Amazon | 76469 | 83761 | 966680 | $1.51 \times 10^{-4}$ |

## 4 Evaluation

We conduct extensive experiments on our LightGNN framework, aiming to answer the following research questions (RQs):

- **RQ1**: How is the performance of LightGNN after the model pruning, compared to existing recommendation methods?
- **RQ2**: How efficient is our pruned GNN, compared to baselines?
- **RQ3**: How do the components of the proposed LightGNN impact the recommendation performance of the pruned GNN?
- **RQ4**: How do the pruning ratios impact the recommendation performance and the efficiency of the pruned GNN?
- **RQ5**: Can the proposed LightGNN framework alleviate the over-smoothing effect with its hierarchical knowledge distillation?
- **RQ6**: Can our LightGNN effectively identify the redundant and noisy information in the user-item interaction graph?

### 4.1 Experimental Settings

*4.1.1* **Datasets.** LightGNN is evaluated using three real-world datasets: Gowalla, Yelp, and Amazon. The **Gowalla** dataset contains user check-in records at geographical locations from January to June 2010, obtained from the Gowalla platform. **Yelp** dataset is obtained from Yelp platform and contains user ratings on venues from January to June 2018. The **Amazon** dataset contains people's ratings of books on the Amazon platform, during 2013. Following [29], we filter out users and items with less than three interactions, and splitting the original datasets into training, validation, and test sets by 70:5:25. Additionally, we convert ratings into binary implicit feedback, following [8]. The data statistics are listed in Table 1.

*4.1.2* **Evaluation Protocols.** We follow common evaluation protocols for recommendation [25, 35]. We rank all uninteracted items with the positive items from test set for each user, a method known as full-rank evaluation. We use two common metrics, *Recall@N* and *NDCG@N* [24, 27] with values of $N = 20$ and 40.

*4.1.3* **Baselines.** We compare LightGNN to 18 baselines from diverse categories, including factorization method (**BiasMF** [13]), deep neural CF methods (**NCF** [9], **AutoR** [19]), graph-based methods (**GCMC** [1], **PinSage** [33], **STGCN** [36], **NGCF** [25], **GCCF** [4], **LightGCN** [8], **DGCF** [26]), self-supervised recommenders (**SLRec** [32], **SGL** [27], **NCL** [15], **SimGCL** [34], **HCCF** [30]), and compressed CF approaches (**GLT** [5], **UnKD** [3], **SimRec** [29]).

*4.1.4* **Hyperparameter Settings.** We implement LightGNN with PyTorch, using Adam optimizer and Xavier initializer with default parameters. For all models, the training batch size is set to 4096 and the embedding size is 32 by default. For all GNN-based models, we set the layer number to 2. Weights $\lambda_0, \lambda_1, \lambda_2$ in LightGNNare tuned from $\{1e^{-k}|k = 0, 1, ..., 4\}$. And $\lambda_3$ is tuned in a wider range which additionally contains $\{1e^{-5}, 1e^{-6}\}$. The weight $\lambda_4$ for weight-decay regularization is selected from $\{1e^{-k}|k = 3, 4, ..., 9\}$. All temperature coefficients are chosen from $\{1e^{-k}, 3e^{-k}, 5e^{-k}|k = -1, 0, 1, 2\}$.

Baseline methods are implemented using their released code with grid search for hyperparameter tuning. The efficiency test is conducted on a device with an NVIDIA GeForce RTX 3090 GPU.

### 4.2 Performance Comparison (RQ1)

We first compare LightGNN to baselines on recommendation accuracy. The results are in Table 2. We make the following observations:

- **Superior performance of LightGNN**: The proposed model LightGNN surpasses all baselines across different categories, including simple neural CF, graph-based recommenders, self-supervised methods, and compression methods. This superiority in performance demonstrates that our learnable pruning framework and hierarchical distillation paradigm not only maintain prediction accuracy after model compression but also enhance existing recommendation frameworks. The effective elimination of noise and redundancy in the interaction graph and embedding parameters contributes to these performance improvements.

- **Drawbacks of CF without model compression**: When comparing the best-performing CF methods, such as self-supervised CF techniques like SGL, HCCF, and SimGCL, to compression methods like UnKD and SimRec, it is evident that CF methods without model compression fall short in terms of recommendation accuracy. This discrepancy can be attributed to the debiasing and anti-over-smoothing effects embedded in the knowledge distillation process of UnKD and SimRec. This suggests that model compression techniques, such as knowledge distillation, can go beyond improving model efficiency. They can also address adverse factors present in observed data and modeling frameworks, such as data bias, noise, and over-smoothing effects.

- **Importance of explicit noise elimination**: While UnKD and SimRec refine the distilled model by addressing bias and over-smoothing effects in GNN-based CF, they rely solely on high-level supervision methods. In contrast, our LightGNN explicitly identifies and eliminates fine-grained noisy and redundant elements within the model, such as edges and embedding entries. This empowers our LightGNN with notable strength in recommender refinement, leading to significant performance superiority.

### 4.3 Efficiency Test (RQ2)

To assess the model efficiency, we evaluate the memory and computational costs of LightGNN and baselines. The compared baselines include NGCF, GCCF, HCCF, and existing GNN compression method UnKD. Our LightGNN is tested with different preservation ratios. In Figure 3, the results are presented relative to the performance of NGCF. We deduce the following observations:

- **Simplified GNNs**. Despite simplifying the GNN architecture by removing transformations and activations, some GNN methods like GCCF fail to significantly reduce memory and time costs related to graph storage and information propagation. Consequently, the costs of GCCF remain comparable to those of NGCF. This demonstrates the limitation of architectural simplifications in improving efficiency for graph-based recommendation.

- **SSL-enhanced GNNs**. SSL techniques have been utilized to enhance graph recommenders by generating self-supervision signals. However, it is important to note that these methods may

**Table 2: Overall performance comparison on Gowalla, Yelp, and Amazon datasets in terms of *Recall@N* and *NDCG@N***

| Data | Metric | BiasMF | NCF | AutoR | PinSage | STGCN | GCMC | NGCF | GCCF | LightGCN | DGCF | SLRec | NCL | SGL | HCCF | SimGCL | GLT | UnKD | SimRec | Ours |
|------|--------|--------|-----|-------|---------|-------|------|------|------|----------|------|-------|-----|-----|------|--------|-----|------|--------|------|
| Amazon | Recall@20 | 0.0324 | 0.0367 | 0.0525 | 0.0486 | 0.0583 | 0.0837 | 0.0551 | 0.0772 | 0.0868 | 0.0617 | 0.0742 | 0.0955 | 0.0874 | 0.0885 | 0.0921 | 0.0901 | 0.0947 | 0.1067 | **0.1189** |
| | NDCG@20 | 0.0211 | 0.0234 | 0.0318 | 0.0317 | 0.0377 | 0.0579 | 0.0353 | 0.0501 | 0.0571 | 0.0372 | 0.0480 | 0.0623 | 0.5690 | 0.0578 | 0.0605 | 0.0585 | 0.0607 | 0.0734 | **0.0820** |
| | Recall@40 | 0.0578 | 0.0600 | 0.0826 | 0.0773 | 0.0908 | 0.1196 | 0.0876 | 0.1175 | 0.0912 | 0.1123 | 0.1409 | 0.1312 | 0.1335 | 0.1367 | 0.1355 | 0.1376 | 0.1535 | **0.1677** |
| | NDCG@40 | 0.0293 | 0.0306 | 0.0415 | 0.0402 | 0.0478 | 0.0692 | 0.0454 | 0.0625 | 0.0697 | 0.0468 | 0.0598 | 0.0764 | 0.0704 | 0.0716 | 0.0730 | 0.0725 | 0.0745 | 0.0879 | **0.0969** |
| Gowalla | Recall@20 | 0.0867 | 0.1019 | 0.1477 | 0.0985 | 0.1574 | 0.1863 | 0.1757 | 0.2012 | 0.2230 | 0.2055 | 0.2001 | 0.2283 | 0.2332 | 0.2293 | 0.2328 | 0.2324 | 0.2331 | 0.2434 | **0.2610** |
| | NDCG@20 | 0.0579 | 0.0674 | 0.0690 | 0.0809 | 0.1042 | 0.1151 | 0.1135 | 0.1282 | 0.1433 | 0.1312 | 0.1298 | 0.1478 | 0.1509 | 0.1482 | 0.1506 | 0.1464 | 0.1496 | 0.1592 | **0.1684** |
| | Recall@40 | 0.1269 | 0.1563 | 0.2511 | 0.1882 | 0.2318 | 0.2627 | 0.2586 | 0.2903 | 0.3181 | 0.2929 | 0.2863 | 0.3232 | 0.3251 | 0.3258 | 0.3276 | 0.3269 | 0.3301 | 0.3399 | **0.3597** |
| | NDCG@40 | 0.0695 | 0.0833 | 0.0985 | 0.0994 | 0.1252 | 0.1390 | 0.1367 | 0.1532 | 0.1670 | 0.1555 | 0.1540 | 0.1745 | 0.1780 | 0.1751 | 0.1772 | 0.1730 | 0.1766 | 0.1865 | **0.1962** |
| Yelp | Recall@20 | 0.0198 | 0.0304 | 0.0491 | 0.0510 | 0.0562 | 0.0584 | 0.0681 | 0.0742 | 0.0761 | 0.0700 | 0.0665 | 0.0806 | 0.0803 | 0.0789 | 0.0788 | 0.0812 | 0.0819 | 0.0823 | **0.0879** |
| | NDCG@20 | 0.0094 | 0.0143 | 0.0222 | 0.0245 | 0.0282 | 0.0280 | 0.0336 | 0.0365 | 0.0373 | 0.0347 | 0.0327 | 0.0402 | 0.0398 | 0.0391 | 0.0395 | 0.0400 | 0.0392 | 0.0414 | **0.0443** |
| | Recall@40 | 0.0307 | 0.0487 | 0.0692 | 0.0743 | 0.0856 | 0.0891 | 0.1019 | 0.1151 | 0.1175 | 0.1072 | 0.1032 | 0.1230 | 0.1226 | 0.1210 | 0.1213 | 0.1249 | 0.1202 | 0.1251 | **0.1328** |
| | NDCG@40 | 0.0120 | 0.0187 | 0.0268 | 0.0315 | 0.0355 | 0.0360 | 0.0419 | 0.0466 | 0.0474 | 0.0437 | 0.0418 | 0.0505 | 0.0502 | 0.0492 | 0.0498 | 0.0507 | 0.0493 | 0.0519 | **0.0553** |



(a) Storage costs.　　　　(b) Time costs.

**Figure 3: Disk storage and time costs of baselines and our LightGNN under different preservation ratios (*e.g.* .33, .44 denote preserving 33% embedding entries and 44% edges).**

introduce additional operations, leading to increased memory and time costs. This is exemplified by the performance of HCCF, where utilizing extra hypergraph propagation necessitates more FLOPs and yields a noticeable increase in computational time.

- **Existing compressed GNNs**. UnKD has been successful in achieving efficiency improvements, particularly in terms of computational time. However, when comparing UnKD to LightGNN, a significant disadvantage becomes evident. This limitation arises from UnKD's lack of explicit identification and removal of redundancy and noise in the GNN model. As a result, UnKD is unable to prune a larger portion of the GNN to achieve superior efficiency improvements like our LightGNN framework does.

- **Efficiency of LightGNN**. The results demonstrate a significant memory reduction of 70% in LightGNN, considering both the parameter number and storage size. Moreover, there is an impressive reduction of over 90% in FLOPs during forward propagation and an over 50% reduction in physical prediction time. These efficiency optimizations can be attributed to two key aspects. **Firstly**, the learnable GNN pruning paradigm accurately removes redundant and noisy information from the GNN. This facilitates efficient utilization of computational resources. **Secondly**, our learnable pruning mechanism is supervised by the hierarchical KD, which incorporates multi-dimensional alignment and high-order structure augmentation. This maximizes the retention of performance, allowing for more extensive pruning of parameters.

## 4.4　Ablation Study (RQ3)

We investigate the effectiveness of LightGNN's technical designs using Gowalla and Yelp data, with different pruning ratios. The results are shown in Table 3. We make the following observations.

**Table 3: Ablation study of LightGNN measured by Recall@20.**

| Dataset | | Gowalla | | | | Yelp | | | |
|---------|----|---------|---------|---------|---------|---------|---------|---------|---------|
| Ratio $\mathbf{E}/\mathcal{E}$ | | .33/.44 | .26/.37 | .11/.19 | .08/.16 | .33/.77 | .26/.74 | .11/.60 | .08/.57 |
| Prn | ~EmbP | 0.2197 | 0.1946 | 0.0741 | 0.0586 | 0.0809 | 0.0737 | 0.0434 | 0.0357 |
| | ~EdgeP | 0.2418 | 0.2255 | 0.1341 | 0.1133 | 0.0867 | 0.0850 | 0.0745 | 0.0709 |
| | ~BothP | 0.1800 | 0.1434 | 0.0556 | 0.0421 | 0.0736 | 0.0661 | 0.0311 | 0.0231 |
| | BnEdge | 0.2210 | 0.2021 | 0.1280 | 0.1165 | 0.0872 | 0.0858 | 0.0775 | 0.0754 |
| KD | -BiAln | 0.2350 | 0.2281 | 0.1934 | 0.1812 | 0.0810 | 0.0801 | 0.0666 | 0.0650 |
| | -IntKD | 0.2607 | 0.2571 | 0.2100 | 0.1890 | 0.0822 | 0.0825 | 0.0788 | 0.0769 |
| | -ImpD | 0.2593 | 0.2564 | 0.2135 | 0.1923 | 0.0862 | 0.0861 | 0.0808 | 0.0797 |
| LightGNN | | 0.2610 | 0.2578 | 0.2162 | 0.1966 | 0.0879 | 0.0877 | 0.0856 | 0.0842 |

**Effectiveness of the GNN pruning techniques**.

- **~EmbP**, **~EdgeP**, **~BothP**: We replace the learnable pruning with random dropping. The three variants replace embedding pruning, edge pruning, and both, respectively. Significant performance drop can be observed under different pruning ratios, indicating the effectiveness of our learnable pruning in identifying the essential embedding entries and edges. Especially, when dropping with high ratios (*e.g.* preserving only 11% and 8% entries), the prediction ability of the random variants experiences a destructive (over 70%) decay, while LightGNN preserves most of its accuracy.

- **BnEdge**: To study the effect of learned edge weights $\mathbf{W}^s$, BnEdge uses binary edge weights instead of $\mathbf{W}^s$ during GNN propagation. Though it maintains the learnable pruning process unchanged, a noticeable degradation can be observed. This suggests the crucial role of learned weights. They not only identify which edges to prune, but also effectively preserve the pruned information.

**Effectiveness of knowledge distillation**.

- **-BiAln**: To assess the significance of the KD constraints for effective pruning, we remove the bilevel alignment, including the prediction-level and embedding-level KD. The notable performance drop verifies the importance of aligning the teacher model with the pruned model, to effectively retain model performance.

- **-IntKD**: This variant removes the intermediate KD layer in LightGNN. As a result, its performance notably deteriorates, particularly on the Yelp dataset. The increased importance of this module for Yelp can be attributed to the higher sparsity of the dataset. In such cases, the intermediate KD layer is able to seek more edges from high-order relations to enrich the small edge set.

- **-ImpD**: This variant removes the importance distillation, and the results confirm the benefits of incorporating learned edge weights and predictions from the intermediate KD layer model into the decision-making process of edge dropping.
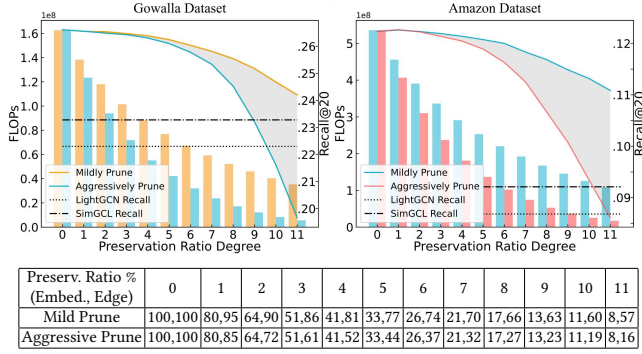
**Figure 4: Computational FLOPs (bars) and recommendation performance (lines) *w.r.t* different preservation ratios.**

| Preserv. Ratio % (Embed., Edge) | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mild Prune | 100,100 | 80,95 | 64,90 | 51,86 | 41,81 | 33,77 | 26,74 | 21,70 | 17,66 | 13,63 | 11,60 | 8,57 |
| Aggressive Prune | 100,100 | 80,85 | 64,72 | 51,61 | 41,52 | 33,44 | 26,37 | 21,32 | 17,27 | 13,23 | 11,19 | 8,16 |

## 4.5 Influence of Pruning Ratios (RQ4)

In this experiment, we investigate the impact of pruning ratios for edges and embedding entries on both model performance and efficiency. Figure 4 shows the evaluated model performance and computing FLOPs (floating point operations) during forward propagation across various preservation ratios. We present two pruning schemes: a mild pruning scheme that removes fewer graph edges in the GNN, and an aggressive pruning scheme that removes more edges. Based on the results, we draw the following observations:

- **Performance change**. As we discard a larger number of embedding entries and graph edges, we observe a continuous decline in performance. However, it is noteworthy that even when a substantial portion of the GNN model is removed, our LightGNN consistently maintains a high level of recommendation performance compared to SimGCL and LightGCN. This resilience can be attributed to the hierarchical KD, which effectively aligns the predictions of the student model with those of the well-performing teacher model through bilevel alignment, and the importance distillation that gives the optimal dropping strategies. Additionally, the intermediate KD layer with structure augmentation further enhances the recommendation ability by incorporating more edges sampled from high-order relations. These features collectively contribute to the robust performance of LightGNN.

- **Efficiency change**. As the pruning ratio increases, LightGNN exhibits a significant decrease in FLOPs. This confirms the effectiveness of enhancing GNN efficiency by pruning embeddings and structures. Specifically, our LightGNN achieves a FLOPs reduction of 90% during forward propagation while maintaining comparable performance to SimGCL, and a FLOPs reduction of 95% while performing similarly to LightGCN. These substantial reductions in FLOPs highlight the effectiveness of our learnable pruning strategy in minimizing computational operations.

- **Differences across datasets**. Furthermore, it is worth mentioning that our LightGNN demonstrates better preservation of recommendation performance when pruning the same proportion of information on the Amazon dataset compared to the Gowalla dataset. This observation suggests the presence of more redundancy or noise in the Amazon data, which aligns with the larger number of edges and users/items present in the Amazon dataset.
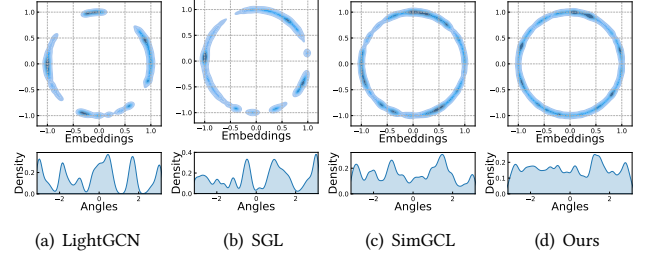


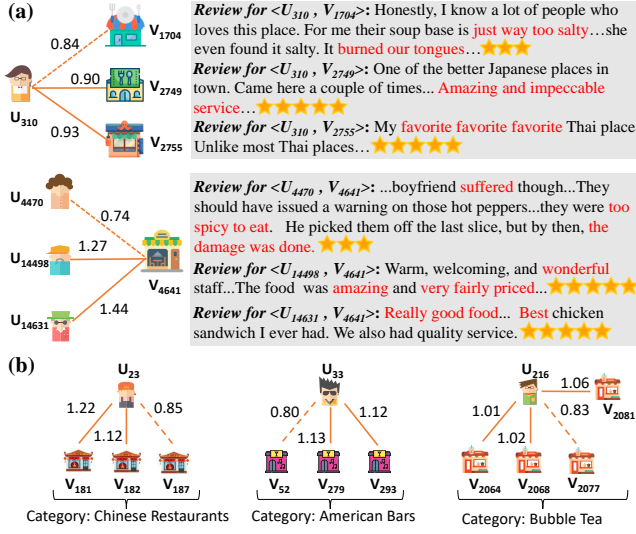**Figure 5: Embedding distributions in the 2-D space and in the 1-D angle space for Yelp dataset, estimated by KDE.**

**Table 4: MAD among popular nodes from Yelp and Gowalla.**

| Datasets | GCCF | LightGCN | NCL | SGL | SimGCL | SimRec | Ours |
|---|---|---|---|---|---|---|---|
| Yelp | 0.8747 | 0.8819 | 0.8929 | 0.8643 | 0.9103 | 0.9272 | **0.9404** |
| Gowalla | 0.8206 | 0.8269 | 0.8236 | 0.7760 | 0.8595 | 0.8406 | **0.8742** |

## 4.6 Anti-Over-Smoothing Effect Study (RQ5)

To assess the ability of LightGNN to mitigate the over-smoothing effect of GNNs during the pruning process, we compare the distribution uniformity of our model's embeddings with those of baseline methods. This comparison is conducted in two dimensions.

- **Visualization of Embedding Distribution**. From the embedding distributions plotted in Figure 5, we can observe that: **i)** The clustering effect observed in both 2-D plots and angle plots is notably stronger for LightGCN, demonstrating the severe over-smoothing effect resulting from the iterative embedding smoothing paradigm. **ii)** To address this issue, SGL and SimGCL incorporate contrastive learning to enhance the distribution uniformity of embeddings. Both methods exhibit higher uniformity in the estimated distribution compared to LightGCN, with SimGCL exhibiting some superiority due to its less-random augmentation design. **iii)** Compared to SimGCL, our LightGNN exhibits even fewer dark regions in the embedding distribution ring, indicating higher uniformity. This advantage becomes more apparent in the angle-based plot, where the low probabilities are much closer to the high ones in LightGNN. This observation strongly indicates a higher anti-over-smoothing ability of our LightGNN, which can be ascribed to the sparsification effect caused by embedding pruning, and the uniformity constraint in our LightGNN.

- **Mean Average Distance (MAD) Values**. We further evaluate the MAD values [2, 29] in Table 4, from which we draw the following observations: **i)** The GNN-based CF paradigms GCCF and LightGCN generally exhibit lower MAD values compared to other methods that employ contrastive learning. This observation highlights the inherent over-smoothing issue in propagation-based graph neural encoders. **ii)** For the other baselines, we observe that NCL and SGL show lower MAD values, indicating a stronger over-smoothing effect. This sheds light on the limitations of their random structure augmentation methods, which are susceptible to the influence of data noise. **iii)** The superiority of SimGCL and SimRec validates their effective design of pushing all embeddings apart. In comparison, our LightGNN achieves further advancements by constructing meaningful positive sample pairs using node-wise similarity in embedding pruning. This technique effectively enhances positive relation learning in a learnable manner.

**Figure 6: Investigation on the capability of (a) noise pruning and (b) redundancy pruning for our LightGNN framework.**

## 4.7 Noise and Redundancy Identification (RQ6)

We explore LightGNN's capacity to trim noise and redundancy in interaction data. The results are detailed in Figure 6.

**Noise Pruning**. In Figure 6(a), two sets of decision weights in $\bar{\mathbf{W}}^s$ for left-side edges are depicted alongside users' text reviews and ratings for corresponding items on the right. Notably, these reviews and ratings were not exposed to our LightGNN. Our results show that LightGNN assigns low weights to interactions like $< u_{310} , v_{1704} >$ and $< u_{4470} , v_{4641} >$, aligning with users' negative feedback (e.g., "too salty."). In the context of graph CF, such negative feedback instances are viewed as regular user-item interactions, possibly adversely affecting user preference modeling. Frequent similar observations in our results show that LightGNN effectively identifies and addresses noise in the graph structure, thereby improving the pruning effect of GNN-based recommendation.

**Redundancy Pruning**. In Figure 6(b), some representative cases demonstrate the efficacy of redundancy pruning in LightGNN, where three users interact with multiple venues sharing the same categories like Chinese restaurants and American bars, reflecting redundant user interest information. Despite being category-agnostic, LightGNN identifies these similarities, assigning lower weights to some of the redundant items. This encourages the pruning algorithm to eliminate the redundancy, thereby enhancing model efficiency. Moreover, thanks to the learnable edge weights in the intermediate KD layer, LightGNN preserves preference strength for each interest, rather than relying on item counts of each interest.

## 5 Related Work

### 5.1 Graph Neural Recommender Systems

Graph neural networks (GNNs) have emerged as foundational architectures for recommendation systems. Early works such as NGCF [25] and GCMC [1] introduced graph convolutional networks (GCNs) for collaborative recommendation. Subsequent studies include STGCN [36], which integrates an autoencoding architecture

within the GNN encoder, and DGCF [26], which incorporates a representation disentanglement module into graph-based collaborative filtering. LightGCN [8] and GCCF [4] emphasize the redundancy in prior graph neural architectures and achieve improved performance by eliminating both non-linear and linear mappings.

Recently, self-supervised learning (SSL) has gained attention for its ability to generate rich supervision signals and address the data sparsity problem in recommendation. Contrastive learning (CL)-based graph CF (*e.g.* SGL [27], SimGCL [34], DirectAU [22], AdaGCL [12]) is a popular SSL technique that effectively learns a uniform distribution to counter the over-smoothing effect of GNNs. HCCF [30] and NCL [15] introduce additional encoding views to enrich graph CL. In addition, graph-based recommendation has also been enhanced with generative SSL techniques based on masked autoencoding, such as AutoCF [28] and DGMAE [17].

Despite the substantial enhancements in recommendation performance due to GNN advancements, an inherent limitation remains in the inefficiency of GNN's extensive information propagation and node-specific parameters. In this context, our LightGNN aims to effectively prune redundant and noisy components of GNNs while preserving high performance through distillation constraints.

### 5.2 Model Compression for Graph Models

To enhance the scalability of GNNs, prior works have utilized random node and edge sampling techniques for large graphs (e.g., PinSAGE [33], HGT [11]). However, these random strategies do not ensure the preservation of crucial information and may significantly affect model performance. In response, several approaches have emerged to better retain important patterns from the original model. GLT [5] advocates for preserving only the essential edges by learning their importance to downstream task performance. Other studies improve compression supervision through knowledge distillation. GLNN [37] and SimRec [29] propose distilling efficient student models based on MLP from heavier GNNs. UnKD [3] further mitigates bias in the KD process using a stratified distillation strategy. Additionally, KD has been applied to compress recommenders based on non-GNN architectures (e.g., [21, 31]).

In contrast to previous approaches that broadly reduce model complexity by substituting GNNs with simpler architectures, our LightGNN preserves robust topology extraction capabilities of GNNs. It achieves efficiency by explicitly identifying and eliminating redundancy and noise within GNN structures and embeddings. This strategy effectively mitigates misinformation in the graph while enhancing interpretability through pruned information.

## 6 Conclusion

This paper introduces a novel pruning framework, LightGNN, aimed at addressing scalability and robustness challenges in GNN-based collaborative filtering. LightGNN explicitly models the probabilities of redundancy and noise for each edge and embedding parameter within the GNN recommender, enabling precise pruning of misinformation. It is driven by innovative hierarchical distillation objectives that leverage high-order relations and multi-level distillation to enhance performance retention. Extensive experiments demonstrate that LightGNN outperforms baselines in recommendation performance, compression efficiency, and robustness.

# References

[1] R. v. d. Berg, T. N. Kipf, and M. Welling. Graph convolutional matrix completion. *arXiv preprint arXiv:1706.02263*, 2017.

[2] D. Chen, Y. Lin, W. Li, P. Li, J. Zhou, and X. Sun. Measuring and relieving the over-smoothing problem for graph neural networks from the topological view. In *AAAI Conference on Artificial Intelligence (AAAI)*, volume 34, pages 3438–3445, 2020.

[3] G. Chen, J. Chen, F. Feng, S. Zhou, and X. He. Unbiased knowledge distillation for recommendation. In *ACM International Conference on Web Search and Data Mining (WSDM)*, pages 976–984, 2023.

[4] L. Chen, L. Wu, R. Hong, K. Zhang, and M. Wang. Revisiting graph based collaborative filtering: A linear residual graph convolutional network approach. In *AAAI Conference on Artificial Intelligence (AAAI)*, volume 34, pages 27–34, 2020.

[5] T. Chen, Y. Sui, X. Chen, A. Zhang, and Z. Wang. A unified lottery ticket hypothesis for graph neural networks. In *International Conference on Machine Learning (ICML)*, pages 1695–1706. PMLR, 2021.

[6] J. Frankle and M. Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations (ICLR)*, 2018.

[7] C. Gao, Y. Zheng, N. Li, Y. Li, Y. Qin, J. Piao, Y. Quan, J. Chang, D. Jin, X. He, et al. A survey of graph neural networks for recommender systems: Challenges, methods, and directions. *ACM Transactions on Recommender Systems (TRS)*, 1(1):1–51, 2023.

[8] X. He, K. Deng, X. Wang, Y. Li, Y. Zhang, and M. Wang. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *International ACM SIGIR conference on research and development in Information Retrieval (SIGIR)*, pages 639–648, 2020.

[9] X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T.-S. Chua. Neural collaborative filtering. In *The ACM Web Conference (WWW)*, pages 173–182, 2017.

[10] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

[11] Z. Hu, Y. Dong, K. Wang, and Y. Sun. Heterogeneous graph transformer. In *The ACM Web Conference (WWW)*, pages 2704–2710, 2020.

[12] Y. Jiang, C. Huang, and L. Huang. Adaptive graph contrastive learning for recommendation. In *International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 4252–4261, 2023.

[13] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.

[14] Y. Koren, S. Rendle, and R. Bell. Advances in collaborative filtering. *Recommender systems handbook*, pages 91–142, 2021.

[15] Z. Lin, C. Tian, Y. Hou, and W. X. Zhao. Improving graph collaborative filtering with neighborhood-enriched contrastive learning. In *The ACM Web Conference (WWW)*, pages 2320–2329, 2022.

[16] A. v. d. Oord, Y. Li, and O. Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

[17] Y. Ren, Z. Haonan, L. Fu, X. Wang, and C. Zhou. Distillation-enhanced graph masked autoencoders for bundle recommendation. In *International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 1660–1669, 2023.

[18] S. Rendle, C. Freudenthaler, Z. Gantner, and L. Schmidt-Thieme. Bpr: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618*, 2012.

[19] S. Sedhain, A. K. Menon, S. Sanner, and L. Xie. Autorec: Autoencoders meet collaborative filtering. In *The ACM Web Conference (WWW)*, pages 111–112, 2015.

[20] X. Su and T. M. Khoshgoftaar. A survey of collaborative filtering techniques. *Advances in Artificial Intelligence*, 2009, 2009.

[21] J. Tang and K. Wang. Ranking distillation: Learning compact ranking models with high performance for recommender system. In *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, pages 2289–2298, 2018.

[22] C. Wang, Y. Yu, W. Ma, M. Zhang, C. Chen, Y. Liu, and S. Ma. Towards representation alignment and uniformity in collaborative filtering. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pages 1816–1825, 2022.

[23] W. Wang, F. Feng, X. He, L. Nie, and T.-S. Chua. Denoising implicit feedback for recommendation. In *ACM International Conference on Web Wearch and Data Mining (WSDM)*, pages 373–381, 2021.

[24] W. Wang, Y. Xu, F. Feng, X. Lin, X. He, and T.-S. Chua. Diffusion recommender model. *International ACM SIGIR conference on research and development in Information Retrieval (SIGIR)*, 2023.

[25] X. Wang, X. He, M. Wang, F. Feng, and T.-S. Chua. Neural graph collaborative filtering. In *International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 165–174, 2019.

[26] X. Wang, H. Jin, A. Zhang, X. He, T. Xu, and T.-S. Chua. Disentangled graph collaborative filtering. In *International ACM SIGIR conference on research and development in information retrieval (SIGIR)*, pages 1001–1010, 2020.

[27] J. Wu, X. Wang, F. Feng, X. He, L. Chen, J. Lian, and X. Xie. Self-supervised graph learning for recommendation. In *International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 726–735, 2021.

[28] L. Xia, C. Huang, C. Huang, K. Lin, T. Yu, and B. Kao. Automated self-supervised learning for recommendation. In *The ACM Web Conference (WWW)*, pages 992–1002, 2023.

[29] L. Xia, C. Huang, J. Shi, and Y. Xu. Graph-less collaborative filtering. In *The ACM Web Conference (WWW)*, pages 17–27, 2023.

[30] L. Xia, C. Huang, Y. Xu, J. Zhao, D. Yin, and J. Huang. Hypergraph contrastive collaborative filtering. In *International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 70–79, 2022.

[31] X. Xia, H. Yin, J. Yu, Q. Wang, G. Xu, and Q. V. H. Nguyen. On-device next-item recommendation with self-supervised knowledge distillation. In *International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 546–555, 2022.

[32] T. Yao, X. Yi, D. Z. Cheng, F. Yu, T. Chen, A. Menon, L. Hong, E. H. Chi, S. Tjoa, J. Kang, et al. Self-supervised learning for large-scale item recommendations. In *ACM International Conference on Information & Knowledge Management (CIKM)*, pages 4321–4330, 2021.

[33] R. Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, and J. Leskovec. Graph convolutional neural networks for web-scale recommender systems. In *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, pages 974–983, 2018.

[34] J. Yu, H. Yin, X. Xia, T. Chen, L. Cui, and Q. V. H. Nguyen. Are graph augmentations necessary? simple graph contrastive learning for recommendation. In *International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 1294–1303, 2022.

[35] A. Zhang, W. Ma, X. Wang, and T.-S. Chua. Incorporating bias-aware margins into contrastive loss for collaborative filtering. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:7866–7878, 2022.

[36] J. Zhang, X. Shi, S. Zhao, and I. King. Star-gcn: stacked and reconstructed graph convolutional networks for recommender systems. In *International Joint Conference on Artificial Intelligence (IJCAI)*, pages 4264–4270, 2019.

[37] S. Zhang, Y. Liu, Y. Sun, and N. Shah. Graph-less neural networks: Teaching old mlps new tricks via distillation. In *International Conference on Learning Representations (ICLR)*, 2021.

[38] S. Zhang, L. Yao, A. Sun, and Y. Tay. Deep learning based recommender system: A survey and new perspectives. *ACM computing surveys (CSUR)*, 52(1):1–38, 2019.

## A  Ethical Considerations

### A.1  Ethical Implications

The proposed research on LightGNN, a distillation-based GNN pruning framework, introduces innovative techniques for model compression in Graph Neural Networks (GNNs) to reduce model complexity while preserving recommendation accuracy. While the advancements in this field are promising, there are ethical implications that need to be considered since graph-based recommendation systems often rely on sensitive user interaction data.

- **Privacy Considerations**. GNN's utilization of user interaction data raises concerns about privacy. The pruning process must safeguard against unauthorized access to sensitive user information contained within the graph data. Besides, the removal of edges and embedding entries during compression should be conducted in a manner that does not inadvertently expose or retain identifiable user information.

- **Security and Safety**. Pruning components based on learnable algorithms may introduce vulnerabilities that could compromise the integrity of the recommendation system, potentially leading to data breaches or manipulation. Moreover, aggressive pruning to achieve high compression rates might compromise the robustness of the GNN model, making it more susceptible to adversarial attacks or unexpected behaviors.

### A.2  Mitigation Strategies

Below, we introduce some possible mitigation strategies.

- **Privacy-Preserving Techniques**. Implement encryption and anonymization methods to protect user data while ensuring that the pruning process does not compromise individual privacy.

- **Security Audits**. Conduct thorough security assessments to identify and address potential vulnerabilities introduced by the pruning framework, ensuring data integrity and system security.

- **Transparency and Accountability**. Maintain transparency in the pruning process, providing clear explanations of how components are pruned and enabling users to understand and challenge the recommendations made by the system.

In conclusion, while LightGNN shows promise in reducing model complexity while retaining recommendation performance, it is important for researchers and developers to prioritize ethical considerations to mitigate potential negative societal impacts and uphold the integrity and fairness of AI systems in recommendation.