

Long-tail Augmented Graph Contrastive Learning for Recommendation

Qian Zhao, Zhengwei Wu, Zhiqiang Zhang, and Jun Zhou[✉]

Ant Group, Hangzhou, China
{zq317110, zejun.wzw, lingyao.zzq, jun.zhoujun}@antgroup.com

Abstract. Graph Convolutional Networks (GCNs) has demonstrated promising results for recommender systems, as they can effectively leverage high-order relationship. However, these methods usually encounter data sparsity issue in real-world scenarios. To address this issue, GCN-based recommendation methods employ contrastive learning to introduce self-supervised signals. Despite their effectiveness, these methods lack consideration of the significant degree disparity between head and tail nodes. This can lead to non-uniform representation distribution, which is a crucial factor for the performance of contrastive learning methods. To tackle the above issue, we propose a novel **Long-tail Augmented Graph Contrastive Learning (LAGCL)** method for recommendation. Specifically, we introduce a learnable long-tail augmentation approach to enhance tail nodes by supplementing predicted neighbor information, and generate contrastive views based on the resulting augmented graph. To make the data augmentation schema learnable, we design an auto drop module to generate pseudo-tail nodes from head nodes and a knowledge transfer module to reconstruct the head nodes from pseudo-tail nodes. Additionally, we employ generative adversarial networks to ensure that the distribution of the generated tail/head nodes matches that of the original tail/head nodes. Extensive experiments conducted on three benchmark datasets demonstrate the significant improvement in performance of our model over the state-of-the-arts. Further analyses demonstrate the uniformity of learned representations and the superiority of LAGCL on long-tail performance.

Keywords: Recommender system · Graph neural networks · Contrastive learning · Self-supervised learning.

1 Introduction

Recommender systems are a critical component of numerous online services, ranging from e-commerce to online advertising. As a classic approach, collaborative filtering (CF) [11,19] plays a vital role in personalized preference prediction by representing user and item embeddings from observed user-item interactions such as clicks and conversions. Recently, enhanced by the powerful Graph Convolutional Networks (GCNs) [13], GCN-based recommendation methods [10,22] have demonstrated significant potential in improving recommendation accuracy.

GCN-based recommendation methods represent interaction data as graphs, such as the user-item interaction graph, and iteratively propagate neighborhood information to learn effective node representations. Compared to traditional CF methods, GCN-based recommendation methods are better equipped to capture higher-order collaborative signals, leading to improved user and item embedding learning.

Despite the effectiveness, GCN-based recommendation methods still face data sparsity issue in real-world scenarios. Most existing models follow the supervised learning paradigm [1,10,22,27], where the supervision signal is derived from the observed user-item interactions. However, the observed interactions are considerably sparse in contrast to the entire interaction space [24]. As a result, they may not be sufficient to acquire effective representations. Although several recent studies have attempted to alleviate the data sparsity of interaction data through contrastive learning [24,28], they generally rely on pre-defined data augmentation strategies, such as uniformly dropping edges or shuffling features. These methods lack consideration of the significant disparity in the graph structure between head and tail nodes and lack the ability to construct adaptive data augmentation tailored for various recommendation datasets. This can lead to non-uniform representation distribution, which is a crucial factor for the performance of contrastive learning methods [21,28].

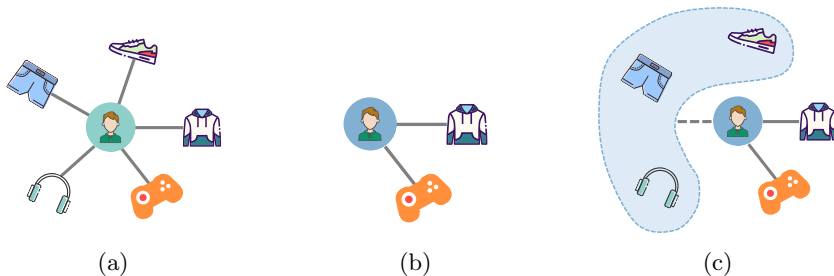


Fig. 1: An illustrated example depicting (a) a head node, (b) a tail node, and (c) a tail node augmented with predicted neighbor information.

In the light of the above limitations and challenges, we propose a novel **Long-tail Augmented Graph Contrastive Learning (LAGCL)** method for recommendation. To illustrate, consider the head user in Fig. 1(a) and the tail user in Fig. 1(b) who share similar preference. Our approach aims to extract informative transition patterns from head users and adapt to tail users effectively, as shown in Fig. 1(c). Specifically, we first design an auto drop module to convert head nodes into pseudo-tail nodes that mimic the patterns of real-tail nodes. Next, we leverage the dropped neighbor information of pseudo-tail nodes to learn the knowledge transfer module, which then augments the tail nodes as pseudo-head nodes by adding predicted neighbor information. These modules are updated using the adversarial training mechanism to ensure the distribution match between

pseudo-head/tail nodes and real-head/tail nodes. Finally, we use an effective and efficient approach to generate contrastive views by injecting random noise into the graph embedding of both head nodes and augmented tail nodes, capitalizing on their uniformity to yield better contrastive performance.

The main contributions of this paper are summarized as follows:

- We propose a novel graph contrastive learning method, which encourages the model to learn the knowledge between head and tail nodes and generate uniform representation distribution for improving the GCN-based recommendation.
- We designed a learnable data augmentation scheme that can adaptively enhance tail nodes representation and easily generalize to different GCN-based recommendation scenarios.
- Extensive experiments are conducted on three public datasets, demonstrating our approach is consistently better than a number of competitive baselines, including GCN-based and graph contrastive learning-based recommendation methods.

2 Preliminaries

2.1 GCN-based Recommendation

Given the user set $\mathcal{U} = \{u\}$ and the item set $\mathcal{I} = \{i\}$, the observed user-item interaction data is denoted as $\mathbf{R} \in \mathbb{R}^{|\mathcal{U}| \times |\mathcal{I}|}$, where each entry $r_{u,i} = 1$ if there exists an interaction between user u and item i , otherwise $r_{u,i} = 0$. The number of nodes is $n = |\mathcal{U}| + |\mathcal{I}|$. GCN-based recommendation methods formulate the available data as a user-item bipartite graph $\mathcal{G} = (\mathcal{V}, \mathbf{A})$, where $\mathcal{V} = \mathcal{U} \cup \mathcal{I}$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$ is the adjacent matrix defined as

$$\mathbf{A} = \begin{bmatrix} \mathbf{0}^{|\mathcal{U}| \times |\mathcal{U}|} & \mathbf{R} \\ \mathbf{R}^T & \mathbf{0}^{|\mathcal{I}| \times |\mathcal{I}|} \end{bmatrix}. \quad (1)$$

With a slight abuse of notation, we use $|\mathbf{A}_i|$ to refer to $\sum_{j \in \mathcal{N}_i} \mathbf{A}_{ij}$, where \mathcal{N}_i denotes the neighbor set of node i . GCN-based recommendation methods utilize graph structure information to aggregate and produce the embedding of users and items on bipartite graph \mathcal{G} through Eq. (2).

$$\mathbf{H}^{(l)} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(l-1)}, \quad (2)$$

where $\mathbf{D} \in \mathbb{R}^{n \times n}$ is the diagonal degree matrix of \mathcal{G} , in which each entry \mathbf{D}_{ii} denotes the number of non-zeros in the i -th row of the matrix \mathbf{A} . $\mathbf{H}^{(l)} \in \mathbb{R}^{n \times d}$ denotes the d -dimensional node embedding matrix after l -th graph convolution layer, and $\mathbf{H}^{(0)}$ is the initial node embedding matrix that need to be learned. Finally, we combine all the L layers output node embeddings, *i.e.*, $\mathbf{H} = f_{\text{readout}}([\mathbf{H}^{(0)}; \mathbf{H}^{(1)}; \dots; \mathbf{H}^{(L)}])$, to generate preference scores between users and items for recommendation, while $f_{\text{readout}}(\cdot)$ is the mean pooling operation here.

2.2 Long-Tail Distribution in the Graph

Graph convolutional networks heavily rely on rich structural information to achieve high performance. However, for nodes with low degrees, their number of neighbors is typically very small, leading to unsatisfactory performance for these nodes. In this paper, in order to investigate the long-tail distribution in the graph, we partition nodes into head nodes \mathcal{V}_{head} and tail nodes \mathcal{V}_{tail} based on their degree with a predetermined threshold value k , i.e., $\mathcal{V}_{head} = \{i : \mathbf{D}_{ii} > k\}$ and $\mathcal{V}_{tail} = \{i : \mathbf{D}_{ii} \leq k\}$. Specifically, we have $\mathcal{V}_{head} \cap \mathcal{V}_{tail} = \emptyset$ and $\mathcal{V}_{head} \cup \mathcal{V}_{tail} = \mathcal{V}$.

3 Methodology

In this section, we introduce the proposed long-tail augmented graph contrastive learning (LAGCL) model for recommendation. First, we briefly exhibit the overall architecture of LAGCL. Then, we elaborate the detail implementation of LAGCL.

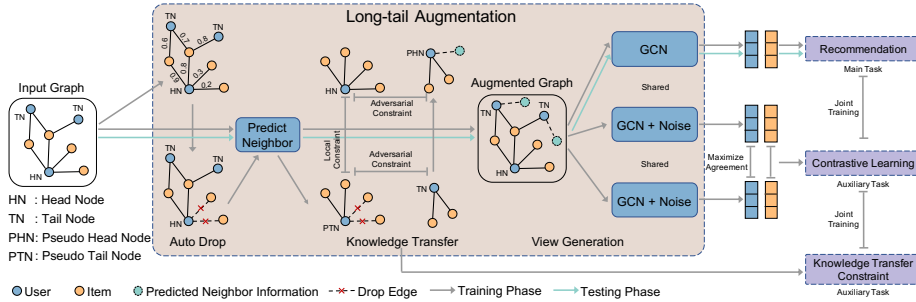


Fig. 2: Overall framework of our proposed long-tail augmented graph contrastive learning method for recommendation. To enhance readability, only user-side augmentation is shown. The core of the model is long-tail augmentation, which extracts informative transition from head nodes and augments tail nodes using the knowledge transfer module. The augmented graph is then used to generate the view pairs for contrastive learning. Finally, we train recommendation task, contrastive learning task and knowledge transfer constraints jointly using the multi-task training strategy.

3.1 Overview of LAGCL

Our framework follows the general contrastive learning paradigm, which aims to achieve maximum agreement between representations of different views. Fig. 2 illustrates the overall framework. Firstly, we enhance tail node by a knowledge transfer module and then generate two graph views on the augmented graph.

Next, we employ a contrastive loss function to conduct the contrastive learning task, which encourages representations of the same nodes in the two different views to be similar, while representations of different nodes in those views to be distinct. Finally, we adopt a multi-task training strategy to improve recommendation performance. Specifically, the main task is the recommendation task, while contrastive learning task and knowledge transfer constraints serve as auxiliary tasks.

3.2 Long-tail Augmentation

Assuming that each node has a ground truth neighbor set, the difference between head and tail nodes lies in the former having a fully observed neighbor set while the latter has only partially observed neighbor information. This incomplete information adversely affects the performance of the model. Therefore, the technique of long-tail augmentation aims to supplement the tail nodes with missing neighbor information, with the aim of enforcing representation uniformity between head and tail nodes. However, conducting long-tail augmentation is a challenging task due to unavailability of ground truth labels for the missing neighbors of tail nodes. The intuition behind this approach is that head nodes have fully observed neighbor information. If we could generate pseudo-tail nodes by dropping the neighbor information of head nodes to mimic tail nodes, we could leverage the dropped neighbor information of pseudo-tail nodes to learn the long-tail augmentation strategy. Our long-tail augmentation boils down to three key modules: auto drop module, knowledge transfer module and generative adversarial learning module.

Auto Drop Module. An intuitive approach is to randomly sample neighbors from head nodes to generate pseudo-tail nodes. However, this method lacks an effective supervisory signal to guide the sampling process, leading to a distribution deviation between the enhanced data samples and the actual tail node samples. To improve the data augmentation process, we propose an auto drop module equipped with a trainable dropout strategy to minimize the distribution deviation between the enhanced data samples and the actual tail node samples. Specifically, the dropout strategy consists of two main steps: sorting by edge weight and selecting top- K neighbors based on their importance weight, which is calculated between node pairs as:

$$\mathbf{S} = (\mathbf{H}^{(0)} \mathbf{W}_s \mathbf{H}^{(0)T}) \odot \mathbf{A}, \quad (3)$$

where $\mathbf{S} \in \mathbb{R}^{n \times n}$, $\mathbf{W}_s \in \mathbb{R}^{d \times d}$ is trainable parameters and \odot denotes the element-wise product. Then, to mimic real-tail nodes, we randomly choose the neighbor size k_i of head nodes uniformly from the range of tail node neighbor sizes $[1, k]$:

$$k_i = \begin{cases} \text{Uniform}[1, k], & \mathbf{D}_{ii} > k, \\ \mathbf{D}_{ii}, & \mathbf{D}_{ii} \leq k. \end{cases} \quad (4)$$

Finally, inspired by the top-rank selection method in [5,15], the new adjacency matrix $\hat{\mathbf{A}}$ is constructed by selecting the top- K neighbors based on their importance weight, which will be used in the graph aggregation of tail nodes. The new adjacency matrix $\hat{\mathbf{A}}$ is defined as:

$$\hat{\mathbf{A}}_{ij} = \begin{cases} \frac{\exp(\delta \mathbf{S}_{ij})}{1 + \exp(\delta \mathbf{S}_{ij})}, & \mathbf{S}_{ij} \in \text{top-}k_i(\mathbf{S}_i), \\ 0, & \mathbf{S}_{ij} \notin \text{top-}k_i(\mathbf{S}_i), \end{cases} \quad (5)$$

where δ is a hyperparameter that controls the smoothness of the adjacency matrix $\hat{\mathbf{A}}$, and the corresponding neighborhood of node i is denoted as $\hat{\mathcal{N}}_i$.

Knowledge Transfer Module. After constructing the auto drop module, we can design a knowledge transfer module that leverages the dropped neighbor information of pseudo-tail nodes for predicting the missing neighbor information of real-tail nodes. Specifically, we use a multi-layer perceptron (MLP) function, denoted by $f_t(\mathbf{h}_i^{(l)}, \mathbf{h}_{\mathcal{N}_i}^{(l)}; \theta_t^{(l)}) = \mathbf{m}_i^{(l)}$, to predict the missing neighbor information based on the center node features and observed neighbor information. Here, $\mathbf{h}_i^{(l)}$ represents the l -layer graph embedding of node i , $\mathbf{h}_{\mathcal{N}_i}^{(l)}$ is the mean pooling representation of observable neighbors. Then, the predicted information is added to the neighbor aggregation process of real-tail node.

We can calculate the node embedding of head node i in the sparse graph $\hat{\mathbf{A}}$ with predicted neighbor information by the knowledge transfer function:

$$\hat{\mathbf{h}}_i^{(l)} = \sum_{j \in \hat{\mathcal{N}}_i} \frac{1}{\sqrt{|\hat{\mathbf{A}}_i|} \sqrt{|\hat{\mathbf{A}}_j|}} \hat{\mathbf{h}}_j^{(l-1)} + f_t(\hat{\mathbf{h}}_i^{(l-1)}, \hat{\mathbf{h}}_{\hat{\mathcal{N}}_i}^{(l-1)}, \theta_t^{(l-1)}), \quad (6)$$

where $\hat{\mathbf{h}}^{(0)} = \mathbf{h}^{(0)}$. The embedding of head node i in the original graph \mathbf{A} is

$$\mathbf{h}_i^{(l)} = \sum_{j \in \mathcal{N}_i} \frac{1}{\sqrt{|\mathbf{A}_i|} \sqrt{|\mathbf{A}_j|}} \mathbf{h}_j^{(l-1)}. \quad (7)$$

In order to train the knowledge transfer function, we define the translation loss is defined as follows:

$$\mathcal{L}_{trans} = \sum_{i \in \mathcal{V}_{head}} \sum_{l=1}^L \|\mathbf{h}_i^{(l)} - \hat{\mathbf{h}}_i^{(l)}\|_2^2. \quad (8)$$

Generative Adversarial Learning Module. To learn effective long-tail augmentation, the representation distribution of real-tail nodes should match the representation distribution of pseudo-tail nodes generated by the auto drop module, which is calculated as follows:

$$\tilde{\mathbf{h}}_i^{(l)} = \sum_{j \in \hat{\mathcal{N}}_i} \frac{1}{\sqrt{|\hat{\mathbf{A}}_i|} \sqrt{|\hat{\mathbf{A}}_j|}} \tilde{\mathbf{h}}_j^{(l-1)}, \quad (9)$$

where $\tilde{\mathbf{h}}^{(0)} = \mathbf{h}^{(0)}$. Additionally, the distribution of pseudo-head nodes augmented by the knowledge transfer module should match that of real-head nodes. To achieve this, we use Generative Adversarial Networks [7]. The discriminator distinguishes pseudo-head/tail nodes from real-head/tail nodes based on the node representations, while the generator aims to provide information that is consistently classified as real nodes by the discriminator. Here we regard the output layer of LAGCL as the generator, which contests with the discriminator in the learning process. In particular, we use the following loss for the tail nodes adversarial constraint:

$$\begin{aligned} \mathcal{L}_{tail-disc} = & \sum_{i \in \mathcal{V}} \mathbb{1}(i \notin \mathcal{V}_{tail}) \text{CROSSENT}(\mathbf{0}, f_d(\tilde{\mathbf{h}}_i; \theta_d)) \\ & + \mathbb{1}(i \in \mathcal{V}_{tail}) \text{CROSSENT}(\mathbf{1}, f_d(\mathbf{h}_i; \theta_d)), \end{aligned} \quad (10)$$

and use the following loss for the head nodes adversarial constraint:

$$\begin{aligned} \mathcal{L}_{head-disc} = & \sum_{i \in \mathcal{V}} \text{CROSSENT}(\mathbf{0}, f_d(\hat{\mathbf{h}}_i; \theta_d)) \\ & + \mathbb{1}(i \in \mathcal{V}_{head}) \text{CROSSENT}(\mathbf{1}, f_d(\mathbf{h}_i; \theta_d)), \end{aligned} \quad (11)$$

where $\text{CROSSENT}(\cdot)$ is the cross entropy function, $\mathbb{1}(\cdot)$ is the indicator function, $f_d(\cdot; \theta_d)$ is the discriminator function parameterized by θ_d , which calculates the probability of a node being a head node, as

$$f_d(\mathbf{h}_i; \theta_d) = \sigma\left(\mathbf{w}_d^\top \text{LEAKYRELU}(\mathbf{W}_d \mathbf{h}_i + \mathbf{b}_d)\right), \quad (12)$$

where $\text{LEAKYRELU}(\cdot)$ is used as the activation function, $\sigma(\cdot)$ is the sigmoid function, $\theta_d = \{\mathbf{W}_d \in \mathbb{R}^{d \times d}, \mathbf{b}_d \in \mathbb{R}^{d \times 1}, \mathbf{w}_d \in \mathbb{R}^{d \times 1}\}$ contains the learnable parameters of the discriminator f_d .

3.3 Contrastive Learning

View Generation. With the knowledge transfer function mentioned above, we can obtain the augmented tail node embedding, as follows:

$$\mathbf{h}_i^{(l)} = \sum_{j \in \mathcal{N}_i} \frac{1}{\sqrt{|\mathbf{A}_i|} \sqrt{|\mathbf{A}_j|}} \mathbf{h}_j^{(l-1)} + f_t(\mathbf{h}_i^{(l-1)}, \mathbf{h}_{\mathcal{N}_i}^{(l-1)}; \theta_t^{(l-1)}). \quad (13)$$

Then, we follow the approach of SimGCL [28] and generate different views by slightly rotating refined node embeddings in space. This method retains original information while introducing the InfoNCE loss as an additional self-supervised task to improve robustness, as follows:

$$\mathbf{h}_i^{(l)'} = \mathbf{h}_i^{(l)} + \Delta_i^{(l)'}, \mathbf{h}_i^{(l)''} = \mathbf{h}_i^{(l)} + \Delta_i^{(l)''}, \quad (14)$$

where the noise vectors $\Delta_i^{(l)'}$ and $\Delta_i^{(l)''}$ are subject to $\|\Delta\|_2 = \epsilon$ and $\Delta = \bar{\Delta} \odot \text{sign}(\mathbf{h}_i^{(l)})$, $\bar{\Delta} \in \mathbb{R}^d \sim U(0, 1)$. We can use ϵ to control the rotation angle of

$\mathbf{h}_i^{(l)'}$, $\mathbf{h}_i^{(l)''}$ compared to $\mathbf{h}_i^{(l)}$. Since $\mathbf{h}_i^{(l)}$, $\mathbf{h}_i^{(l)'}$, $\mathbf{h}_i^{(l)''}$ always belong to the same hyperoctant, so adding noise will not cause significant deviation. The noise is injected into each convolution layer and we average each layer output as the final node embedding. To simplify the notation, we denote \mathbf{h}_i as the final node embedding after L layers, \mathbf{h}'_i and \mathbf{h}''_i as two generated views.

Contrastive Loss. After obtaining the generated views of each node, we utilize the contrastive loss, InfoNCE [8], to maximize the agreement of positive pairs and minimize that of negative pairs:

$$\mathcal{L}_{cl}^U = \sum_{u \in \mathcal{U}} -\log \frac{\exp(s(\mathbf{h}'_u, \mathbf{h}''_u)/\tau)}{\sum_{v \in \mathcal{U}} \exp(s(\mathbf{h}'_u, \mathbf{h}''_v)/\tau)}, \quad (15)$$

where $s(\cdot)$ measures the similarity between two vectors, which is set as cosine similarity function; τ is the hyper-parameter, known as the temperature in softmax. Analogously, we obtain the contrastive loss of the item side \mathcal{L}_{cl}^I . Combining these two losses, we get the objective function of self-supervised task as $\mathcal{L}_{cl} = \mathcal{L}_{cl}^U + \mathcal{L}_{cl}^I$.

3.4 Multi-task Training

We leverage a multi-task training strategy to optimize the main recommendation task and the auxiliary tasks including translation task, discrimination task and contrastive learning task jointly:

$$\mathcal{L} = \mathcal{L}_{rec} + \lambda_1 \mathcal{L}_{trans} + \lambda_2 \mathcal{L}_{disc} + \lambda_3 \mathcal{L}_{cl} + \lambda_4 \|\Theta\|^2, \quad (16)$$

where Θ is the set of model parameters, $\lambda_1, \lambda_2, \lambda_3, \lambda_4$ are hyperparameters to control the strengths of the diversity preserving loss. \mathcal{L}_{rec} is the loss function of the main recommendation task. In this work, we adopt Bayesian Personalized Ranking (BPR) loss [18]:

$$\mathcal{L}_{rec} = \sum_{u, i, j \in \mathcal{O}} -\log \sigma(\hat{y}_{u, i} - \hat{y}_{u, j}), \quad (17)$$

where $\hat{y}_{u, i} = \mathbf{h}_u^\top \mathbf{h}_i$ is the preference score. $\sigma(\cdot)$ denotes the sigmoid function. $\mathcal{O} = \{(u, i, j) | (u, i) \in \mathcal{O}^+, (u, j) \in \mathcal{O}^-\}$ denotes the training data, and \mathcal{O}^- is the unobserved interactions.

4 Experiments

To verify the effectiveness of the proposed LAGCL, we conduct extensive experiments and report detailed analysis results.

4.1 Experimental Setups

Datasets. We evaluate the LAGCL using three widely-used public benchmark datasets: Yelp2018¹, Amazon-Book² and Movielens-25M³. The statistics of these datasets are presented in Table 1.

Evaluation Metrics. Due to our focus on Top-N recommendation, following the convention in the previous research[29], we discard ratings less than 4 in Movielens-25M, and reset the rest to 1. We split the interactions into training, validation, and testing set with a ratio of 7:1:2. In the test set, we evaluate the performance of each model using the relevancy-based metric Recall@20 and the ranking-aware metric NDCG@20.

Table 1: The statistics of three datasets.

Dataset	#Users	#Items	#Interactions	Density
Yelp2018	31,668	38,048	1,561,406	0.1296%
Amazon-Book	52,643	91,599	2,984,108	0.0619%
Movielens-25M	155,002	27,133	3,612,474	0.0859%

Baselines. We compare LAGCL with other GCN-based recommendation methods, including:

- **LightGCN**[10] designs a light graph convolution to improve training efficiency and representation ability.
- **SGL**[24] designs an auxiliary tasks via perturbation to the graph structure (such as edge dropout), which achieves greater benefits in long-tail recommendation.
- **NCL**[16] improves the recommendation performance by clustering similar nodes to provide semantic neighbors and structural neighbors.
- **RGCF**[20] integrates a graph denoising module and a diversity preserving module to enhance the robustness of GCN-based recommendation.
- **SimGCL**[28] proves the positive correlation between the uniformity of representations and the ability to debias through feature-level perturbation and contrastive learning, and achieved greater long-tail performance than SGL.

Settings and Hyperparameters. We develop the model using the open-source SELFRec⁴ [29]. For a fair comparison, we prioritize the hyperparameters reported in the original papers for each baseline when feasible. In cases

¹ <https://www.yelp.com/dataset>

² https://cseweb.ucsd.edu/~jmcauley/datasets.html#amazon_reviews

³ <https://grouplens.org/datasets/movielens/25m/>

⁴ <https://github.com/Coder-Yu/SELFRec>

where this information is not provided, we conduct a grid search to adjust the hyperparameters. As for the general settings of all the baselines, the Xavier initialization[6] is used on all embeddings. The embedding size is 64, the parameter for L_2 regularization is 10^{-4} and the batch size is 2048. We use Adam[12] with the learning rate 0.001 to optimize all the models. More settings can be found in <https://github.com/im0qianqian/LAGCL>.

4.2 Performance Comparison

Table 2 shows our comparison results with other baselines in three datasets. We have the following observations: (1) Our proposed LAGCL consistently outperforms all baselines in different datasets and metrics. Specifically, LAGCL achieves a relative improvement of 25.6%, 61.8%, and 9.6% on Recall@20 compared to LightGCN on the Yelp2018, Amazon Book, and Movielens-25M datasets, respectively. Compared to the strongest baseline (SimGCL), LAGCL also achieves better performance, e.g, about 1.81%, 2.35%, 0.73% performance improvement of Recall@20 on the same datasets. (2) All graph contrastive learning based methods *e.g.*, SGL, NCL, RGCF, SimGCL, show significant improvement compared to LightGCN on three datasets, which verifies the effectiveness of contrastive learning for collaborative filtering. SimGCL achieves better performance than other baselines, demonstrating that feature-level augmentation is more suitable for collaborative filtering tasks than structure-level augmentation. It is worth noting that our method predicts neighborhood information for tail nodes while preserving the original graph structure. LAGCL incorporates the advantages of feature-level augmentation and avoids the possibility of drastic changes in tail node information due to graph structural changes caused by structure-level augmentation. As a result, our method LAGCL achieves the state-of-the-art performance.

Table 2: Overall performance comparison. The percentage in brackets denote the relative performance improvement over LightGCN. The best results are bolded and the best results of baseline are underlined.

Method	Yelp2018		Amazon Book		MovieLens-25M	
	Recall@20	NDCG@20	Recall@20	NDCG@20	Recall@20	NDCG@20
LightGCN	0.0583	0.0486	0.0323	0.0254	0.3267	0.2276
SGL	0.0659(+13.0%)	0.0541(+11.4%)	0.0443(+37.0%)	0.0352(+38.5%)	0.3471(+6.2%)	0.2440(+7.2%)
NCL	0.0663(+13.7%)	0.0547(+12.5%)	0.0426(+32.0%)	0.0331(+30.2%)	0.3292(+0.8%)	0.2306(+1.3%)
RGCF	0.0591(+1.5%)	0.0487(+0.1%)	0.0345(+6.9%)	0.0274(+7.9%)	0.3137(-4.0%)	0.2060(-9.5%)
SimGCL	0.0719(+23.4%)	0.0600(+23.4%)	0.0510(+57.9%)	0.0406(+59.8%)	0.3553(+8.8%)	0.2468(+8.4%)
LAGCL	0.0732(+25.6%)	0.0604(+24.3%)	0.0522(+61.8%)	0.0415(+63.4%)	0.3579(+9.6%)	0.2509(+10.2%)

4.3 Ablation study

We conduct an ablation study to compare several ablated variants, including the “w/o AD” variant that uses random dropout instead of the auto drop module. We also consider the “w/o KT” variant that do not use knowledge transfer

module to augment the tail nodes, and the “w/o GAN” variants, which are generative adversarial networks proposed in Section 3.2. We have the following observations from Fig. 3. (1) Removing any of the components leads to a performance decrease, with the knowledge transfer module contributing the most. This demonstrates the importance of using knowledge transfer module to augment tail nodes is crucial in GCN-based recommendation scenarios. (2) Using random dropout instead of the auto drop module results in a decrease in performance. It indicates that auto drop module can help to extract informative transition patterns from head nodes that the random dropout strategy cannot learn. (3) Removing the generative adversarial networks significantly decreases the performance. This demonstrates that we cannot achieve meaningful data augmentation without the generative adversarial networks to ensure the distribution match between pseudo-head/tail nodes and real-head/tail nodes.

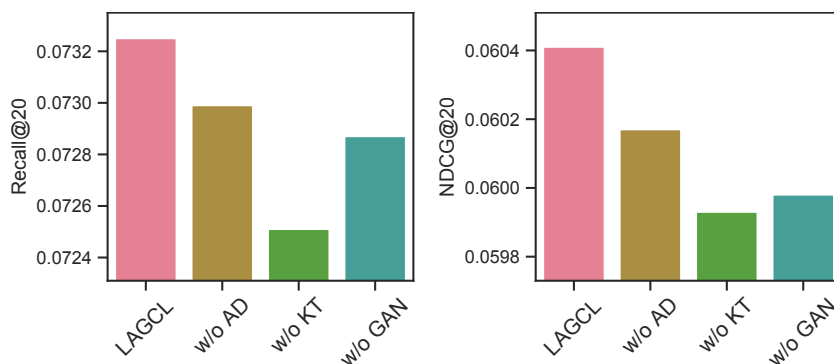
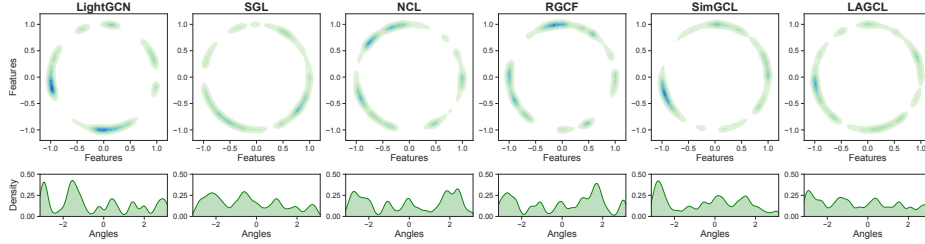


Fig. 3: Ablation study in the Yelp2018.

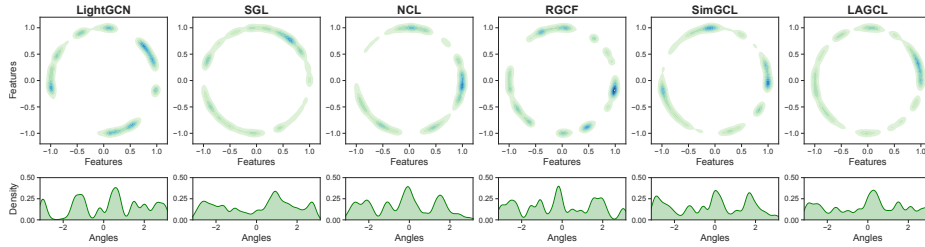
4.4 Further Analysis of LAGCL

Distribution uniformity analysis. A key contribution of the proposed LAGCL is the utilization of long-tail augmentation to supplement the neighbor information of tail nodes for GCN-based recommendation. To better understand the benefits brought by LAGCL, we visualize the learned embeddings in Fig 4 to illustrate how the proposed approach affects representation learning. We use Gaussian kernel density estimation (KDE) [2] to plot user and item embedding distributions in two-dimensional space. Additionally, we visualize each node’s density estimations on angles on the unit hypersphere \mathcal{S}^1 (i.e., circle with radius 1). We can see that, the embeddings learned by LightGCN fall into several clusters located on narrow arcs. Graph contrastive learning-based methods exhibit a more uniform distribution than LightGCN, where SimGCL has a more uniform distribution than other structure-based graph contrastive learning-based methods (SGL, NCL and RGCF). When compared to the best baseline SimGCL, the

LAGCL distribution has fewer dark areas on the circle, indicating that LAGCL has a more uniform distribution that benefits the model performance. Previous studies [21,28] have shown a strong correlation between contrastive learning and the uniformity of learned representations. Thus, we speculate that a more uniform distribution of embeddings could endow the model with better capacity to capture diverse user preferences and item characteristics.



(a) Distribution of user representations learned from Yelp2018 dataset.



(b) Distribution of item representations learned from Yelp2018 dataset.

Fig. 4: The distribution of user and item representations learned from the Yelp2018 dataset. The top of each figure plots the Gaussian Kernel Density Estimation (KDE) in \mathbb{R}^2 (the darker the color is, the more points fall in that area). The bottom of each figure plots the KDE on angles (i.e., $\text{atan2}(y, x)$) for each point $(x, y) \in \mathcal{S}^1$

Long Tail Degree Analysis. To verify whether LAGCL can provide additional performance gains for tail nodes, we divide each user into 10 groups of equally size based on their node degree in the user-item bipartite graph, as shown in Fig. 5. The smaller the Group Id, the lower node degree, and the lower user activity. We compare LAGCL with other baselines that alleviate long-tail distribution in graph, such as SGL employs structure-level augmentation and SimGCL utilizes feature-level augmentation. The results show that the performance of these methods are very similar for high-activity users, while for low-activity users, LAGCL exhibits better performance. This proves that our method has significant gains for modeling tail users.

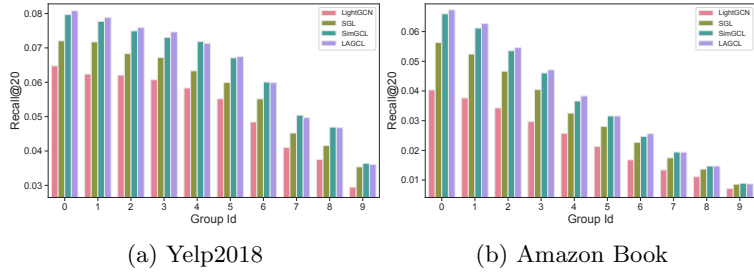


Fig. 5: Performance comparison of different user groups. In addition to the classic method LightGCN, we also select two different approaches to alleviate the long-tail distribution in the graph, namely SGL and SimGCL.

Degree threshold k . As shown in Fig. 6, we conduct a parameter sensitivity experiment on the threshold for dividing the head and tail nodes using the Yelp2018 dataset. The results show that our model achieve the best performance at $k = 20$. We conclude that: (1) LAGCL can learn a transfer strategy through the head users to help bring benefits to tail users. (2) A larger k will result in fewer number of the head users, which may cause the knowledge transfer module to be insufficiently learned. Conversely, a smaller k will result in significant quality differences within the head users, causing inaccurate learning in the knowledge transfer module. Therefore, choosing a suitable threshold for division can maximize the overall benefit.

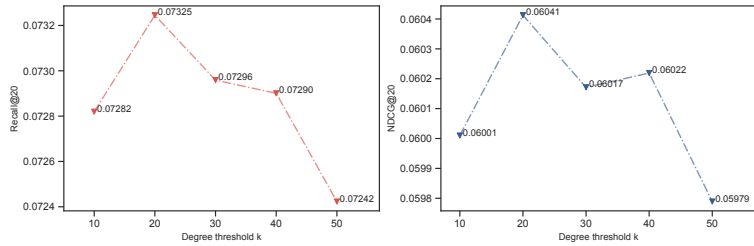


Fig. 6: Performance of LAGCL in the Yelp2018 when adjusting the degree threshold k .

5 Related Work

5.1 GCN-based Recommendation

GCN-based recommendation methods have treated interaction data as a bipartite graph and applied graph neural networks to obtain embeddings while cap-

turing high-order information. For example, GC-MC [1] applies the Graph Convolution Network (GCN) on the user-item graph and employs one convolutional layer to exploit the direct connections between users and items. PinSage [27] is an industrial solution that leverages efficient random walks and graph convolutions to generate item embeddings in Pinterest. NGCF [22] models high-order connectivities by effectively injecting the collaborative signal into the embedding process explicitly. LightGCN [10] simplifies the design of GCN by removing feature transformation and nonlinear activation, making it more effective for recommendation purpose.

However, all of these works perform training in the supervised learning paradigm, which heavily relies on labeled data. Since the observed interactions are considerably sparse compared to the entire interaction space, they may not be sufficient to acquire effective representations.

5.2 Self-supervised Learning in Recommender Systems

Self-supervised learning is an emerging technique that leverages unlabeled data to achieve impressive success in various fields, including computer vision [3,9], natural language processing [4,14] and graph learning [25,17,26]. Inspired by the success of these works, self-supervised learning has also been applied to recommender systems, where it has shown great improvement [29]. For instance, S^3 -Rec utilizes four self-supervised objectives to learn the correlations among attribute, item, subsequence, and sequence by maximizing the mutual information [30]. Similarly, CLCRec further utilizes contrastive learning to mitigate the cold-start issue by maximizing the mutual dependencies between item content and collaborative signals [23]. Moreover, a promising line of research has incorporated contrastive learning into graph-based recommenders to tackle the label sparsity issue with self-supervision signals. For example, SGL [24] uses node dropout, edge dropout and random walk on the user-item interaction graph to generate different views and maximize the agreement between them. NCL [16] incorporates structural and semantic neighbors to enhance graph-based collaborative filtering. SimGCL [28] generates contrastive views by adding uniform noise in the embedding space instead of graph augmentations.

6 Conclusion

In this work, we propose a novel long-tail augmented graph contrastive learning (LAGCL) method for recommendation. Specifically, we introduce a learnable long-tail augmentation schema that enhances tail nodes by supplementing them with predicted neighbor information. To make the data augmentation schema learnable, we design an auto-drop strategy to generate pseudo-tail nodes from head nodes and a knowledge translation module to reconstruct the head nodes from pseudo-tail nodes. To ensure the effectiveness of the data augmentation, we leverage adversarial constraints to distinguish between pseudo-tail and real-tail nodes, as well as between augmented tail nodes and real-head nodes. Comprehensive experiments demonstrate the effectiveness of our proposed LAGCL.

Ethics Statement

Our work aims to improve the accuracy of recommendations using only exposure and click data, without collecting or processing any personal information. We recognize the potential ethical implications of recommendation systems and their impact on user privacy, and we have taken steps to ensure that our work adheres to ethical standards.

We understand that recommendation systems have the potential to influence user behavior and may raise concerns related to user privacy. To address these concerns, we have designed our approach to rely solely on non-personal data, ensuring that our methods do not infringe on user privacy or rights.

In summary, our work focuses on improving recommendation accuracy while upholding ethical standards related to user privacy. We believe that our approach can contribute to the development of recommendation systems that are both effective and ethical.

References

1. Berg, R.v.d., Kipf, T.N., Welling, M.: Graph convolutional matrix completion. arXiv preprint arXiv:1706.02263 (2017)
2. Botev, Z.I., Grotowski, J.F., Kroese, D.P.: Kernel density estimation via diffusion. *The Annals of Statistics* **38**(5), 2916 – 2957 (2010). <https://doi.org/10.1214/10-AOS799>, <https://doi.org/10.1214/10-AOS799>
3. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: International conference on machine learning. pp. 1597–1607. PMLR (2020)
4. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
5. Gao, H., Ji, S.: Graph u-nets. In: international conference on machine learning. pp. 2083–2092. PMLR (2019)
6. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: Teh, Y.W., Titterton, M. (eds.) Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics. Proceedings of Machine Learning Research, vol. 9, pp. 249–256. PMLR, Chia Laguna Resort, Sardinia, Italy (13–15 May 2010), <https://proceedings.mlr.press/v9/glorot10a.html>
7. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N., Weinberger, K. (eds.) Advances in Neural Information Processing Systems. vol. 27. Curran Associates, Inc. (2014)
8. Gutmann, M., Hyvärinen, A.: Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In: Proceedings of the thirteenth international conference on artificial intelligence and statistics. pp. 297–304. JMLR Workshop and Conference Proceedings (2010)
9. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 9729–9738 (2020)

10. He, X., Deng, K., Wang, X., Li, Y., Zhang, Y., Wang, M.: Lightgcn: Simplifying and powering graph convolution network for recommendation. In: Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval. pp. 639–648 (2020)
11. He, X., Liao, L., Zhang, H., Nie, L., Hu, X., Chua, T.S.: Neural collaborative filtering. In: Proceedings of the 26th international conference on world wide web. pp. 173–182 (2017)
12. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: International Conference on Learning Representations (ICLR) (2015)
13. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907 (2016)
14. Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., Soricut, R.: Albert: A lite bert for self-supervised learning of language representations. arXiv preprint arXiv:1909.11942 (2019)
15. Lee, J., Lee, I., Kang, J.: Self-attention graph pooling. In: International conference on machine learning. pp. 3734–3743. PMLR (2019)
16. Lin, Z., Tian, C., Hou, Y., Zhao, W.X.: Improving graph collaborative filtering with neighborhood-enriched contrastive learning. In: Proceedings of the ACM Web Conference 2022. pp. 2320–2329 (2022)
17. Liu, Y., Jin, M., Pan, S., Zhou, C., Zheng, Y., Xia, F., Yu, P.: Graph self-supervised learning: A survey. *IEEE Transactions on Knowledge and Data Engineering* (2022)
18. Rendle, S., Freudenthaler, C., Gantner, Z., Schmidt-Thieme, L.: Bpr: Bayesian personalized ranking from implicit feedback. arXiv preprint arXiv:1205.2618 (2012)
19. Sarwar, B., Karypis, G., Konstan, J., Riedl, J.: Item-based collaborative filtering recommendation algorithms. In: Proceedings of the 10th international conference on World Wide Web. pp. 285–295 (2001)
20. Tian, C., Xie, Y., Li, Y., Yang, N., Zhao, W.X.: Learning to denoise unreliable interactions for graph collaborative filtering. In: Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 122–132 (2022)
21. Wang, T., Isola, P.: Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In: International Conference on Machine Learning. pp. 9929–9939. PMLR (2020)
22. Wang, X., He, X., Wang, M., Feng, F., Chua, T.S.: Neural graph collaborative filtering. In: Proceedings of the 42nd international ACM SIGIR conference on Research and development in Information Retrieval. pp. 165–174 (2019)
23. Wei, Y., Wang, X., Li, Q., Nie, L., Li, Y., Li, X., Chua, T.S.: Contrastive learning for cold-start recommendation. In: Proceedings of the 29th ACM International Conference on Multimedia. pp. 5382–5390 (2021)
24. Wu, J., Wang, X., Feng, F., He, X., Chen, L., Lian, J., Xie, X.: Self-supervised graph learning for recommendation. In: Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval. pp. 726–735 (2021)
25. Wu, L., Lin, H., Tan, C., Gao, Z., Li, S.Z.: Self-supervised learning on graphs: Contrastive, generative, or predictive. *IEEE Transactions on Knowledge and Data Engineering* (2021)
26. Xie, Y., Xu, Z., Zhang, J., Wang, Z., Ji, S.: Self-supervised learning of graph neural networks: A unified review. *IEEE transactions on pattern analysis and machine intelligence* (2022)

27. Ying, R., He, R., Chen, K., Eksombatchai, P., Hamilton, W.L., Leskovec, J.: Graph convolutional neural networks for web-scale recommender systems. In: SIGKDD. pp. 974–983 (2018)
28. Yu, J., Yin, H., Xia, X., Chen, T., Cui, L., Nguyen, Q.V.H.: Are graph augmentations necessary? simple graph contrastive learning for recommendation. In: Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 1294–1303 (2022)
29. Yu, J., Yin, H., Xia, X., Chen, T., Li, J., Huang, Z.: Self-supervised learning for recommender systems: A survey. arXiv preprint arXiv:2203.15876 (2022)
30. Zhou, K., Wang, H., Zhao, W.X., Zhu, Y., Wang, S., Zhang, F., Wang, Z., Wen, J.R.: S3-rec: Self-supervised learning for sequential recommendation with mutual information maximization. In: Proceedings of the 29th ACM international conference on information & knowledge management. pp. 1893–1902 (2020)