

Contrastive Self-supervised Learning in Recommender Systems: A Survey

MENGYUAN JING, Shanghai Jiao Tong University, China

YANMIN ZHU*, Shanghai Jiao Tong University, China

TIANZI ZANG, Shanghai Jiao Tong University, China

KE WANG, Shanghai Jiao Tong University, China

Deep learning-based recommender systems have achieved remarkable success in recent years. However, these methods usually heavily rely on labeled data (i.e., user-item interactions), suffering from problems such as data sparsity and cold-start. Self-supervised learning, an emerging paradigm that extracts information from unlabeled data, provides insights into addressing these problems. Specifically, contrastive self-supervised learning, due to its flexibility and promising performance, has attracted considerable interest and recently become a dominant branch in self-supervised learning-based recommendation methods. In this survey, we provide an up-to-date and comprehensive review of current contrastive self-supervised learning-based recommendation methods. Firstly, we propose a unified framework for these methods. We then introduce a taxonomy based on the key components of the framework, including view generation strategy, contrastive task, and contrastive objective. For each component, we provide detailed descriptions and discussions to guide the choice of the appropriate method. Finally, we outline open issues and promising directions for future research.

CCS Concepts: • **Information systems** → **Recommender systems**.

Additional Key Words and Phrases: contrastive learning, self-supervised learning, unsupervised learning, survey, deep learning

ACM Reference Format:

Mengyuan Jing, Yanmin Zhu, Tianzi Zang, and Ke Wang. 2018. Contrastive Self-supervised Learning in Recommender Systems: A Survey. *J. ACM* 37, 4, Article 111 (August 2018), 39 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Recommender systems, as the most effective way to alleviate information overloading, have been an indispensable tool in daily life [20, 140]. They are intensively employed in a broad range of online services such as e-commerce platforms, social media, and music platforms. Owing to the ability to effectively capture the user-item relationships, deep learning techniques have been widely used in recommender systems [17, 40]. Despite their effectiveness, most deep learning-based methods

*corresponding author

Authors' addresses: Mengyuan Jing, jingmy@sjtu.edu.cn, Shanghai Jiao Tong University, No.800 Dongchuan Road, Minhang District, Shanghai, Shanghai, China, 200240; Yanmin Zhu, yzhu@sjtu.edu.cn, Shanghai Jiao Tong University, No.800 Dongchuan Road, Minhang District, Shanghai, Shanghai, China, 200240; Tianzi Zang, zangtianzi@sjtu.edu.cn, Shanghai Jiao Tong University, No.800 Dongchuan Road, Minhang District, Shanghai, Shanghai, China, 200240; Ke Wang, onecall@sjtu.edu.cn, Shanghai Jiao Tong University, No.800 Dongchuan Road, Minhang District, Shanghai, Shanghai, China, 200240.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Association for Computing Machinery.

0004-5411/2018/8-ART111 \$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

focus on supervised learning settings. The recommendation model is trained with abundant labeled data (i.e., user-item interactions). However, user-item interaction records are very sparse compared to the interaction space [2, 38]. Hence, these methods usually suffer from the problem of data sparsity [104]. Meanwhile, these methods are prone to the problem of over-fitting and generalization error [56].

Self-supervised learning (SSL) [56], as a novel learning paradigm, provides new insights to overcome aforementioned problems. The basic idea of SSL is to acquire transferable knowledge from the data itself without the need for manually annotated labels. This is achieved by solving auxiliary tasks (named pretext tasks). The acquired knowledge is then used in downstream tasks. Due to its efficiency, SSL has been widely used in many fields such as computer vision (CV) [13, 42, 68], natural language processing (NLP) [26, 61] and graph learning [88, 108]. Inspired by the success of SSL in other fields, there is growing interest in applying SSL to the area of recommendation.

Existing SSL-based recommendation methods can be classified into generative, contrastive, and predictive methods [135]. However, generative self-supervised learning is memory-consuming when trained on large-scale datasets. Predictive self-supervised learning often requires domain knowledge to generate labels, leading to increased costs and reduced generalization performance. In contrast, contrastive self-supervised learning (CL for brevity) has lightweight models and flexible designs since it is independent of the encoder structure and typically requires no domain knowledge. As a result, CL-based methods have witnessed significant growth in recent years, emerging as the dominant approach among SSL-based recommendation methods. Furthermore, the number of publications related to CL-based recommendations exceeds 50% of the total number of SSL-based recommendation publications in the ACM Digital Library¹. Considering this increasing trend, we aim to provide a timely and comprehensive review to summarize these CL-based methods in this paper.

Although there have been several reviews [44, 49] on contrastive learning, they mainly focus on methods in CV and NLP without reviewing CL-based recommendation methods. However, due to the uniqueness of the recommendation, it is difficult to apply existing CL-based methods from other fields to recommendation. Specifically, in CV/NLP, models usually deal with dense input data and treat each data instance as isolated. However, in recommender systems, the input data are extremely sparse (e.g., one-hot ID and categorical features of users/items) and there is homophily between users or items. Moreover, various recommendation tasks are unique to recommender systems, such as bundle recommendation and multi-behavior recommendation. Furthermore, several reviews on SSL-based graph learning [57, 106, 117] also include some CL-based recommendation methods. However, these reviews only provide a limited introduction and lack in-depth discussions on recommendation. Therefore, these reviews cannot provide sufficient insights into CL-based recommendation. Considering the unique characteristics of recommendation, a comprehensive survey is necessary to thoroughly review CL-based recommendation methods.

In the field of recommendation, the most relevant survey is [135]. This survey reviews SSL-based recommendation methods, including some CL-based methods. Compared to [135], our survey purely focuses on CL-based recommendation and provides a more comprehensive and detailed analysis of this topic. Specifically, our survey has the following differences. Firstly, we present a more rational, fine-grained, and comprehensive taxonomy. For instance, we add model-based augmentation methods and methods without augmentation to view generation strategies and categorize contrastive tasks based on the characteristics of the contrastive instances. Secondly, we provide in-depth analyses of different options for key components of CL-based recommendation methods, guiding the selection of these components. This critical discussion is not present in [135].

¹<https://dl.acm.org/>

Finally, because of the increasing popularity of CL-based recommendation methods, we provide a more up-to-date review that summarizes recently published studies that were not included in [135].

To sum up, the key contributions of this paper are summarized as follows:

- We propose a general framework to unify the CL-based methods for recommendation. Based on the framework, we review existing research according to three key components: view generation, pretext task, and contrastive objective.
- We provide an up-to-date and comprehensive review of CL-based recommendation methods. We provide detailed descriptions and discussions for each key component to guide the choice of the appropriate method. We also introduce the relevant background knowledge to help readers easily understand CL-based recommendation.
- We identify the limitations of existing research and propose promising future directions for CL-based recommendation to inspire new research.

Paper Collection. We first adopt Google Scholar as the main search engine to collect related papers. Then, we search for related work from top-tier conferences and journals, such as SIGIR, KDD, WWW, AAAI, IJCAI, WSDM, CIKM, NuerIPS, ICML, TKDE, TOIS, etc. Specifically, we search with keywords including "self-supervised", "contrastive" in combination with "recommend", "collaborative filtering". To prevent omissions of relevant work, we further look through the references of each paper.

Survey Organization. The remainder of the survey is organized as follows. In Section 2, we introduce background knowledge. We then introduce the unified framework and taxonomy in Section 3. Section 4, Section 5 and Section 6 are the main contents, which review contrastive learning in recommender systems. In Section 7, we discuss the open problems and future directions. Finally, we conclude the survey in Section 8.

2 BACKGROUND

In this section, we introduce essential background knowledge about CL-based recommendation. First, we provide the definitions of relevant concepts. Then, we give a brief introduction to contrastive learning. At last, we introduce training strategies used in CL-based recommendation methods. In addition, we summarize the notations used in this survey in Table. 1.

2.1 Term Definitions

2.1.1 Supervised Learning, Unsupervised Learning, and Self-supervised Learning. Supervised learning refers to a learning paradigm that trains models with manually annotated labels. In contrast, unsupervised learning indicates the learning paradigm that trains models without using manually annotated labels. Self-supervised learning can be viewed as a subset of unsupervised learning as it requires no manually annotated labels. However, unlike other unsupervised learning methods (e.g., clustering) that concentrate on mining data patterns, self-supervised learning aims to generate supervision signals from the data itself, and models are still trained in supervised settings.

2.1.2 Pretext Tasks Versus Downstream Tasks. Pretext tasks are pre-designed tasks to be solved by models (e.g., node self-discrimination [104]). By learning the objective functions of the pretext tasks, models learn more generalized representations from unlabeled data, thus benefiting downstream tasks. Downstream tasks refer to tasks used to evaluate the quality of representations learned by models. Specifically, in recommender systems, downstream tasks are the recommendation tasks such as sequential recommendation and social recommendation. In general, solving downstream tasks requires manually annotated labels.

Table 1. Key notations.

| Notations | Discriptions |
|--|--------------------------------------|
| \mathcal{U} | The set of users |
| \mathcal{I} | The set of items |
| \mathcal{B} | The mini-batch |
| $\mathbf{h}, \mathbf{c}, \mathbf{g}, \mathbf{z}$ | The learned representations |
| f_θ | The encoder to learn representations |
| p_ω | The pretext decoder |
| q_ϕ | The downstream decoder |
| $\theta, \omega, \phi, \xi, \psi$ | Learnable parameters |
| λ, ρ, ϵ | The hyperparameter |
| \mathcal{T} | Data-based view generation strategy |
| \mathbf{L} | The location matrix |
| \mathbf{A} | The adjacent matrix of graph |
| \mathbf{X} | The feature matrix |
| \mathbf{H} | The representation matrix |
| \mathcal{G} | The graph |
| \mathcal{V}/\mathcal{E} | The set of the graph nodes/edges |
| s_u | The interaction sequence of user u |
| MI | Mutual information function |
| \parallel | Concatenation operation |
| \circ | The Hadamard product |
| $ \cdot $ | The length of a set |

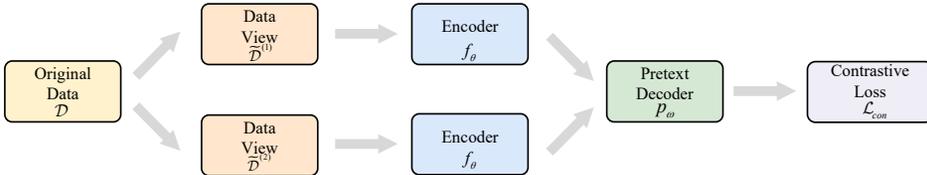


Fig. 1. Pipeline of contrastive learning.

2.2 Contrastive Learning

The core idea of contrasting learning (CL) is to maximize agreement between different views, where the agreement is usually measured by Mutual Information (MI). The general pipeline of CL is shown in Fig.1. In specific, two different data views are generated using view generation strategies. Then, representations in different views are generated by an encoder, which is usually shared by the two views. Finally, the model is optimized by contrastive loss to maximize the agreement between positive pairs and minimize the agreement between negative pairs. In general, positive pairs are the same instances from different views, while negative pairs are different instances from different views. Formally, contrastive self-supervised learning can be formulated as:

$$\theta^*, \omega^* = \arg \min_{\theta, \omega} \mathcal{L}_{con} \left(p_\omega \left(f_\theta \left(\tilde{\mathcal{D}}^{(1)} \right), f_\theta \left(\tilde{\mathcal{D}}^{(2)} \right) \right) \right) \quad (1)$$

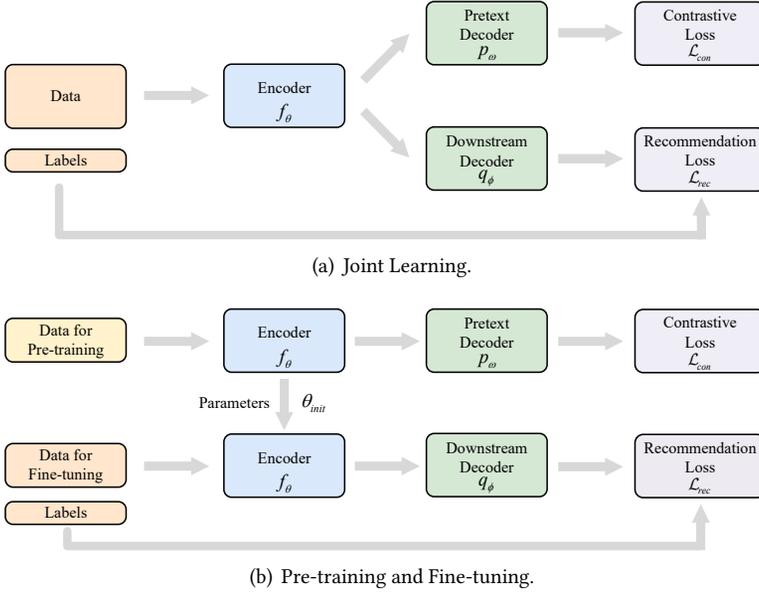


Fig. 2. Two types of training strategies for CL-based recommendation.

where $\tilde{D}^{(1)}$ and $\tilde{D}^{(2)}$ are two generated data views. $f_\theta(\cdot)$ is the (shared) encoder to learn representations of instances in different views. $p_\omega(\cdot)$ is the pretext decoder that estimates the agreement between two instances. \mathcal{L}_{con} denotes the contrastive loss.

2.3 Training Strategy

Currently, CL-based recommendation methods employ two typical training strategies: Pre-training and Fine-tuning, and Joint Learning. The detailed workflow of them is shown in Fig. 2.

2.3.1 Pre-training and Fine-tuning (P&F). In this strategy, the model is trained in two stages. In the pre-training stage, the encoder $f_\theta(\cdot)$ is first pre-trained with contrastive tasks. In addition, the pre-trained parameter θ_{init} is then used as the initialization parameter for the encoder $f_{\theta_{init}}(\cdot)$. In the fine-tuning stage, the pre-trained encoder $f_{\theta_{init}}(\cdot)$ is fine-tuned with the downstream decoder $q_\phi(\cdot)$ supervised by the recommendation task. The formulation of this strategy can be defined as:

$$\begin{aligned} \theta_{init}, \omega^* &= \arg \min_{\theta, \omega} \mathcal{L}_{con}(f_\theta, p_\omega) \\ \theta^*, \phi^* &= \arg \min_{\theta_{init}, \phi} \mathcal{L}_{rec}(f_{\theta_{init}}, q_\phi) \end{aligned} \quad (2)$$

where \mathcal{L}_{con} is the contrastive loss and \mathcal{L}_{rec} is the recommendation loss. q_ϕ is the downstream decoder.

2.3.2 Joint Learning (JL). In this strategy, the encoder $f_\theta(\cdot)$ is jointly trained with the pretext tasks and downstream tasks (i.e., recommendation tasks). Moreover, the encoder is usually shared by pretext and recommendation tasks. This strategy can be considered a type of multi-tasking learning strategy, in which the contrastive pretext task is the auxiliary task to regularize the recommendation task. The loss function consists of both contrastive loss and recommendation loss. The learning

objective can be formalized as:

$$\theta^*, \omega^*, \phi^* = \arg \min_{\theta, \omega, \phi} [\mathcal{L}_{rec}(f_{\theta}, q_{\phi}) + \lambda \mathcal{L}_{con}(f_{\theta}, p_{\omega})] \quad (3)$$

where λ is a trade-off hyperparameter that controls the contribution of \mathcal{L}_{con} .

2.3.3 Discussion. Compared to JL, P&F has better generalizability. Specifically, for different recommendation tasks, P&F only requires fine-tuning, while JL requires re-training. Additionally, when data for the target recommendation task is limited, P&F allows for pre-training with data from other recommendation tasks. However, due to the two-step training process, P&F is more complex compared to JL. Moreover, since the pre-trained model is trained without labels, it may not explicitly learn the features of a specific task, which could result in compromised performance on certain tasks.

When recommendation tasks are determined, JL usually achieves better recommendation performance and is easier to be implemented. Therefore, most of the existing CL-based approaches adopt JL. However, JL requires high computational resources since multiple tasks need to be trained simultaneously. Additionally, careful balancing of loss functions for different tasks is necessary to avoid instability or performance degradation. To summarize, if the primary goal is to improve recommendation performance on a specific task, it is recommended to choose JL. Conversely, if the goal is to achieve good performance on different recommendation tasks, P&F is the better choice.

3 TAXONOMY

In this section, we first propose a unified framework of CL-based recommendation methods. Then we introduce our proposed taxonomy with three perspectives.

3.1 Unified Framework

As introduced in Section 2, the general framework of CL-based methods is first to perform view generation strategies to obtain multiple views and then maximize the agreement of positive pairs in these views by conducting the contrastive pretext task. Specifically, given the data \mathcal{D} , K data views $\{\tilde{\mathcal{D}}^{(k)}\}_{k=1}^K$ are obtained through K data-based augmentations $\{\mathcal{T}_k(\cdot)\}_{k=1}^K$, which can be formulated as:

$$\tilde{\mathcal{D}}^{(k)} = \mathcal{T}_k(\mathcal{D}), k = 1, \dots, K \quad (4)$$

Then, encoders $\{f_{\theta_k}(\cdot)\}_{k=1}^K$ are applied to generate representations $\{\mathbf{h}_k\}_{k=1}^K$ for each data view. Formally, we have

$$\mathbf{h}_k = f_{\theta_k}(\tilde{\mathcal{D}}^{(k)}), k = 1, \dots, K \quad (5)$$

In addition, $\{\mathbf{h}_k\}_{k=1}^K$ may have different scales depending on the type of pretext tasks. For example, it can be a representation of an item or a representation of a sequence that consists of multiple items.

During training, contrastive learning is to maximize the agreement between representations of positive pairs $(\mathbf{h}_i, \mathbf{h}_j)$ in two views. Moreover, the mutual information $\mathcal{MI}(\mathbf{h}_i, \mathbf{h}_j)$ is usually applied to measure the agreement. The contrastive objective can be defined as:

$$\max_{\{\theta\}_{i=1}^K} \sum_i \sum_{i \neq j} \lambda_{ij} \mathcal{MI}(\mathbf{h}_i, \mathbf{h}_j) \quad (6)$$

where $\lambda \in \{0, 1\}$, if the mutual information between \mathbf{h}_i and \mathbf{h}_j is calculated then $\lambda = 1$, otherwise $\lambda = 0$.

Since it is difficult to directly calculate mutual information, mutual information estimators are usually used instead. The estimation is calculated based on the discriminator $p_{\omega}(\cdot)$ (i.e., the pretext

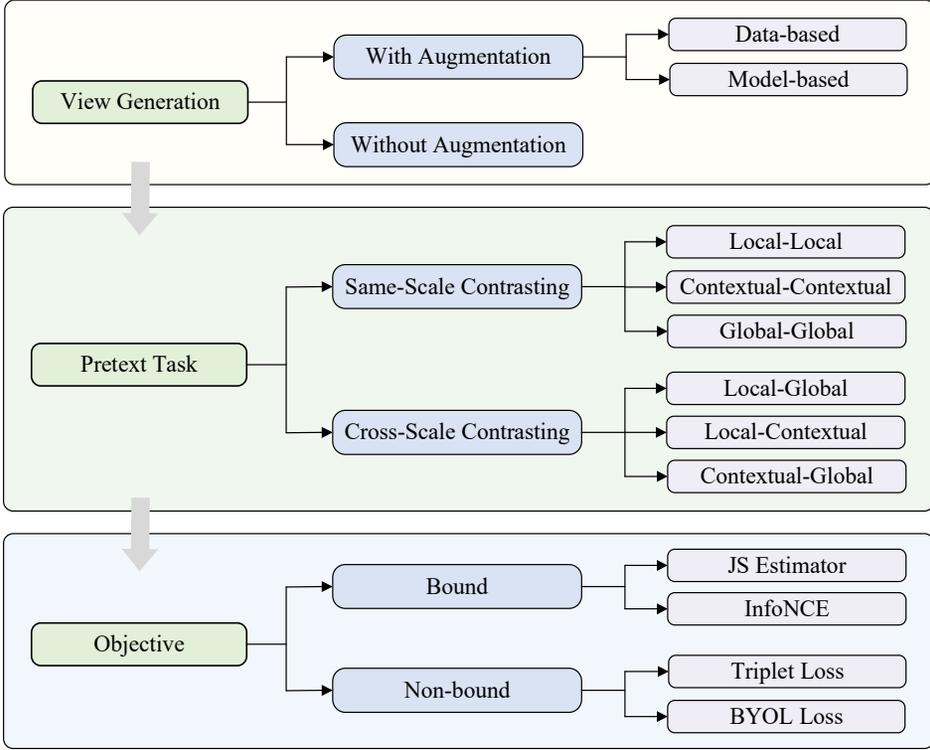


Fig. 3. Taxonomy of contrastive learning-based recommendation. "JS" refers to Jensen-Shannon. "InfoNCE" loss refers to the loss proposed by Oord et al. [68], where "NCE" stands for Noise-Contrastive Estimator. "BYOL" (Bootstrap Your Own Latent) loss refers to the contrastive loss proposed by Grill et al. [28].

decoder). Moreover, projection heads [18] can be optionally applied to $\{\mathbf{h}_k\}_{k=1}^K$, defined as:

$$\mathbf{z}_k = g_{\xi_k}(\mathbf{h}_k), k = 1, \dots, K \quad (7)$$

where $g_{\xi_i}(\cdot)$ is a projection head, which can be the Multi-Layer Perceptron (MLP) or linear projection. For the sake of convenience, we treat the projection head as part of the pretext decoder $p_{\omega}(\cdot)$. Then $f_{\theta^*}(\cdot)$ and $p_{\omega^*}(\cdot)$ can be obtained by learning Eq.(6). Furthermore, by utilizing $f_{\theta^*}(\cdot)$, the generated representations can be used for recommendation tasks. The recommendation task can be formulated as:

$$\theta^{**}, \phi^* = \arg \min_{\theta^*, \phi} \mathcal{L}_{rec}(q_{\phi}(f_{\theta^*}(\mathcal{D})), y) \quad (8)$$

where y denotes the labels. \mathcal{L}_{rec} is the supervised loss for recommendation tasks such as the cross-entropy (CE) loss.

3.2 Proposed Taxonomy

The differences among contrastive learning methods lie in three key components: view generation strategies, pretext tasks, and contrastive objectives. A CL-based recommendation method can be determined by specifying these components. Note that the encoder $f_{\theta}(\cdot)$ is not included in our taxonomy, as it is not the focus of CL-based recommendation and is determined by the specific recommendation tasks. Therefore, we propose a taxonomy based on these components as shown in Fig. 3. Table.2 shows representative works of CL-based recommendation.

Table 2. A summary of CL-based recommendation methods. For alphabets in "Pretext Task", L means Local; C means Contextual; G means Global. For acronyms used, "BPR" refers to Bayesian Personalized Ranking loss; "JL" refers to Joint Learning; "P&F" refers to Pre-training and Fine-tuning; "CF" refers to Collaborative Filtering; "KG" refers to Knowledge Graph.

| Model | Venue | Year | View Generation | Pretext Task | Objective | Training Strategy | Recommendation Task |
|----------------------------|--------|------|--|--------------|-----------|-------------------|---------------------|
| SGL [104] | SIGIR | 2021 | Node/Edge Dropout /Random Walk | L-L | InfoNCE | JL | Graph-based CF |
| DCL [60] | Arxiv | 2021 | Edge Dropout | L-L | InfoNCE | JL | Graph-based CF |
| GDCL [139] | DASFAA | 2022 | Graph Diffusion | L-L | InfoNCE | JL | Graph-based CF |
| SimGCL [134] | SIGIR | 2022 | Embedding Noise | L-L | InfoNCE | JL | Graph-based CF |
| XSimGCL [131] | Arxiv | 2022 | Embedding Noise | L-L | InfoNCE | JL | Graph-based CF |
| RocSE [130] | TOIS | 2023 | Embedding Noise | L-L | InfoNCE | JL | Graph-based CF |
| LightGCL [7] | ICLR | 2023 | SVD-based Augmentation | L-L | InfoNCE | JL | Graph-based CF |
| RGCF [85] | SIGIR | 2022 | Edge Perturbation | L-L | InfoNCE | JL | Graph-based CF |
| DCCF [74] | SIGIR | 2023 | Edge Dropout | L-L | InfoNCE | JL | Graph-based CF |
| GCARec [46] | PKDD | 2022 | Edge Dropout | L-L | InfoNCE | JL | Graph-based CF |
| LDA_GCL [43] | DASFAA | 2023 | Edge Perturbation | L-L | InfoNCE | JL | Graph-based CF |
| AdaGCL [45] | KDD | 2023 | Edge Perturbation | L-L | InfoNCE | JL | Graph-based CF |
| VGCL [126] | SIGIR | 2023 | Model-based Augmentation | L-L | InfoNCE | JL | Graph-based CF |
| SimRec [110] | WWW | 2023 | Model-based Augmentation | L-L | InfoNCE | JL | Graph-based CF |
| RGCL [76] | SIGIR | 2022 | Node Dropout | L-L | InfoNCE | JL | Review-based |
| MCLSR [99] | CIKM | 2022 | Without Augmentation | L-L | InfoNCE | JL | Sequential |
| HCCF [111] | SIGIR | 2022 | Edge Dropout | L-L | InfoNCE | JL | Graph-based CF |
| SGCCL [50] | WSDM | 2023 | Edge/Feature Dropout | L-L | InfoNCE | JL | Graph-based CF |
| LWC_KD [98] | CIKM | 2021 | Model-based Augmentation | L-L | InfoNCE | JL | Graph-based CF |
| MPT [35] | TOIS | 2023 | Node Masking/Substituting/Deleting | L-L | InfoNCE | P&F | Cold-start |
| MCCLK [146] | SIGIR | 2022 | Without Augmentation | L-L | InfoNCE | JL | KG-based |
| KACL [91] | WSDM | 2023 | Edge Dropout | L-L | InfoNCE | JL | KG-based |
| KGCL [124] | SIGIR | 2022 | Edge Dropout | L-L | InfoNCE | JL | KG-based |
| KGRec [122] | KDD | 2023 | Edge Dropout | L-L | InfoNCE | JL | KG-based |
| S-MBRec [29] | IJCAI | 2022 | Without Augmentation | L-L | InfoNCE | JL | Multi-behavior |
| MMCLR [107] | DASFAA | 2022 | Without Augmentation | L-L | BPR | JL | Multi-behavior |
| KMCLR [120] | WSDM | 2023 | Edge Dropout | L-L | InfoNCE | JL | Multi-behavior |
| SEPT [132] | KDD | 2022 | Edge Dropout/Predicted amples | L-L | InfoNCE | JL | Social |
| HGCL_S [12] | WSDM | 2023 | Without Augmentation | L-L | InfoNCE | JL | Social |
| CCDR [115] | KDD | 2022 | Subgraph Sampling | L-L | InfoNCE | JL | Cross-domain |
| ML-SAT [142] | CIKM | 2022 | Without Augmentation | L-L | InfoNCE | P&F | Cross-domain |
| DR-MTCDR [31] | TOIS | 2022 | Edge/Node Dropout | L-L | InfoNCE | JL | Cross-domain |
| COTREC [112] | CIKM | 2022 | Predicted Samples | L-L | InfoNCE | JL | Session-based |
| S ² -HHGR [138] | CIKM | 2021 | Node Dropout | L-L | JS | JL | Group |
| SGGCF [52] | WSDM | 2023 | Edge/Node Dropout | L-L | InfoNCE | JL | Group |
| CrossCBR [63] | KDD | 2022 | Edge/Message Dropout | L-L | InfoNCE | JL | Bundle |
| DCRec [123] | WWW | 2023 | Edge/Message Dropout | L-L | InfoNCE | JL | Sequential |
| HMG-CR [121] | Arxiv | 2021 | Without Augmentation | C-C | InfoNCE | JL | Multi-behavior |
| CHEST [92] | TOIS | 2023 | Subgraph Sampling | C-C | InfoNCE | P&F | HIN-based |
| KGIC [147] | CIKM | 2022 | /Path Dropout/Inserting/Substituting | C-C | InfoNCE | JL | KG-based |
| GCL4SR [141] | IJCAI | 2022 | Subgraph Sampling | C-C | InfoNCE | JL | Sequential |
| MISS [30] | ICDE | 2022 | Feature Extractor | C-C | InfoNCE | JL | Sequential |
| CL4SRec [116] | ICDE | 2022 | Item Masking/Shuffling/Cropping | G-G | InfoNCE | JL | Sequential |
| H ² SeqRec [54] | CIKM | 2021 | Item Masking/Cropping | G-G | InfoNCE | P&F | Sequential |
| CoSeRec [59] | Arxiv | 2021 | Item Shuffling/Cropping | G-G | InfoNCE | JL | Sequential |
| ContraRec [89] | TOIS | 2022 | /Masking/Substituting/Inserting | G-G | InfoNCE | JL | Sequential |
| ContraRec [89] | TOIS | 2022 | Item Masking/Shuffling/Overlapping | G-G | InfoNCE | JL | Sequential |
| TiCoseRec [21] | AAAI | 2023 | Ti-crop/Ti-mask/Ti-reorder /Ti-substitute/Ti-insert | G-G | InfoNCE | JL | Sequential |
| IOCR [53] | WSDM | 2023 | Item Shuffling/Cropping /Masking/Substituting/Inserting | G-G | InfoNCE | JL | Sequential |
| MCCM [93] | WSDM | 2023 | Feature Extractor /FItem Masking/Insttuting | L-L/G-G | InfoNCE | JL | News |
| CCL [5] | CIKM | 2021 | Item Masking/Sequence Generator | G-G | InfoNCE | P&F | Sequential |
| MIC [66] | CIKM | 2022 | Feature/Message Dropout | G-G | InfoNCE | JL | Sequential |
| EC4SRec [94] | CIKM | 2022 | Item Cropping/Masking/Shuffling | G-G | InfoNCE | JL | Sequential |
| DuoRec [73] | WSDM | 2022 | Message Dropout | G-G | InfoNCE | JL | Sequential |
| CBiT [23] | CIKM | 2022 | Item Masking/Message Dropout | G-G | InfoNCE | JL | Sequential |
| ContrastVAE [97] | CIKM | 2022 | Item Masking/Cropping/Shuffling /Message Dropout/Variational Dropout | G-G | InfoNCE | JL | Sequential |
| CLUE [18] | CIKM | 2022 | Item Masking/Cropping/Shuffling /Message Dropout | G-G | BYOL | P&F | Sequential |
| FDSA_CL [37] | TKDE | 2023 | Message Dropout | G-G | InfoNCE | JL | Sequential |
| EMKD [24] | SIGIR | 2023 | Item Masking | G-G | InfoNCE | JL | Sequential |
| MCLRec [71] | SIGIR | 2023 | /Model-based Augmentation | G-G | InfoNCE | JL | Sequential |
| DHCN [113] | AAAI | 2021 | Sequence Augmentor | G-G | InfoNCE | JL | Sequential |
| OD-Rec [114] | AAAI | 2021 | Feature Shuffling | G-G | JS | JL | Session-based |
| | SIGIR | 2022 | Without Augmentation | L-C | InfoNCE | JL | Session-based |

Continued on next page

Table 2. A summary of CL-based recommendation methods. (Continued)

| Model | Venue | Year | View Generation | Pretext Task | Objective | Training Strategy | Recommendation Task |
|---------------------------|--------|------|---|--------------|-----------|-------------------|---------------------|
| CGL [69] | TOIS | 2022 | Without Augmentation | G-G | JS | JL | Session-based |
| CFM [129] | CIKM | 2021 | Feature Dropout | G-G | InfoNCE | JL | Feature-based |
| CL4CTR [90] | WSDM | 2023 | Message/feature Dropout /Dimension Masking | G-G | BYOL | JL | CTR Prediction |
| CLCRec [103] | ACM MM | 2021 | Without Augmentation | G-G | InfoNCE | JL | Cold-start |
| NCL [55] | WWW | 2022 | Clustering | L-C | InfoNCE | JL | Graph-based CF |
| ICL [16] | WWW | 2022 | Clustering | L-C | InfoNCE | JL | Sequential |
| SITN [80] | AAAI | 2023 | Clustering | L-C | InfoNCE | P&F | Cross-domain |
| MHCN [133] | WWW | 2021 | Subgraph Sampling/Feature Shuffling | L-C | Triplet | JL | Social |
| SMIN [62] | CIKM | 2021 | Graph Diffusion | L-C | JS | JL | Social |
| CubeRec [15] | SIGIR | 2022 | Without Augmentation | L-C | Triplet | JL | Group |
| EGLN [125] | SIGIR | 2021 | Edge Dropout/Adding /Feature Shuffling | L-G | JS | JL | Graph-based CF |
| HGCL [6] | TMM | 2022 | Feature Shuffling | L-G | JS | JL | Micro-video |
| GroupIM [75] | SIGIR | 2020 | Without Augmentation | L-G | JS | JL | Group |
| BiGI [9] | WSDM | 2021 | Subgraph Sampling | C-G | JS | JL | Graph-based CF |
| MMSSL [102] | WWW | 2023 | Without Augmentation | C-G | InfoNCE | JL | Graph-based CF |
| C ² DSR [8] | CIKM | 2022 | Item Substituting | C-G | JS | JL | Cross-domain |
| SSI [136] | IJCAI | 2021 | Item Masking | C-G | InfoNCE | P&F | Sequential |
| SESRec [77] | SIGIR | 2023 | Without Augmentation | C-G | Triplet | P&F | Sequential |
| S ³ -Rec [143] | SIGIR | 2020 | Item Masking/Cropping | L-C/C-G/L-G | InfoNCE | P&F | Sequential |
| TCPSRec [84] | CIKM | 2022 | Sequence Dividing | L-C/L-G/C-C | InfoNCE | P&F | Sequential |

View Generation is the design of how to generate contrastive views. Depending on whether the augmentation is needed, we classify the view generation strategies into view generation with augmentation and without augmentation.

Pretext Task is the design of how to obtain supervision signals. Depending on the scale of the instances being contrasted, we classify the pretext tasks into same-scale contrasting and cross-scale contrasting.

Contrastive Objective is the design of how to measure mutual information. Depending on whether an estimation of lower-bound of mutual information is provided, we classify the contrastive objectives into bound objective and non-bound objective.

Note that the selection of these components depends on the characteristics of the input data and downstream tasks. For instance, sequence-based augmentation methods may not be suitable for graph data. For sequential recommendation, the pretext tasks generally contrast sequence representations as the primary objective is to learn high-quality sequence representations. It is worth noting that the selection of components is not entirely independent. Although the same pretext task can be performed with different view generation strategies or contrastive objectives, some may not be effective. For instance, feature-based augmentation methods may not be effective when the pretext task aims to model the sequential relationships of items. Therefore, when designing a CL-based recommendation method, we can first design the pretext task based on the specific recommendation task and then select view generation strategies and contrastive objectives accordingly.

4 VIEW GENERATION

Recent works [51, 86, 96] in other fields have shown that contrastive learning relies heavily on view generation, as generating multiple views facilitates models to explore richer underlying semantic information. In practice, if multiple data views naturally exist, such as interaction views and social networks in social recommendation, pretext tasks can be performed directly on these views. In addition, multiple views are not available in many scenarios, so augmentations are needed to generate contrastive views from the original data [13, 26, 32]. Therefore, we divide existing view generation strategies into view generation with augmentation and without augmentation.

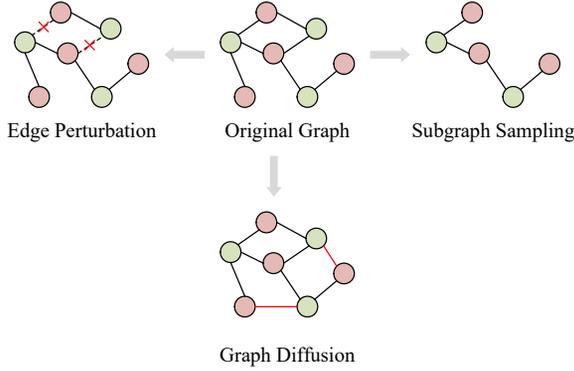


Fig. 4. Graph-based augmentation.

4.1 With Augmentation

Augmentation strategies can be categorized into data-based augmentation and model-based augmentation. The former generates views based on the data, while the latter is based on the model (i.e., the encoder).

4.1.1 Data-based Augmentation. Based on the type of data to be augmented, we classify data-based augmentation into graph-based augmentation, sequence-based augmentation, and feature-based augmentation.

Graph-based Augmentation. This strategy (shown in Fig.4) performs augmentations on the graph (e.g., interaction graph and social graph) to generate multiple views. Note that since the augmentations of node attributes in graphs are similar to feature-based augmentation, under this subcategory we only present the augmentations of the graph structure (shown in Fig. 4). Formally, given a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, graph-based augmentation transforms the adjacent matrix \mathbf{A} of \mathcal{G} , i.e., $\mathcal{T} = \mathcal{T}(\mathbf{A})$.

Edge perturbation. This strategy [31, 50, 52, 60, 63, 91, 92, 104, 111, 124, 125] generates graph views through randomly adding or dropping edges. It can be defined as:

$$\mathcal{T}(\mathbf{A}) = \mathbf{A} \circ (\mathbf{1} - \mathbf{L}) + (\mathbf{1} - \mathbf{A}) \circ \mathbf{L} \quad (9)$$

where \mathbf{L} is the location matrix. If $L_{ij} = 1$, the edge between i and j will be perturbed. Specifically, if $A_{ij} = 1, L_{ij} = 1$, the edge between i and j will be dropped. If $A_{ij} = 0, L_{ij} = 1$, an edge will be added between i and j . \mathbf{L} can be randomly sampled [104, 124] or manually set. Furthermore, \mathbf{L} can also be calculated adaptively [43, 45, 46] to keep important edges while perturbing possibly unimportant ones.

Graph Diffusion. The graph diffusion [62, 139] incorporates the global information to the original graph by creating new edges between nodes. It can be formulated as:

$$\mathcal{T}(\mathbf{A}) = \sum_{k=0}^{\infty} \Theta_k \mathbf{T}^k \quad (10)$$

where Θ_k is the weighting coefficient. \mathbf{T} denotes the generalized transition matrix. For example, SMIN [62] generates a substructure-aware adjacent matrix and injects it into the user-item interaction graph.

Subgraph Sampling. This strategy samples a node subset and corresponding edges to generate a subgraph as the data view. Existing methods usually obtain the node subset \mathcal{V}' by uniform

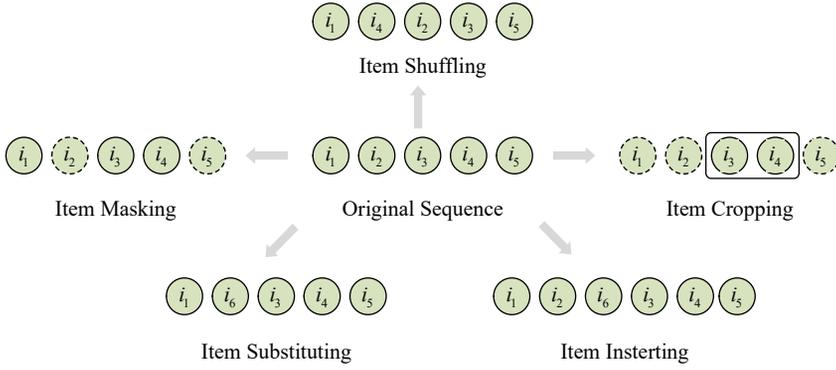


Fig. 5. Sequence-based augmentation.

sampling, ego-net sampling and knowledge-based sampling. *Uniform sampling* [31, 35, 76, 104, 138] uniformly samples a certain portion of nodes and corresponding edges to augment the views. Node dropout belongs to uniform sampling. For example, SGL [104] randomly drops a portion of nodes, which is denoted as \mathcal{V}_d . Therefore, the sampled node subset can be obtained by $\mathcal{V}' = \mathcal{V} - \mathcal{V}_d$. *Ego-net sampling* [9] samples the L -hop neighbors of each node in a graph, also known as the L -ego net. Therefore, the node subset can be represented as $\mathcal{V}' = \{j | d(i, j) \leq L\}$, where $d(v_i, v_j)$ is the shortest distance between node i and j . *Knowledge-based sampling* [92, 133] incorporates domain knowledge when sampling subgraph. For example, MHCN [133] designs three types of triangular motifs based on underlying semantics. Motifs specify high-order relations like "having a mutual friend".

Sequence-based Augmentation. This strategy (shown in Fig.5) performs augmentations on the user interaction sequences. Formally, give the interaction sequence s_u of user u , it can be formulated as $\tilde{s}_u = \mathcal{T}(s_u)$.

Item Shuffling. The item shuffling [53, 59, 89, 94, 97, 116] randomly shuffle a continuous subsequence of the interaction sequence to generate the augmented sequence:

$$\mathcal{T}(s_u) = [i_{u,1}, i_{u,2}, \dots, \tilde{i}_{u,k}, \dots, \tilde{i}_{u,k+l_s-1}, \dots, i_{u,|s_u|}] \quad (11)$$

where $[i_{u,k}, \dots, i_{u,k+l_s-1}]$ is shuffled as $[\tilde{i}_{u,k}, \dots, \tilde{i}_{u,k+l_s-1}]$. $l_c = \lceil \rho_s |s_u| \rceil$ is the length of the subsequence and $\rho_s \in [0, 1]$.

Item Cropping. The item cropping [53, 54, 59, 94, 97, 116] randomly chooses a continuous subsequence of the interaction sequence and can be represented as:

$$\mathcal{T}(s_u) = [i_{u,k}, i_{u,k+1}, \dots, i_{u,k+l_c-1}] \quad (12)$$

where $l_c = \lceil \rho_c |s_u| \rceil$ is the length of the subsequence and $\rho_c \in [0, 1]$ is the hyperparameter.

Item Masking. This strategy [23, 53, 54, 59, 89, 94, 97, 116] randomly chooses a portion of items in the interaction sequence and replaces them with a [mask] token, which can be formulated as:

$$\mathcal{T}(s_u) = [\tilde{i}_{u,1}, \tilde{i}_{u,2}, \dots, \tilde{i}_{u,|s_u|}] \quad (13)$$

where $\tilde{i}_{u,k} = [\text{mask}]$ if $i_{u,k}$ is masked, otherwise $\tilde{i}_{u,k} = i_{u,k}$.

Item Substituting. [53, 59] As dropout-based augmentation methods such as item masking may exacerbate the problem of data sparsity and cold-start, item substituting and item inserting are proposed. The item substituting randomly replaces a portion of items in the sequence with other

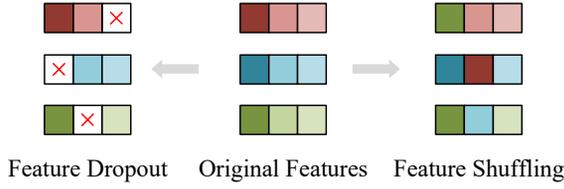


Fig. 6. Feature-based Augmentation.

items, which can be formulated as:

$$\mathcal{T} = [i_{u,1}, i_{u,2}, \dots, \tilde{i}_{u,k}, \dots, i_{u,|s_u|}] \quad (14)$$

where $\tilde{i}_{u,k}$ replaces $i_{u,k}$. Moreover, CoSeRec [59] substitutes items with highly correlated items to maintain the item correlations in the sequences.

Item Inserting. [53, 59] Fewer interactions are recorded in the interaction sequence than the complete behavior of the user, as interaction data from other sources may be missing. Therefore, the comprehensive user preferences and item correlations cannot be captured. To complete the sequence, CoSeRec [59] proposes the item inserting to generate the augmented sequence. Firstly, it randomly samples a portion of items in the sequence. Then, items that correlated to sampled items are inserted around them:

$$\mathcal{T}(s_u) = [i_{u,1}, i_{u,2}, \dots, \tilde{i}_{u,k}, i_{u,k}, \dots, i_{u,|s_u|}] \quad (15)$$

where $i_{u,k}$ is the sampled item and $\tilde{i}_{u,k}$ is the item related to it.

Feature-based Augmentation. Feature-based augmentation (shown in Fig.6) performs augmentations on the feature vectors, which can be categorical features or feature representations (e.g., embeddings). Given feature matrix \mathbf{X} , the augmented view is represented as $\tilde{\mathbf{X}} = \mathcal{T}(\mathbf{X})$.

Feature Dropout. The feature dropout (masking) [66, 90, 129] masks/drops a portion of the features and is formulated as:

$$\mathcal{T}(\mathbf{X}) = \mathbf{X} \circ (\mathbf{1} - \mathbf{L}) \quad (16)$$

where \mathbf{L} is the masking matrix that indicates the masking locations. If the j -th feature of i is masked/dropped, then $L_{ij} = 1$, otherwise $L_{ij} = 0$. Similar to edge perturbation, \mathbf{L} can be uniformly sampled or manually assigned. \circ is the Hadamard product.

Feature Shuffling. The feature shuffling [6, 113, 125, 133] perturbs the feature matrix by row or column. It can be formulated as:

$$\mathcal{T}(\mathbf{X}) = \mathbf{X}[\text{idx}_r, \text{idx}_c] \quad (17)$$

where idx_r and idx_c are the shuffled row index and the shuffled column index, respectively.

4.1.2 Model-based Augmentation. Model-based augmentation strategies (shown in Fig.7) generate views by perturbing the model (i.e., encoder). It is worth noting that unlike data-based augmentation strategies, which first generate different data views and then generate representations for each data view, model-based strategies directly generate different representations for the original data to perform the pretext task. It can be formulated as:

$$(\mathbf{h}, \mathbf{h}') = (f_{\theta}(\mathcal{D}), f'_{\theta'}(\mathcal{D})) \quad (18)$$

where $f_{\theta}(\cdot)$ and $f'_{\theta'}(\cdot)$ are the encoder and perturbed encoder, respectively. \mathbf{h}, \mathbf{h}' are representations output by $f_{\theta}(\cdot)$ and $f'_{\theta'}(\cdot)$, respectively. \mathcal{D} is the original data.

Message Dropout. This strategy [23, 58, 63, 73, 90] randomly masks the neurons in the layers for a certain dropout ratio [26]. Then, by applying different dropout masks, multiple views can be

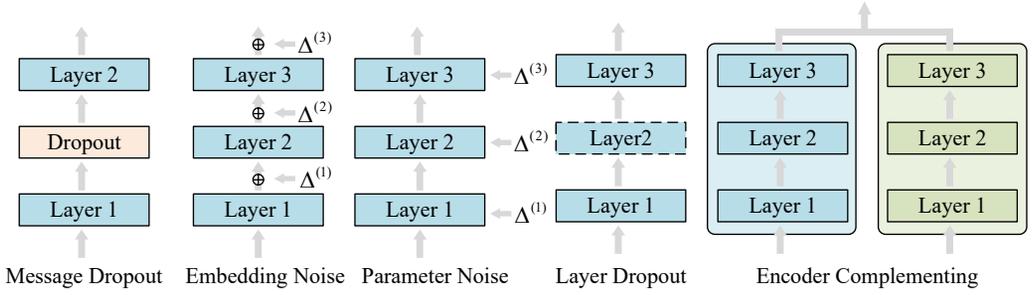


Fig. 7. Model-based augmentation.

obtained with the same input data. For example, DuoRec [73] applies two different dropout masks on the Transformer-based model to generate two different views.

Embedding Noise. This strategy [130, 131, 134] generates different views by adding different noises to original embeddings. Unlike feature-based augmentations that only perturb input embeddings or the final representations, this strategy adds noise to the embeddings at different layers of the encoder. It can be formulated as:

$$\tilde{E}_l = E_l + \Delta_l \quad (19)$$

where E_l is the original embedding and Δ_l is the perturbation noise at the l -th layer. In SimGCL [134], the $\Delta_l \sim U(0, 1)$ is the random uniform noise. In RocSE [130], $\Delta_l = \epsilon \cdot f_{\text{norm}}(f_{\text{shuffle}}(E_l))$, where f_{shuffle} and f_{norm} are the random shuffling and normalization operations, respectively. ϵ is a hyper-parameter.

Parameter Noise. This strategy [109] adds noises to the parameters of the encoder, which is formulated as:

$$\theta'_l = \theta_l + \epsilon \Delta_l \quad (20)$$

where θ_l and θ'_l are the original parameters and perturbed parameters of l -th layer, respectively. The Δ_l is the random noise, that can be sampled from the Gaussian distribution. ϵ is a hyper-parameter.

Architecture Perturbation. Unlike the above strategies that perturb learnable parameters in the model, some works generate different views by changing the model architecture. For example, SRMA [58] proposes *Layer Dropout* and *Encoder Complementing*. Specifically, the *Layer Dropout* randomly drops a portion of layers in the model during training to enable contrastive learning between shallow features and deep features. The *Encoder Complementing* uses a pre-trained encoder to generate representations. These representations are combined with the representations generated by the original encoder for contrastive learning. MA-GCL [27] proposes to perturb the architecture of graph neural network (GNN) encoders by varying the number and permutations of propagation and transformation operators.

4.2 Without Augmentation

The key idea of contrastive learning is to maximize the agreement between different views. Thus, if multiple views naturally exist, these views can be contrasted directly without additional augmentations. For example, in cross-domain recommendation, the two domains can be considered as two views. Therefore, some methods such as CCDR [115] and ML-SAT [142], directly perform contrastive learning between these domains. For knowledge graph-based recommendation, Some methods [91, 146] use the knowledge graph as a contrastive view. For multi-behavior recommendation, views can be constructed based on the auxiliary behavior data. For example, S-MBRec [29] treats each type of behavior as a view. Specifically, HMG-CR [121] build different hyper meta-graphs

Table 3. Comparison between different view generation strategies.

| | With Augmentation | | Without Augmentation |
|-----------------------|-------------------|-------------|----------------------|
| | Data-based | Model-based | |
| Trial-and-errors Free | ✗ | ✓ | ✓ |
| Domain Knowledge Free | ✓ | ✓ | ✗ |
| Generalizability | ✗ | ✓ | ✗ |
| Semantic Preservation | ✗ | ✓ | ✓ |

based on the hyper meta-paths constructed using the distance between auxiliary behavior and target behavior *buy*. In bundle recommendation, user-item interaction and user-bundle interaction can also be contrasted [63].

4.3 Discussion

Table. 3 shows the comparison between different view generation strategies. In specific, most existing CL-based recommendation methods adopt data-based augmentation strategies due to their ease of implementation. However, data-based augmentations are usually selected by manual trial-and-errors, which significantly limits the generalizability of these methods. In addition, some data-based augmentations destroy the semantic information of the original data, potentially harming recommendation performance [134].

Strategies without augmentation do not require trial-and-errors. These strategies typically use domain knowledge to build auxiliary views, which preserves the semantics of the data. However, domain knowledge is expensive and cannot be applied to other domains. Furthermore, since the views are fixed during model training, strategies without augmentation lack the introduction of randomness that helps to learn noise-invariant representations.

Compared to other strategies, model-based augmentations have better generalizability because they vary the learned representations without considering the original data. Although model-based augmentations require no trial-and-error and domain knowledge, settings such as the dropout ratio of messages/layers still require manual tuning. This limits their generalizability to some extent. Additionally, designing architecture-based perturbations is challenging.

Furthermore, many works [23, 63, 90, 97] adopt hybrid methods by combining multiple view generation strategies. In this way, the advantages of different strategies can be combined. However, some disadvantages may still exist. For instance, combining strategies without augmentation with data-based augmentations can be helpful in introducing randomness but data-based augmentation still requires manual trial-and-errors.

To summarize, selecting the appropriate strategy for view generation requires considering various factors. Here, we provide guidance for strategy selection based on typical issues in recommender systems. For the cold-start problem, strategies that drop data should be avoided as it exacerbates the problem [59]. Instead, views can be generated by adding/substituting interactions or using model-based strategies. Incorporating other data, like knowledge graphs or data from different recommendation domains, can also help mitigate data sparsity. To tackle the noise issue, data-based augmentation strategies are often more effective than other strategies as the noise mainly exists in the data. Perturbing data can make models more robust, thus mitigating the impact of noise [79]. Additionally, constructing a denoised data view based on metrics such as edge reliability degrees [85] can be helpful. Similarly, for addressing bias, a debiased view can also be constructed. Introducing other side information is also beneficial in reducing bias [11].

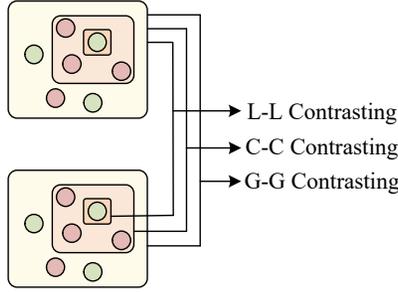


Fig. 8. Illustration of same-scale contrasting.

5 PRETEXT TASK

The goal of contrastive learning is to maximize the agreement between positive pairs (i.e., instances with same semantic information) and minimize the agreement between negative pairs (i.e., instances with unrelated semantic information). According to the scale of instances, we classify existing contrastive pretext tasks into two categories: same-scale contrasting and cross-scale contrasting.

Specifically, there are three contrastive scales: local, contextual, and global. The local scale usually represents the minimum granularity of the input data, while the global scale represents the maximum. For instance, in graph (sequence) data, the local scale represents the node (item/feature), and the global scale represents the whole graph (sequence). The contextual scale is between the local and global scales and represents the subgraph (subsequence).

5.1 Same-Scale Contrasting

Depending on the different scales being contrasted, same-scale contrasting (shown in Fig.8) can be further divided into three sub-types: local-local (L-L) contrasting, contextual-contextual (C-C) contrasting, and global-global (G-G) contrasting. Considering the unique characteristics of the recommendation tasks, we present existing methods based on their recommendation tasks.

5.1.1 Local-Local Contrasting. Methods under this category mainly discriminate the local representations (i.e., representation of users/items) and can be formulated as

$$\theta^*, \omega^* = \arg \min_{\theta, \omega} \mathcal{L}_{con} (p_{\omega} (\mathbf{h}_i, \mathbf{h}_j)) \quad (21)$$

where \mathbf{h}_i and \mathbf{h}_j are the representation of instance i and j in different views respectively. Furthermore, these representations are generated by encoder $f_{\theta}(\cdot)$, which is usually shared by different views.

Graph-based Collaborative Filtering. Depending on the types of graphs being contrasted, methods can be categorized into *contrasting on user-item graph* and *contrasting on different graphs*.

(i) *Contrasting on User-Item Graph.* As only one graph exists, methods under this category should perform augmentations on the user-item interaction graph to generate different views.

SGL [104] first applies contrastive learning to graph-based recommendation. Given a user-item interaction graph \mathcal{G} . It first generates two different graph views $\tilde{\mathcal{G}}^{(1)} = \mathcal{T}(\mathcal{G})$ and $\tilde{\mathcal{G}}^{(2)} = \mathcal{T}(\mathcal{G})$. \mathcal{T} is the data-based view generation strategy. Moreover, it utilizes three data-based augmentations including node dropout, edge dropout, and random walk (apply edge dropout at each layer).

Then, it utilizes LightGCN [39] as graph encoder $f_{\theta}(\cdot)$ to generate node representations $\mathbf{H}^{(1)} = f_{\theta}(\tilde{\mathcal{G}}^{(1)})$ and $\mathbf{H}^{(2)} = f_{\theta}(\tilde{\mathcal{G}}^{(2)})$. Afterward, it performs the node self-discrimination task. Specifically, it makes the representations of the same node (i.e., the positive pair) in different views similar while

making representations of different nodes (i.e., the negative pairs) in different views dissimilar. The contrastive loss of the user side can be formulated as:

$$\mathcal{L}_{con}^{user} = -\log \frac{\exp(p_{\omega}(\mathbf{h}_u^{(1)}, \mathbf{h}_u^{(2)}))}{\exp(p_{\omega}(\mathbf{h}_u^{(1)}, \mathbf{h}_u^{(2)})) + Neg} \quad (22)$$

where $\mathbf{h}_u^{(1)} \in \mathbf{H}^{(1)}$ and $\mathbf{h}_u^{(2)} \in \mathbf{H}^{(2)}$ are representations of user u . $p_{\omega}(\cdot)$ is the cosine similarity with a temperature parameter τ . $p_{\omega}(\mathbf{z}_u^{(1)}, \mathbf{z}_u^{(2)}) = (\mathbf{z}_u^{(1)})^T \mathbf{z}_u^{(2)} / \tau$ and $\mathbf{z}_u^{(1)} = \mathbf{h}_u^{(1)} / \|\mathbf{h}_u^{(1)}\|$. In addition, for efficiency, the in-batch negative sampling can be adopted, i.e., considering only different nodes of the same batch \mathcal{B} instead of using all other nodes as negative samples. Therefore, Neg is defined as

$$Neg = \sum_{v \in \mathcal{B}, v \neq u} \exp(p_{\omega}(\mathbf{h}_u^{(1)}, \mathbf{h}_v^{(2)})) \quad (23)$$

Note that $(\mathbf{h}_u^{(1)}, \mathbf{h}_v^{(2)})$ is the *inter-view* negative pairs. The loss of the item side \mathcal{L}_{con}^{item} can be obtained in the same way. The contrastive loss is $\mathcal{L}_{con} = \mathcal{L}_{con}^{user} + \mathcal{L}_{con}^{item}$. Finally, SGL adopts a joint learning strategy to optimize the contrastive loss and recommendation loss.

Based on the framework of SGL, several works are proposed. The main difference with SGL is in the view generation strategies. **DCL** [60] perturbs the edges in L -ego net of each node to obtain views. **GDCL** [139] generate new graph view using graph diffusion. Moreover, it constructs the *intra-view* negative pairs and the Eq.(23) can be rewritten as

$$Neg = \sum_{v \in \mathcal{U}, v \neq u} \exp(p_{\omega}(\mathbf{h}_u^{(1)}, \mathbf{h}_v^{(2)}) + p_{\omega}(\mathbf{h}_u^{(1)}, \mathbf{h}_v^{(1)})) \quad (24)$$

LightGCL [7] proposes a singular value decomposition (SVD)-based graph augmentation strategy to effectively distill global collaborative signals. In specific, SVD is first performed on the adjacency matrix. Then, the list of singular values is truncated to retain the largest K values and truncated matrices are used to reconstruct the adjacency matrix. The node contrastive learning is performed between the reconstructed graph and the original graph. **RGCL** [76] also performs edge contrastive learning. It maximizes the MI between the review representation and the corresponding interaction representation.

SimGCL [134], **XSimGCL** [131], and **RocSE**[130] generate views by adding uniform noises to node representations. Moreover, to reduce the computational complexity, XSimGCL [131] replaces the final-layer contrast with cross-layer contrasting. It only utilizes one Convolutional Network (GCN)-based encoder and contrasts embeddings of different layers:

$$\mathcal{L}_{con} = -\log \frac{\exp(\mathbf{h}_i^T \mathbf{h}_i^{l^*} / \tau)}{\sum_{j \in \mathcal{B}} \exp(\mathbf{h}_i^T \mathbf{h}_j^{l^*} / \tau)} \quad (25)$$

where \mathbf{h} is the node representation and \mathbf{h}^{l^*} is the representation at the l^* layer.

SimRec [110] proposes contrastive knowledge distillation by incorporating contrastive learning into knowledge distillation. It adopts a GCN-based encoder as the teacher model and an MLP-based model as the student model to generate node (user/item) representations. Furthermore, it maximizes the MI between the representations of the same node learned from the teacher model and the student model.

RGCF [85] constructs contrastive views based on the edge reliability degree. It first obtains the node structural feature by aggregating its one-hop neighbor representations. Then reliability

degree is calculated based on the similarity of the structural feature

$$\begin{aligned} \cos(\mathbf{h}_u^s, \mathbf{h}_i^s) &= \frac{\mathbf{h}_u^{s\top} \mathbf{h}_i^s}{\|\mathbf{h}_u^s\|_2 \cdot \|\mathbf{h}_i^s\|_2} \\ s_{u,i} &= (\cos(\mathbf{h}_u^s, \mathbf{h}_i^s) + 1) / 2 \end{aligned} \quad (26)$$

where \mathbf{h}_u^s and \mathbf{h}_i^s are the structure feature of user u and item i in the user-item interaction graph. $s_{u,i}$ is the reliability degree of the edge between u and i . Then, RGCF constructs a denoised graph and a diversity graph. The denoised graph is constructed by dropping edges with lower reliability degrees while the diversity graph is constructed by randomly adding edges with higher degrees. It maximizes the MI between representations of the same user in the two graphs.

DCCF [74] constructs two relation graphs to perform contrastive tasks. Specifically, it first obtains general representation \mathbf{h}^z and intent-aware representation \mathbf{h}^r for each node. Then two graph relation matrices \mathcal{G}^z and \mathcal{G}^r are generated using them. The calculation of \mathcal{G}^z can be formulated as:

$$\begin{aligned} \mathcal{M}_{u,i}^z &= (\cos(\mathbf{h}_u^z, \mathbf{h}_i^z) + 1) / 2 \\ \mathcal{G}^z &= \mathcal{M}^z \circ \mathbf{A} \end{aligned} \quad (27)$$

where \mathbf{A} is the original user-item interaction graph. \mathcal{G}^r can be obtained similarly. The node self-discrimination task is performed on these three graphs. Furthermore, the augmented graph views become learnable in this process. GCARec [46], LDA_GCL [43], and AdaGCL [45] also use a learnable strategy for generating graph views. Specifically, GCARec applies an MLP to obtain preserving probabilities for edges and uses them to generate graph views on these probabilities by sampling edges. LDA_GCL also generates views by obtaining probabilities but requires pre-trained models to generate representations for calculating probabilities. AdaGCL utilizes a generative model and a denoised model to generate graph views.

Similar to AdaGCL, VGCL [126] generates graph views based on a generative model, but instead of generating graphs, it directly obtains node representations for contrastive learning. VGCL also maximizes the MI between the representations of nodes in the same cluster.

(ii) *Contrasting on Different Graphs*. In addition to the user-item interaction graph, some works construct other graphs using interaction data. That is, views are usually generated without augmentation.

MCLSR [99] constructs three graphs based on the interaction sequences, including a user-item relation graph, an item-item relation graph, and a user-user relation graph. HCCF [111] constructs two views, including a user-item interaction graph and a learnable hypergraph. The node self-discrimination same as SGL is performed in MCLSR and HCCF. SGCCl [50] constructs a user-user graph \mathcal{G}_{uu} and an item-item graph \mathcal{G}_{ii} and performs edge/feature dropout to augment them. Node self-discrimination is conducted on the $\tilde{\mathcal{G}}_{uu}$ and $\tilde{\mathcal{G}}_{ii}$.

LWC_KD [98] incorporates contrastive knowledge distillation into incremental learning. In each time block, a user-item graph, user-user graph, and item-item graph are constructed. It proposes layer-wise structure-aware contrastive learning, which contrasts node representations of the same layer k between adjacent time blocks. It can be formulated as:

$$\mathcal{L}_{con} = \frac{-1}{|\mathcal{N}_i^{t-1}|} \sum_{i' \in \mathcal{N}_i^{t-1}} \log \frac{\exp(\mathbf{h}_{i,k}^{t-1} \cdot \mathbf{h}_{i',k}^t / \tau)}{\exp(\mathbf{h}_{i,k}^{t-1} \cdot \mathbf{h}_{i',k}^t / \tau) + \sum_{i \in \text{Neg}^{t-1}} \exp(\mathbf{h}_{i,k}^{t-1} \cdot \mathbf{h}_{i,k}^t / \tau)} \quad (28)$$

where t denotes the time block. $\mathbf{h}_{i,k}^{t-1}$ is the representation of node i . \mathcal{N}_i^{t-1} denotes the one-hop neighbors of i at $t-1$. Neg^{t-1} denotes nodes randomly selected from the unconnected nodes. MPT [35] extends PT-GNN [36] which performs reconstruction tasks and can only model the

intra-correlations. It leverages contrastive tasks to capture the inter-correlations within the data. Specifically, it samples subgraphs/paths for each user. Node dropout/substitution is applied to augment subgraphs/paths. The MI between representations of the same user in augmented subgraphs/paths is maximized.

Knowledge Graph-based Recommendation. Apart from interaction data, the knowledge graph (KG) is also utilized for CL-based recommendation, as it can bring rich semantic information.

Generally, CL-based recommendation methods using KG generate views by manual design. **MCCLK** [146] constructs three graph views, including user-item graph \mathcal{G}_{ui} , item-entity graph \mathcal{G}_{ie} and user-item-entity graph \mathcal{G}_{uie} . It maximizes the MI between the user representations $\mathbf{h}_u^{(ui)}$, $\mathbf{h}_u^{(uie)}$ in \mathcal{G}_{ui} and \mathcal{G}_{uie} and between the item representations $\mathbf{h}_i^{(ui)}$, $\mathbf{h}_i^{(ie)}$ in \mathcal{G}_{ui} and \mathcal{G}_{ie} . Furthermore, it generates a new representation for each item

$$\mathbf{h}'_i = \mathbf{h}_i^{(ui)} \parallel \mathbf{h}_i^{(ie)} \quad (29)$$

where \parallel is the concatenation operation. Then the MI between \mathbf{h}'_i and $\mathbf{h}_i^{(uie)}$ is maximized by performing node-self discrimination. **KACL** [91] maximizes the MI between representations of the same item in the augmented user-item graph and the augmented knowledge graph. Moreover, the augmented graphs are generated by automatically dropping unimportant edges.

KG can also be used to guide the generation of different user-item graph views. For example, **KGCL** [124] performs stochastic augmentation on the knowledge graph to generate two different views $\tilde{\mathcal{G}}_k^{(1)}$ and $\tilde{\mathcal{G}}_k^{(2)}$. The item consistency is $c_i = \cos(\mathbf{h}_i^{(1)}, \mathbf{h}_i^{(2)})$. $\cos(\cdot)$ is the cosine similarity function. $\mathbf{h}_i^{(1)}$ and $\mathbf{h}_i^{(2)}$ are item representations in $\tilde{\mathcal{G}}_k^{(1)}$ and $\tilde{\mathcal{G}}_k^{(2)}$, respectively. Then, the user-item interaction graph \mathcal{G}_{ui} is augmented using knowledge-guided data augmentation. Specifically, edge dropout is performed based on the probability calculated as follows:

$$\begin{aligned} w_{ui} &= \exp(c_i) \\ p'_{ui} &= \max\left(\frac{w_{ui} - w^{\max}}{w^{\max} - w^{\min}}, p_\tau\right) \\ p_{ui} &= p_a \cdot \mu_{p'} \cdot p'_{ui} \end{aligned} \quad (30)$$

where p_{ui} is the dropout probability of edge (u, i) in \mathcal{G}_{ui} . p_τ is the threshold. p'_{ui} is an intermediate variable that is integrated with the mean value $\mu_{p'}$. p_a is a strength controller. With the p_{ui} , masking vectors $\mathcal{M}_i \in \{0, 1\}$ are generated based on the Bernoulli distribution [65]. The augmented graphs is $\tilde{\mathcal{G}}_{ui}^{(i)} = (\mathcal{V}, \mathcal{E} \circ \mathcal{M}_i)$. Moreover, KGCL performs both intra-view contrasting and inter-view contrasting. **KGRec** [122] generates rationale scores of knowledge triplets based on attention mechanism to augment the KG and user-item graph. Besides, it only performs inter-view contrasting.

Multi-behavior Recommendation. For enhancing user intention modeling, multi-behavior recommendation methods incorporate multiple types of user behavior data, which can be used to build contrastive views. **S-MBRec** [29] adopts a star-style contrastive task, i.e., it only performs contrastive learning between the target behavior (usually the buy) and each auxiliary behavior. It samples positive samples based on the similarity under target behavior. The similarity is calculated by point-wise mutual information [127].

$$\begin{aligned} PMI(u, u') &= \log \frac{p(u, u')}{p(u)p(u')}, \\ p(u) &= \frac{|\mathcal{I}(u)|}{|\mathcal{I}|}, \\ p(u, u') &= \frac{|\mathcal{I}(u) \cap \mathcal{I}(u')|}{|\mathcal{I}|}, \end{aligned} \quad (31)$$

where $\mathcal{I}(u)$ is items that user u has interacted. \mathcal{I} is the item set and $|\mathcal{I}|$ is the number of items. If similarity $PMI(u, u') > t$, (u, u') are considered as positive pairs. t is the threshold. The positive samples of the items are selected in a similar way. Moreover, negative samples are selected randomly.

MMCLR [107] constructs a graph view (user-item graph) \mathcal{G} and a sequence view (multi-behavior sequence) \mathcal{S} . For each view, different behavior representations of the same users (e.g., \mathbf{h}_{u,b_1}^g and \mathbf{h}_{u,b_2}^g) are treated as positive pairs. It also maximizes the MI between overall representations of the same user (e.g., \mathbf{h}_u^g and \mathbf{h}_u^s) in different views. **KMCLR** [120] maximizes the MI between different behaviors of the same user. In addition, it performs knowledge-aware contrastive learning. It leverages a knowledge graph to guide the augmentation of the user-item graph under the target behavior $\mathcal{G}_{ui}^{b_i}$. It first calculates the consistency c_i of each item i like KGCL. Then the edge dropout probability is obtained by

$$\hat{p}_{ui} = \sigma(\mathbf{h}_u^T \mathbf{h}_i) \circ c_i, \quad (32)$$

$$p_{ui} = (1 - \text{Min_Max}(\hat{p}_{ui})) a + \text{Min_Max}(\hat{p}_{ui}) b$$

where \mathbf{h}_u and \mathbf{h}_i are the user representation and item representation in $\mathcal{G}_{ui}^{b_i}$, respectively. $\text{Min_Max}(\cdot)$ is the min-max normalization function. a and b are hyperparameters that control the value interval of p_{ui} . Then it performs edge dropout in a way similar to KGCL. Moreover, KMCLR adopts the node self-discrimination task.

Social Recommendation. For social recommendation, social networks are utilized to improve recommendation performance. Similar to KG-based recommendation, views can be generated based on manual design.

SEPT [132] constructs three graph views based on the interaction data and social network, including a preference view \mathcal{G}_r , friend view \mathcal{G}_f , and sharing view \mathcal{G}_s . \mathcal{G}_r is the user-item interaction graph. Other views are constructed based on two types of triangle motifs. Moreover, it leverages tri-training [145] to predict positive samples for each view. The Eq.(22) is changed as follows

$$\mathcal{L}_{con} = - \sum_{v \in \{r,s,f\}} \log \frac{\sum_{p \in \mathcal{P}_{u+}^v} p_{\omega}(\mathbf{h}_u^v, \tilde{\mathbf{h}}_p)}{\sum_{p \in \mathcal{P}_{u+}^v} p_{\omega}(\mathbf{h}_u^v, \tilde{\mathbf{h}}_p) + \sum_{j \in \mathcal{U} / \mathcal{P}_{u+}^v} p_{\omega}(\mathbf{h}_u^v, \tilde{\mathbf{h}}_j)} \quad (33)$$

where \mathcal{P}_{u+}^v is the set of predict positive samples. \mathbf{h}_u^v is the user representation in view v . p_{ω} is the cosine similarity with temperature parameter τ . $\tilde{\mathbf{h}}_p$ is the representation of user p in $\tilde{\mathcal{G}}$, which is obtained by performing random edge dropout on the joint graph of the user-item interaction graph and the social network.

HGCL_S [12] constructs three types of graphs: user-item graph, user-user graph, and item-item graph. It performs the node self-discrimination task on the corresponding graphs. Moreover, when generating node representations in the user-user graph and item-item graph, HGCL_S utilizes a meta-network.

Cross-domain Recommendation. Different domains in cross-domain recommendation can be considered as different views. For cross-domain recommendation, two types of contrastive tasks can be conducted, i.e., *single-domain* contrasting and *cross-domain* contrasting.

CCDR [115] performs both single-domain contrasting and cross-domain contrasting. For the single-domain contrasting, it samples two subgraphs of each node to generate two node representations. Then the MI between these two representations is maximized. For the cross-domain contrasting, It maximizes the MI between the representations of the same node in the source domain and the target domain. Besides, to extract more cross-domain knowledge between unaligned nodes,

CCDR maximizes the MI between the representation of an aligned node in the source domain and the representations of its target-domain neighbors.

ML-SAT [142] studies the multi-scenario problem, which can be viewed as a multi-domain recommendation problem. Moreover, it only performs the cross-domain contrasting between two different scenarios. It treats the representations of the same users/items in different domains as positive pairs and representations of other users/items in both scenarios as negative pairs.

DR-MTCDD [31] only performs sing-domain contrasting. It augments the user-item graph in each domain by edge/node dropout. For each domain, it generates K channel node representations. For each channel, MI between representations of the same node is maximized.

Group Recommendation. Group recommendation aims to recommend items that can fulfill the preferences of a collective of users. To improve group recommendation performance, **S²-HHGR** [138] builds a hierarchical hypergraph based on the user-item, group-item, and user-group interactions. It applies double-scale node dropout strategies including coarse- and fine-grained dropout on the hypergraph. The former drops users in all groups. The latter only drops some nodes in a specific group while other groups still contain these users. It performs node self-discrimination on the coarse- and fine-grained user representations as follows:

$$\begin{aligned} \mathcal{L}_{con} &= -\log \sigma(p_\omega(\mathbf{h}'_u, \mathbf{h}''_u)) - \sum_{j=1}^n \left[\log \sigma(1 - p_\omega(\mathbf{h}'_j, \mathbf{h}''_u)) \right], \\ p_\omega(\mathbf{h}'_u, \mathbf{h}''_u) &= \sigma(\mathbf{h}'_u \mathbf{W} \mathbf{h}''_u{}^T) \end{aligned} \quad (34)$$

where \mathbf{h}'_u and \mathbf{h}''_u are the coarse-grained user representation and fine-grained user representation, respectively. n is the number of negative samples.

SGGCF [52] performs user node dropout and edge dropout on the user-item-group graph. Moreover, representations of the same nodes in original and augmented graphs are viewed as positive pairs. Besides, it performs cross-layer contrasting. The MI between initial node embedding and l -th (l is an even number) layer embedding is maximized.

Sequential Recommendation. DCRec [123] constructs sequential and collaborative views based on interaction sequences. It builds an item transition graph and an item co-interaction graph as collaborative views. Additionally, DCRec adjusts the strength of contrastive regularization by disentangling user conformity and actual interest, further enhancing its performance, which is formulated as:

$$\begin{aligned} \mathcal{L}_{con} &= - \sum_{i \in s_u} \omega_{u,i} \log \frac{\exp(\cos(\mathbf{h}_i^s, \mathbf{h}_i^t) / \tau)}{\sum_{i' \in \mathcal{I}} \exp(\cos(\mathbf{h}_i^s, \mathbf{h}_{i'}^t) / \tau)} \\ &\quad - (1 - \omega_{u,i}) \log \frac{\exp(\cos(\mathbf{h}_i^t, \mathbf{h}_i^c) / \tau)}{\sum_{i' \in \mathcal{I}} \exp(\cos(\mathbf{h}_i^t, \mathbf{h}_{i'}^c) / \tau)} \end{aligned} \quad (35)$$

where s_u is the interaction sequence of user u . \mathbf{h}_i^s , \mathbf{h}_i^t , \mathbf{h}_i^c are representations of item i in the sequence view, item transition graph and item co-interaction graph, respectively. $\omega_{u,i}$ is the conformity degree of (u, i) .

Session-based Recommendation. Based on the session data, **COTREC** [112] constructs two graph views, including item view (item-item graph) and session view (session-session graph). Inspired by SEPT, it utilizes co-training [14] to predict the positive and negative samples for each session. Note that the positive and negative samples are items. Furthermore, it maximizes the agreement between the representation of the last item in the session and representations of the predicted positive samples. The agreement between the representation of the last item and representations of the negative samples is minimized.

Bundle Recommendation. Based on user-item interaction, user-bundle interaction, and item-bundle affiliation, **CrossCBR** [63] constructs a bundle view (user-bundle graph) and an item view (user-item graph and bundle-item graph). It performs edge dropout and message dropout on these views and generates user representation and bundle representations. Moreover, MI between representations of the same user/bundle in corresponding views is maximized.

5.1.2 Contextual-Contextual Contrasting. For contextual-contextual contrasting, discrimination is performed on the contextual representations. It can be formulated as:

$$\theta^*, \omega^* = \arg \min_{\theta, \omega} \mathcal{L}_{con} (p_{\omega} (\mathbf{c}_i, \mathbf{c}_j)) \quad (36)$$

where \mathbf{c}_i and \mathbf{c}_j are contextual representations denoting data with similar contextual information.

HIN-based Recommendation. **CHEST** [92] extracts subgraphs from heterogeneous information networks (HIN) by considering path relevance and then applies data-based augmentations to them. Augmented subgraphs generated from the same original subgraphs are considered positive samples, while subgraphs connecting the same user to other items are considered negative samples. Furthermore, it employs generative pretext tasks, such as masked node/edge prediction, to capture local information. The model leverages curriculum learning [4] by pre-training in a progressive manner, with contrastive pretext tasks serving as the advanced course and generative pretext tasks as the elementary course. **KGIC** [147] constructs local and non-local graphs by combining user-item interactions and the knowledge graph. Moreover, *intra-view* contrasting and *inter-view* contrasting are performed among these graphs.

Multi-behavior Recommendation. To capture different user behavior patterns, **HMG-CR** [121] performs the hyper meta-graph discrimination task. It first constructs different hyper meta-graphs $\{\mathcal{G}_u^{(i)}\}_{i=1}^K$ for each user based on hyper meta-paths. In specific, the hyper meta-path is constructed based on the distance to the target behavior. Then representations of these hyper meta-graphs are generated through different encoders, i.e., $\mathbf{c}_u^{(i)} = f_{\theta}^i(\mathcal{G}_u^{(i)})$. Furthermore, it treats the representations of adjacent hyper meta-graphs ($\mathbf{c}_u^{(i-1)}, \mathbf{c}_u^{(i)}$) as negative pairs. The positive samples are generated by feeding the current hyper meta-graph into the encoder of the adjacent hyper meta-graph. For example, the positive sample of $\mathbf{c}_u^{(i)}$ is $\mathbf{c}_p = f_{\theta}^{i-1}(\mathcal{G}_u^{(i)})$.

Sequential Recommendation. **GCL4SR** [141] obtains subgraphs based on uniform sampling. Specifically, it first constructs a transition graph based on all user interaction sequences. For each sequence node, it randomly samples two different subgraphs. The subgraphs of the same sequence are positive pairs, and those of different sequences are viewed as negative pairs.

Unlike other works that perform contrasting on the graph, **MISS** [30] performs it on the feature vectors, which consist of categorical and sequential features. It leverages two CNN-based models to extract multiple interests contained in each feature vector. Moreover, it makes representations of the same interest similar and representations of different interests dissimilar.

5.1.3 Global-Global Contrasting. Methods under this category discriminate the global representations, which can be represented as:

$$\theta^*, \omega^* = \arg \min_{\theta, \omega} \mathcal{L}_{con} (p_{\omega} (\mathbf{g}_i, \mathbf{g}_j)) \quad (37)$$

where \mathbf{g}_i and \mathbf{g}_j are the global representation. Moreover, global-global contrasting is typically used in sequential recommendation and session-based recommendation, where \mathbf{g} represents a sequence or a session.

Sequential Recommendation. We introduce the methods for sequential recommendation according to the view generation strategies. Generally, existing methods adopt *data-based* augmentation and *model-based* augmentation.

(i) *Methods using Data-based Augmentation.* **CL4SRec** [116] utilizes three types of sequence-based augmentation: sequence cropping, sequence shuffling, and item masking. Given an interaction sequence of user u , it applies augmentation \mathcal{T} that is randomly sampled from three augmentations to generate different sequence views. $\tilde{s}_u^{(1)} = \mathcal{T}(s_u)$ and $\tilde{s}_u^{(2)} = \mathcal{T}(s_u)$. SASRec [48] is used as sequence encoder $f_\theta(\cdot)$ to generate sequence (global-level) representations $\mathbf{g}_u^{(1)} = f_\theta(\tilde{s}_u^{(1)})$ and $\mathbf{g}_u^{(2)} = f_\theta(\tilde{s}_u^{(2)})$. It performs the sequence self-discrimination task. Specifically, it makes the representations of augmented sequences from the same sequence (i.e., positive pairs) to be similar and those from different sequences (i.e., negative pairs) to be dissimilar.

$$\mathcal{L}_{con} = -\log \frac{\exp(p_\omega(\mathbf{g}_u^{(1)}, \mathbf{g}_u^{(2)}))}{\exp(p_\omega(\mathbf{g}_u^{(1)}, \mathbf{g}_u^{(2)})) + Neg} \quad (38)$$

where $p_\omega(\cdot)$ is the cosine similarity with temperature parameter τ . CL4SRec adopts in-batch negative sampling and the Neg is defined as

$$Neg = \sum_{v \in \mathcal{B}, v \neq u} \exp(p_\omega((\mathbf{g}_u^{(1)}, \mathbf{g}_v^{(1)})) + p_\omega((\mathbf{g}_u^{(1)}, \mathbf{g}_v^{(2)}))) \quad (39)$$

Similar to CL4SRec, **H²SeqRec** [54] contrasts sequence representations. Moreover, the model is pre-trained with contrastive tasks. **CoSeRec** [59] and **ContraRec** [89] adopt the same framework and objective as CL4SRec. Moreover, CoSeRec proposes two robust augmentation strategies, i.e., *item substituting* and *item inserting*. In addition to the augmented sequence from the same sequence, ContraRec also treats the sequences that have the same target item as positive samples. Based on CoseRec, **TiCoseRec** [21] proposes five data augmentation based on time intervals to generate uniform sequences. Moreover, a uniform (un-uniform) sequence is one in which the standard deviation value of its time interval series is relatively small (large). **IOCRec** [53] generates K intention representations for each augmented sequence. It maximizes MI between representations of the same sequence with the same intent, which can be formulated as:

$$\mathcal{L}_{con}^k = -\log \frac{\exp(p_\omega(\mathbf{g}_{u,k}^{(1)}, \mathbf{g}_{u,k}^{(2)}))}{\sum_{g \in \mathcal{N}} \exp(p_\omega(\mathbf{g}_{u,k}^{(1)}, \mathbf{g}))} \quad (40)$$

where $\{\mathbf{g}_{u,k}^{(1)}\}_{k=1}^K$ and $\{\mathbf{g}_{u,k}^{(2)}\}_{k=1}^K$ denote the intention representations of augmentation sequences. $p_\omega(\cdot)$ is the dot product. \mathcal{N} is the set of negative samples, including different intention representations of the same user and all intention representations of different users.

MCCM [93] incorporates contrastive learning into news recommendation. It augments sequences by performing item masking/substituting. Moreover, the masking/substituting probability is calculated based on the frequency of news.

$$p_i = \frac{\log(\text{count}(i))}{\max_{j \in \mathcal{I}} \log(\text{count}(j))} (p_{\max} - p_{\min}) + p_{\min} \quad (41)$$

where p_{\max} and p_{\min} are the predefined boundaries of the probability. $\text{count}(i)$ is the frequency of news i in the dataset. \mathcal{I} is the news set. It performs sequence self-discrimination similar to CL4SRec.

CCL [5] augments sequences by leveraging a data generator based on the mask-and-fill operation. In specific, it first masks a portion of items in each sequence. Then the data generator recovers the original sequence. The original sequence and recovered sequences from the same user are treated as

positive pairs and other recovered sequences from different users are negative samples. Moreover, CCL utilizes curriculum learning to conduct contrastive learning via an easy-to-difficult process.

In addition to interaction data, **MIC** [66] utilizes attributes of user/items. It constructs user/item sequences, which consist of user/item attributes and interaction records. Then, feature dropout is applied to generate different user/item sequences. Moreover, it treats k-nearest neighbors of the user/item as positive samples. It also clusters the users/items using k-means++ and treats users/items that are from different clusters as negative samples.

EC4SRec [94] further incorporates explanation methods into data-based augmentation. It obtains the importance scores of items to guide the augmentation based on explanation methods. The importance scores of items in sequence s_u are calculated as

$$score(s_u) = F_e(y_u, s_u, f_\theta), \quad (42)$$

where F_e is any explanation method. f_θ is the sequence encoder. y_u is the prediction probability for the next item. $score(s_u) = [score(i_{u,1}), \dots, score(i_{u,|s_u|})]$ and $score(i_{u,1})$ is the important score of item $i_{u,1}$. Then, it crops/masks/shuffles the items with the lowest scores to generate positive samples. It also adopts supervised positive sampling to sample sequences like ContraRec and takes the sequences with higher important scores among them as the positive samples. To generate negative samples, it masks the items with the highest scores. The cropped items also form negative samples. In addition to making positive samples from the same users similar, it also makes negative samples from different users more similar than positive samples from all users.

(ii) *Methods using Model-based Augmentation.* **DuoRec** [73] obtains different sequence representations based on message dropout. Specifically, the sequence is fed twice with different dropout masks in the Transformer-based encoder. Then it performs sequence self-discrimination on these representations. Similar to ContraRec, DuoRec also uses supervised positive sampling, where an interaction sequence with the same target item is randomly selected as a positive sample.

Several methods also apply dropout to generate different sequence representations. Besides message dropout, **CBiT** [23] also uses item masking. Specifically, it first generates K different sequences by item masking. Then, these sequences are fed into bidirectional Transformers with different dropout masks to generate representations. In addition, CBiT proposes multi-pair contrastive learning. All K augmented sequences are treated as positive samples. It defined the loss as

$$\mathcal{L}_{con} = - \sum_{x=1}^K \sum_{y=1}^K \mathbf{1}_{[x \neq y]} \log \frac{\exp(p_\omega(\mathbf{g}_u^{(x)}, \mathbf{g}_u^{(y)}))}{\exp(p_\omega(\mathbf{g}_u^{(x)}, \mathbf{g}_u^{(y)})) + \sum_{\mathbf{g} \in \mathcal{N}} \exp(\mathbf{g}_u^{(x)}, \mathbf{g})} \quad (43)$$

where \mathcal{N} is the set of negative samples, which are the sequence representations of different users.

ContrastVAE [97] incorporates contrastive learning into the Variational AutoEncoder. Besides message dropout and data-based augmentations, it utilizes the variational dropout that introduces a learnable Gaussian dropout rate during the sampling step. The selection of negative and positive pairs is similar to CL4SRec. **CLUE** [18] also adopts both message dropout and sequence-based augmentations. The main difference is that CLUE does not use negative samples. **FDSA_CL** [37] constructs feature sequences consisting of item features. It maximizes the MI between the representations of feature sequences and interaction sequences. In addition, it applies message dropout to generate K representations for each feature sequence.

EMKD [24] combines contrastive knowledge distillation with ensemble modeling. Specifically, it first generates K sequences by item masking. Then, it uses N networks that have the same structure but different initializations as an ensemble of sequence encoders. Each parallel network is a bidirectional Transformer encoder. EMKD conducts both intra- and inter-network contrasting. The intra-network contrasting is conducted by maximizing the MI between representations of the

original sequence and the augmented sequence generated by the same network. The inter-network contrasting is conducted by maximizing the MI between sequence representations generated by the different networks.

MCLRec [71] utilizes a learnable model-based augmentation method for view generation. It employs two different MLP-based augmenters to obtain representations to perform contrastive learning. Specifically, it first generates two sequences using data-based augmentation methods such as item masking and obtains their original representations through a shared encoder $f_{\theta}(\cdot)$. Then, these representations are fed into the augmenters to obtain augmented representations. MCLRec consequently contrasts the four representations. Additionally, it performs a meta-learning strategy to train the augmenters.

Session-based Recommendation. DHCN [113] constructs two hypergraphs (i.e., \mathcal{G}_h and \mathcal{G}_l) to represent intra- and inter-session information, respectively. In DHCN, representations of the same session (e.g., \mathbf{g}_s^h and \mathbf{g}_s^l) is the positive pair. Moreover, it perturbs the representations matrix \mathbf{H}^h with row- and column-wise shuffling. The negative samples are the perturbed representation $\tilde{\mathbf{g}}_s^h \in \tilde{\mathbf{H}}$ of the same session. $\tilde{\mathbf{H}}$ is the perturbed representation matrix. It maximizes MI between positive pairs and minimizes MI between negative pairs by

$$\mathcal{L}_{con} = -\log \sigma(p_{\omega}(\mathbf{g}_s^h, \mathbf{g}_s^l)) - \log \sigma(1 - p_{\omega}(\tilde{\mathbf{g}}_s^h, \mathbf{g}_s^l)) \quad (44)$$

where $p_{\omega}(\cdot)$ is the dot product and $\sigma(\cdot)$ is the sigmoid function.

OD-Rec [114] incorporates contrastive learning into knowledge distillation for session-based recommendation. Specifically, it maximizes the mutual information between the representations of the same session s that learned from the teacher model and student model (i.e., \mathbf{g}_s^{tea} , and \mathbf{g}_s^{stu}). The negative pairs are the $(\mathbf{g}_s^{tea}, \mathbf{g}_{s'}^{tea})$. CGL [69] constructs a global graph based on the similarity of sessions. In the graph, each session node is connected with their M most similar sessions. It maximizes the MI between a session node and its neighbors and minimizes the MI of unconnected sessions. Moreover, the session representation is obtained by aggregating the representations of items in it.

Feature-based Recommendation. CFM [129] adopts feature dropout to augment feature vectors. The vectors generated from the same feature vector are treated as positive pairs, while those generated from different feature vectors are treated as negative pairs. Moreover, to make the pretext task more difficult, it masks/drops the related features. CL4CTR [90] performs feature-based augmentation strategies to generate two different feature vectors of each sample. It minimizes the expected distance between representations of the same samples and does not use negative samples. Besides using feature representations, CLCRec [103] also generates item collaborative representations based on the interaction data. Therefore, it can be viewed as *hybrid recommendation* [1]. The feature representation and collaborative representation of the same item are treated as positive pairs and those of different items are treated as negative pairs.

5.2 Cross-Scale Contrasting

In the cross-scale contrasting (shown in Fig.9), the contrasting is conducted across different scales. Furthermore, according to the scale, we further divide this branch of methods into three sub-types: local-contextual (L-C) contrasting, local-global (L-G) contrasting, and contextual-global (C-G) contrasting.

5.2.1 Local-Contextual Contrasting. The local-contextual contrasting can be formulated as:

$$\theta^*, \omega^* = \arg \min_{\theta, \omega} \mathcal{L}_{con}(p_{\omega}(\mathbf{h}_i, \mathbf{c}_j)) \quad (45)$$

where \mathbf{h}_i is the local representation and \mathbf{c}_j is the contextual representation.

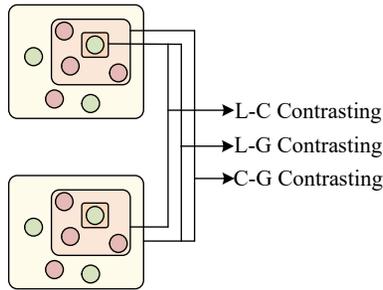


Fig. 9. Illustration of cross-scale contrasting.

Graph-based Collaborative Filtering. To capture contextual information, NCL [55] proposes a prototype-contrastive objective. In specific, for each item/user, the positive sample is the prototype of the cluster it belongs to, and the negative sample is the prototype of other clusters. The prototype is the representation of the cluster center. Moreover, the prototype-contrastive objective is learned with Expectation-Maximization (EM) algorithm. NCL also performs cross-layer contrasting. For each user/item, corresponding representations output from the even-numbered layer GNN are treated as positive samples.

Sequential Recommendation. ICL [16] relies on a framework similar to that of NCL [55]. The difference is that NCL clusters the representations of users or items, while ICL clusters the representations of sequences. Specifically, in ICL, the representation of a sequence can be viewed as a local view of the cluster it belongs to. In addition, each prototype represents the user intent. Besides, ICL also places contrasting between sequences like sequential recommendation methods in global-global contrasting.

Cross-domain Recommendation. SITN [80] only performs cross-domain contrasting. Like ICL, it clusters user sequence representations to represent the user interests and maximizes the MI between the representations of users and clusters. Specifically, the positive pairs are representations of a user and its corresponding cluster in different domains. The negative samples are new clusters. Additionally, SITN also maximizes the MI of user representations in different domains.

Social Recommendation. In addition to using clustering algorithms, contextual information is also modeled based on human prior knowledge. MHCN [133] designs three types of triangle motifs based on social relations and models them with a multi-channel hypergraph encoder. Moreover, a multi-channel hypergraph is proposed to model the information. In each channel, MHCN hierarchically maximizes the mutual information between the user representation, the user-centered sub-hypergraph representation, and the hypergraph representation. SMIN [62] constructs the context by generating a substructure-aware adjacent matrix based on the addition operations of different order adjacent matrices.

Group Recommendation. CubeRec [15] maximizes the MI between the representation of group intersection and the representations of users belonging to that intersection. Specifically, it generates a hypercube representation for each group. The group intersection representations can be viewed as contextual representations. It considers the representations of users within the intersection and the representation of the intersection as positive pairs. As such, CubeRec enhances the representations, allowing for better modeling of the common interests among different groups.

5.2.2 Local-Global Contrasting. This branch of methods conducts contrastive tasks between the local representation and global representation, which can be presented as

$$\theta^*, \omega^* = \arg \min_{\theta, \omega} \mathcal{L}_{con} (p_{\omega} (\mathbf{h}_i, \mathbf{g}_j)) \quad (46)$$

where \mathbf{h}_i is the local representation and \mathbf{g}_j is the global representation.

Graph-based Collaborative Filtering. **EGLN** [125] places the contrasting across the edge representation (i.e., the concatenation of representations of its connected nodes) and the global graph representations (i.e., the average of all edge representations). Specifically, the positive samples of \mathbf{g}_1 (representation of $\mathcal{G}^{(1)}$) are the edge representations in the graph $\mathcal{G}^{(1)}$ and the negative samples are the edge representations in the augmented graph $\mathcal{G}^{(2)}$. **HGCL** [6] constructs node-type specific homogeneous graphs to preserve the heterogeneity. Following DGI [88], it maximizes the MI between a node representation and corresponding graph presentation. Moreover, to incorporate the relationship between different node types, HGCL also designs a cross-type contrasting object. For each node type pair (t_1, t_2) , given the node-type specific homogeneous graph $\mathcal{G}^{(t_2)}$, the positive samples of it are the node representations in $\mathcal{G}^{(t_2)}$, and the negative samples are the node representation in the augmented graph of $\mathcal{G}^{(t_1)}$.

Group Recommendation. For group recommendation, there are typically two views: the user view and the group view, to perform contrastive tasks. **GroupIM** [75] maximizes the MI between the representations of group members and the representations of groups. In specific, it treats the user representation \mathbf{h}_i and group representation \mathbf{g}_j as positive pairs, where user i belongs to group j . Negative samples $\mathbf{h}_{\tilde{u}}$ are sampled from non-member user representations. Moreover, GroupIM introduces a preference-biased negative user sampling distribution $\mathcal{P}_{\mathcal{N}}(\tilde{u} | j)$. This distribution gives a higher likelihood to non-member users who have interacted with similar items as the target group j . The sampling distribution is defined as:

$$\mathcal{P}_{\mathcal{N}}(\tilde{u} | j) \propto \eta \mathcal{I} (\mathbf{x}_{\tilde{u}}^T \cdot \mathbf{x}_j > 0) + (1 - \eta) \frac{1}{|\mathcal{U}|} \quad (47)$$

where $\mathbf{x}_{\tilde{u}}$ and \mathbf{x}_j are the interacted items of user \tilde{u} and group j , respectively. $\mathcal{I}(\cdot)$ is the indicator function. η is the hyperparameter controlling the sampling bias, and $|\mathcal{U}|$ is the number of users.

5.2.3 Contextual-Global Contrasting. These methods contrast the contextual representation with global representation, which can be defined as:

$$\theta^*, \omega^* = \arg \min_{\theta, \omega} \mathcal{L}_{con} (p_{\omega} (\mathbf{c}_i, \mathbf{g}_j)) \quad (48)$$

Graph-based Collaborative Filtering. For each edge (i, j) , **BiGI** [9] performs ego-net sampling to get two subgraphs centered at i and j , respectively. Then, it adopts attention mechanism to obtain two contextual representations. The contextual representation of this edge \mathbf{s}_{ij} is the concatenation of contextual representations of i and j . Specifically, positive samples of \mathbf{g}_1 are the contextual representations of edges in $\mathcal{G}^{(1)}$ and negative samples are the contextual representation of edges in the augmented graph $\mathcal{G}^{(2)}$. **MMSL** [102] generates multiple modal-specific representations of users. Moreover, it maximizes MI between the modality-specific representation and the overall representation of the same user. The negative samples of a modality-specific representation \mathbf{c}_u^m are both the modality-specific representation \mathbf{c}_u^m and overall representation \mathbf{g}_u^m of different users.

Cross-domain Recommendation. **C²DSR** [8] obtains cross-domain sequences through merging single-domain sequences in chronological order. Then, it generates two augmented cross-domain sequences based on item substituting. Based on the self-attention mechanism, it generates item representations in the sequences. The representations of cross-domain sequences are obtained by

Table 4. Comparison between different pretext tasks.

| | Same-Scale | | | Cross-Scale | | |
|--------------------------------|------------|-----|-----|-------------|-----|-----|
| | L-L | C-C | G-G | L-C | L-G | C-G |
| Context Extraction Free | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ |
| Summary(Readout) Function Free | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Noise Level | Low | | | High | | |

aggregating the representations of items in each domain. Similarly, the representations of single-domain sequences are obtained. It maximizes the MI between the single-domain sequences and the original cross-domain sequences and minimizes the MI between the single-domain sequences and the augmented cross-domain sequences.

Sequential Recommendation. SSI [136] leverages contrastive learning to capture global consistency in sequential recommendation. Specifically, it samples a subsequence from the interaction sequence and masks the corresponding items in that sequence. It then maximizes the MI between the representation of the subsequence and that of the entire sequence. Negative samples are subsequences sampled from other sequences. **SESRec** [77] contrasts the recommendation sequence with the search sequence. It divides each sequence into a positive subsequence and a negative subsequence. Moreover, subsequences are generated based on similarity scores, which are obtained through the co-attention technique [118]. The higher-scoring elements of the sequence are assigned to the positive subsequence while the lower-scoring elements are assigned to the negative subsequence. Anchor is generated from the original sequence. It makes anchors similar to positive sequences and different from negative sequences through the triplet loss.

Moreover, based on the belonging relationships in interaction sequences, **S³-Rec** [143] and **TCPSRec** [84] conduct multiple contrastive tasks. Specifically, S³-Rec devises four objectives, including sequence-item, sequence-attribute, item-attribute, and sequence-subsequence mutual information maximization. TCPSRec performs item-sequence and item-subsequence contrasting as well as subsequence-subsequence contrasting at both coarse- and fine-grained periodicity levels. In TCPSRec, the subsequences are generated by dividing the interaction sequence when the time interval is greater than a threshold.

5.3 Discussion

Table.4 shows the comparison between different pretext tasks. Most existing methods utilize the same-scale contrasting, as it is only necessary to generate representations of the same scale. In contrast, cross-scale contrasting requires generating the corresponding representations for all the different scales. Furthermore, since existing CL-based methods usually use the shared encoder to generate representations, methods that adopt cross-scale contrasting usually require an additional module (i.e., summary function) to generate large-scale representations after generating small-scale representations. Take local-global contrasting in the graph-based recommendation as an example, it needs to learn the representation of each node first and then aggregate these representations to generate the graph representation using a readout function. The contextual contrasting also tends to have high complexity, because it needs to design the corresponding strategy (i.e., context extraction) to decide which part of the data to generate the contextual representation.

In addition, compared to same-scale contrasting, which usually aims to identify different instances, cross-scale contrasting focuses on modeling the belonging relationship between small and large scales. Considering the complexity, in the cross-scale contrasting, the negative samples are usually

selected from the small-scale representations. Moreover, cross-scale contrasting can introduce more information into the small-scale representations, but this may also introduce more noises, i.e., irrelevant information.

The objective of CL-based recommendation is to enhance recommendation performance by incorporating contrastive pretext tasks as auxiliary tasks. Therefore, the choice of pretext tasks depends on the specific recommendation tasks. For instance, in sequential recommendation where the goal is to learn sequence representations, utilizing (sub)sequence-level contrasting will perform better compared to solely relying on item-level contrasting. In recommendation tasks with inherent multiple views like KG-based recommendation, employing both inter- and intra-view contrasting can further enhance model performance [124, 147]. This is because, under the supervision of the CL, auxiliary views such as knowledge graphs can acquire better representations.

In addition, we can address common issues by designing the sampling strategy for negative and positive samples. For example, to solve the data noise problem, one approach is to divide the data into noise-free data and noisy data and then treat them as positive and negative data for the original data, respectively [72]. When sampling, focusing more on tailed samples can help mitigate the bias problem [11]. Likely, for the cold-start problem, paying more attention to users/items that have a few interactions can be helpful.

6 CONTRASTIVE OBJECTIVE

As introduced in Section. 3.1, the contrastive objective is to maximize the mutual information (MI) between different views. Specifically, given representations $(\mathbf{h}_i, \mathbf{h}_j)$ of instances (i, j) , the MI between them can be represented as:

$$MI(\mathbf{h}_i, \mathbf{h}_j) = KL(P(\mathbf{h}_i, \mathbf{h}_j) \| P(\mathbf{h}_i)P(\mathbf{h}_j)) = \mathbb{E}_{P(\mathbf{h}_i, \mathbf{h}_j)} \left[\log \frac{P(\mathbf{h}_i, \mathbf{h}_j)}{P(\mathbf{h}_i)P(\mathbf{h}_j)} \right] \quad (49)$$

where $KL(\cdot)$ is the Kullback-Leibler (KL) divergence. Contrastive learning aims to maximize the agreement between positive pairs and minimize the agreement between negative pairs. Moreover, the positive pair comes from the joint distribution $P(\mathbf{h}_i, \mathbf{h}_j)$ and the negative pair comes from the product of marginal distributions $P(\mathbf{h}_i)P(\mathbf{h}_j)$.

Depending on whether an estimation of lower-bound of mutual information is provided, we classify the contrastive objective into bound objective and non-bound objective.

6.1 Bound Objective

As calculating MI directly is difficult, lower-bounds are derived to estimate it [42], such as the Donsker-Varadhan estimator MI_{DV} [3, 22], the Jensen-Shannon estimator MI_{JS} [67], and the noise-contrastive estimator (InfoNCE) MI_{NCE} [33, 68]. Therefore, MI can be maximized by maximizing the lower-bound. Moreover, of these three estimators, only MI_{JS} and MI_{NCE} are currently used for CL-based recommendation.

6.1.1 Jensen-Shannon Estimator. Compared to the DV estimator, the Jensen-Shannon (JS) estimator enables a more efficient estimation of MI. It replaces the Kullback-Leibler divergence with the Jensen-Shannon divergence. The contrastive loss based on it can be defined as

$$\mathcal{L}_{JS} = -MI_{JS}(\mathbf{h}_i, \mathbf{h}_j) = -\mathbb{E}_P[\log(p_\omega(\mathbf{h}_i, \mathbf{h}_j))] - \mathbb{E}_{P \times \tilde{P}} \left[\log \left(1 - p_\omega(\mathbf{h}_i, \mathbf{h}'_j) \right) \right] \quad (50)$$

\mathbf{h}_i and \mathbf{h}_j are sampled from distribution P , and \mathbf{h}'_j is sampled from distribution \tilde{P} . $p_\omega(\cdot)$ is the discriminator (i.e., pretext decoder), which generates the agreement score of \mathbf{h}_i and \mathbf{h}_j . Moreover, there may be a projection head $g_\xi(\cdot)$ in $p_\omega(\cdot)$, which map representation \mathbf{h}_i to \mathbf{z}_i . Specifically, $g_\xi(\cdot)$

can be a linear mapping, MLP, or identical mapping. The $p_\omega(\cdot)$ can be inner product $\mathbf{z}_i^T \mathbf{z}_j$, the cosine similarity $\mathbf{z}_i^T \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|)$, or bi-linear transformation $\mathbf{z}_i^T \mathbf{W} \mathbf{z}_j$.

6.1.2 InfoNCE Estimator. InfoNCE is the most popular MI lower-bound adopted in CL-based methods for recommendation. The contrastive loss based on it can be formulated as

$$\mathcal{L}_{NCE} = -\mathcal{M}I_{NCE}(\mathbf{h}_i, \mathbf{h}_j) = -\mathbb{E}_P [p_\omega(\mathbf{h}_i, \mathbf{h}_j)] - \mathbb{E}_{K \sim \tilde{P}^N} \left[\log \frac{1}{N} \sum_{\mathbf{h}'_j \in K} e^{p_\omega(\mathbf{h}_i, \mathbf{h}'_j)} \right] \quad (51)$$

where K is the set of samples that consists of N random variables identically and independently distributed from \tilde{P} . Generally, $p_\omega(\cdot)$ is the cosine similarity with a temperature parameter τ , i.e., $p_\omega(\mathbf{z}_i, \mathbf{z}_j) = \mathbf{z}_i^T \mathbf{z}_j / \tau$ and $\mathbf{z}_i = \mathbf{h}_i / \|\mathbf{h}_i\|$. This is also known as the NT-Xent [78] loss.

In practice, InfoNCE is calculated on a mini-batch \mathcal{B} whose size is $N + 1$. Specifically, for each instance i in \mathcal{B} , the rest N instances are considered as negative samples. The loss based on InfoNCE can be

$$\mathcal{L}_{NCE} = -\frac{1}{N+1} \sum_{i \in \mathcal{B}} \left[\log \frac{e^{p_\omega(\mathbf{h}_i, \mathbf{h}_j)}}{\sum_{j \in \mathcal{B}} e^{p_\omega(\mathbf{h}_i, \mathbf{h}'_j)}} \right]. \quad (52)$$

6.2 Non-Bound Objective

In addition to the lower-bound MI estimators mentioned above, some other objectives are used to optimize contrastive learning, i.e., triplet loss and BYOL loss. However, they have not proven to be the lower-bound of MI and thus minimizing it does not guarantee to maximize mutual information.

6.2.1 Triplet Loss. The triplet loss does not minimize the agreement of the negative pairs but only makes the agreement of the positive pairs greater than that of the negative pairs. It is defined as:

$$\mathcal{L}_{Triplet}(\mathbf{h}_i, \mathbf{h}_j) = \mathbb{E}_{P \times \tilde{P}} \left[\max \left\{ p_\omega(\mathbf{h}_i, \mathbf{h}'_j) - p_\omega(\mathbf{h}_i, \mathbf{h}_j) + \epsilon, 0 \right\} \right] \quad (53)$$

where ϵ is the margin value. \mathbf{h}_i and \mathbf{h}_j are sampled from distribution P , and \mathbf{h}'_j is sampled from \tilde{P} . The discriminator p_ω can be calculated the agreement by $p_\omega(\mathbf{h}_i, \mathbf{h}_j) = \text{sigmoid}(\mathbf{h}_i, \mathbf{h}_j)$ or $p_\omega(\mathbf{h}_i, \mathbf{h}_j) = \|\mathbf{h}_i - \mathbf{h}_j\|$.

6.2.2 BYOL Loss. This objective is proposed by BYOL [28]. It only maximizes the agreement of positive pairs and does not use negative samples. It is defined as:

$$\mathcal{L}_{BYOL}(\mathbf{h}_i, \mathbf{h}_j) = \mathbb{E}_{P \times P} \left[2 - 2 \cdot \frac{[p_\psi(\mathbf{h}_i)]^T \mathbf{h}_j}{\|p_\psi(\mathbf{h}_i)\| \|\mathbf{h}_j\|} \right] \quad (54)$$

where \mathbf{h}_i and \mathbf{h}_j are sampled from P . $p_\psi(\cdot)$ is an online predictor. As it does not use negative samples to prevent collapse, other designs are needed. For example, BYOL [28] utilizes momentum encoders, stop gradient, etc.

6.3 Discussion

Table.5 shows the comparison between different contrastive objectives. Among all the contrastive objectives, InfoNCE is the most widely used due to its good performance. Moreover, both InfoNCE and JS estimate MI based on lower-bound, and Poole et al. [70] demonstrates that InfoNCE has a lower variance of the estimated MI than JS. However, InfoNCE requires a large number of negative samples and thus a large batch size during training. This leads to high computational and time complexity. In contrast, JS can achieve better performance when the batch size is small.

Table 5. Comparison between different contrastive objectives.

| | Bound | | Non-bound | |
|---------------------------|--------------------------|-------------------|--------------|-----------|
| | Jensen-Shannon Estimator | InfoNCE Estimator | Triplet Loss | BYOL Loss |
| Lower-bound MI estimation | ✓ | ✓ | ✗ | ✗ |
| Batch-Size Independence | ✓ | ✗ | ✓ | ✓ |
| Uniformity | ✓ | ✓ | ✓ | ✗ |
| Low Variance | ✗ | ✓ | N.A. | N.A. |

Triplet loss and BYOL loss are also independent of the large batch size. However, they lack theoretical support, i.e. no theory proves that maximizing them will achieve the goal of maximizing mutual information. Moreover, triplet loss just makes the agreement of negative pairs smaller than that of positive pairs. Therefore, selecting informative positive/negative samples that are difficult to discriminate can lead to better performance, whereas using random or easy samples leads to poor performance. Additionally, triplet loss can be sensitive to the choice of margin value, hence it requires careful adjustment.

BYOL loss is the most efficient since it does not require negative samples. However, BYOL loss does not contain the uniformity proposed by Wang and Isola [95], which suggests that normalized representations should be uniformly distributed over the unit hypersphere. Hence, it easily encounters the problem of collapse. Therefore, if BYOL loss is used, additional design is usually required to prevent it.

Overall, InfoNCE is generally a good choice. If the batch size is limited, JS may be a better alternative. Triplet loss can be adopted when positive and negative pairs should not be absolutely discriminated. For the bias issue, Yu et al. [134] demonstrate that InfoNCE can implicitly alleviate the popularity bias by making representations uniformly distributed across the unit hypersphere. This implies JS and Triplet loss can also partially mitigate the popularity. However, they use much fewer negative samples than InfoNCE, and thus may not mitigate bias as effectively. For the noise issue, re-weighting strategies [79] can be used to assign a higher weight to reliable (noise-free) data in the loss.

In addition, there are some contrastive recommendation loss functions such as BC loss [137], SSM loss [82, 105], and CCL loss [64]. However, these works utilize contrastive learning in a supervised manner. In specific, they treat interacted user-item pairs as positive pairs and non-interacted ones as negative pairs. As we focus on contrastive self-supervised learning, we will not discuss these supervised approaches in detail.

7 OPEN ISSUES AND FUTURE DIRECTIONS

While contrastive learning-based recommendation methods have achieved great success, there are still some open issues. In this section, we discuss these issues and outline some potential future research directions.

7.1 View Generation

View generation is a key component of CL-based methods. However, unlike in computer vision, where various data augmentation methods (e.g., resize, rotation, color distortion, etc.) are available, the way of generating views for CL-based recommendation is still not well explored. Specifically, most existing CL-based recommendation methods are limited to randomly removing some interactions or disrupting the order of the interaction sequence. Moreover, these methods are often based on intuitive designs and may not be applicable to downstream recommendation tasks [134].

Therefore, designing more effective view generation strategies is a promising future direction. Generally, the view generation strategies need to have the following properties: (1) Adaptability, the generated views should be adaptive to different tasks, as different tasks may use different types of data and require different information. (2) Efficiency, view generation strategies should not have high computational or time complexity. Moreover, dynamically updating the augmentation strategy during training is also a promising direction.

7.2 Pretext Task

By solving pretext tasks, the model acquires the knowledge from data for downstream tasks. Therefore, extracting useful knowledge is an important issue. For example, CGI [100] proposes an information bottleneck-based method that enables representations to capture the minimum sufficient information for the recommended task. However, it is designed for graph-based recommendation and is difficult to apply to other recommendation tasks. Moreover, as different tasks can capture different information, learning with multiple different pretext tasks can further improve recommendation performance. It is also worthwhile to further investigate the adaptive combination of different pretext tasks for the specific recommended task.

In contrastive pretext tasks, negative samples are essential, but obtaining informative negative samples is challenging. The commonly used uniform sampling strategy, which obtains negative samples by random sampling, suffers from false negatives. Besides, easy negative samples may degrade the performance of contrastive learning as they provide little information. Therefore, effective negative sampling strategies deserve further investigation. Some works [19, 47] explore this problem in computer vision. However, these methods are specifically designed for image data and are difficult to apply to recommendation methods. Moreover, since current methods require a large number of negative samples, efficient negative sample strategies also need to be explored.

7.3 Contrastive Objective

Most CL-based recommendation methods use InfoNCE as their objective function due to its simplicity and effectiveness. Although great success has been achieved, two issues need further exploration. First, the measurement of mutual information in InfoNCE is based on KL divergence. Therefore, it suffers from problems stemming from KL divergence (e.g., asymmetrical estimation and unstable training). Hence, better mutual information measurement is required but a few works [25] investigate this issue. Fan et al. [25] propose the Wasserstein discrepancy measurement based on the 2-Wasserstein distance to measure mutual information. In addition, it has only been applied to sequential recommendation, and its applicability to other recommendation tasks needs to be further explored. Second, current methods use mutual information to measure the agreement. However, mutual information has several shortcomings. Besides being hard to estimate, mutual information can also lead to suboptimal representations [87]. Therefore, exploring alternative measures of the agreement, such as \mathcal{V} -information proposed by Xu et al. [119], is a promising direction.

7.4 Miscellaneous

7.4.1 Meeting Real-World Recommendation. Most existing CL-based recommendation models are trained offline. However, in real-world recommendation scenarios, such as online shopping and news recommendation, large-scale interaction data are continuously generated and user preferences are dynamic. Offline trained models may suffer from the problem of information asymmetry as they rely only on historical user interaction data to make recommendations. Hence, exploring online learning strategies in CL-based recommendation to quickly capture dynamic preference trends would be a potential direction. Moreover, conversational recommendation [81, 144] is also proposed to address the information asymmetry. Specifically, the model makes recommendations based on

the multi-turn interaction (e.g., dialogues) with users. By leveraging real-time user feedback, the users' current preferences can be modeled. Combining contrastive learning with conversational recommendation methods would be an interesting direction to explore.

In addition, there are various types of data noise in recommender systems, like random clicks and false interactions, which influence the effectiveness of recommendation models. Consequently, recommendation denoising has gained considerable attention. However, current CL-based recommendation methods are mainly implicit denoising. They achieve denoising by learning perturbation-invariant representations to enhance the robustness of the model. Explicit strategies to address the issue of data noise are rarely explored. One potential approach is to construct a denoised data view to perform contrastive learning. Moreover, as denoising may impair recommendation diversity, leveraging contrastive learning to balance diversity and accuracy when denoising can also be a promising direction. For instance, RGCF [85] incorporates contrastive learning into recommendation denoising by maximizing the mutual information between the denoised graph and the diversity graph.

Apart from denoising, recommendation debiasing has also been a hot research topic in recent years. DCRec [123] investigates the combination of contrastive learning and debiasing techniques. It addresses the popularity bias by disentangling user conformity and interest to adjust the contrastive regularization strength. Moreover, some existing debiased recommendation models can also be used as backbones in combination with contrastive learning to address the popularity bias issue. For example, UnKD [10], which proposes an unbiased knowledge distillation approach, can be combined with contrastive learning by incorporating contrastive learning into knowledge transfer or by treating partition groups as contrastive views. Additionally, there are various other biases, such as selection bias, and incorporating contrastive learning to mitigate these biases is also worth exploring.

Additionally, in recommender systems, data is often multi-modal, including video, text, images, etc. These modalities contain rich information and are useful for improving recommendation performance, particularly when interaction data is sparse. Therefore, deriving useful knowledge from multi-modal data through contrastive learning is a promising direction. However, only a limited number of studies [34, 83, 102] have explored this area.

7.4.2 Learning with Advanced Techniques. With the rapid development of deep learning, many advanced techniques can be used to improve the performance of CL-based recommendation. One such technique is Knowledge Distillation (KD) [41], which is proposed to address the trade-off between high cost and model performance. Typically, KD first trains a large teacher model using the training set. Then a small student model is trained under supervision from the soft labels generated by the teacher model. Thus, there are naturally two views, namely the teacher view and the student view in KD. It can be easily combined with contrastive learning. Recently, several studies [24, 98, 110] have proposed to incorporate contrastive learning into KD and achieved promising recommendation performance. However, these explorations are still at an early stage, and there is still significant research potential in combining KD with contrastive learning. Moreover, some works [112, 132] leverage semi-supervised learning to obtain more supervision signals. In specific, SEPT [132] and COTREC [112] use tri-training and co-training to obtain more informative samples, respectively. CML [101] unifies contrastive learning and meta-learning by capturing meta-knowledge through contrastive learning. CCL [5] incorporates curriculum learning [4] into contrastive learning.

Additionally, there are still many new techniques that can be utilized for CL-based recommendation. For example, there are various choices for view generation strategies, contrastive tasks,

encoders, etc. in CL-based methods. Therefore, it is a promising direction to use Automated Machine Learning (AutoML) [128] to automatically select the appropriate method to reduce human effort.

8 CONCLUSION

In this survey, we present a comprehensive and systematic review of recent works in contrastive self-supervised learning-based recommendation. We first propose a unified framework and then introduce a taxonomy based on its key components, which include view generation strategy, pretext task, and contrastive objective. For each component, we provide detailed descriptions and discussions to guide the choice of the appropriate method. Finally, we discuss open issues and promising research directions for contrastive self-supervised learning-based recommendation in the future. We hope that this survey can provide both junior and experienced researchers with a comprehensive understanding of contrastive self-supervised learning-based recommendation and inspire future research in this area.

ACKNOWLEDGMENTS

This research is supported in part by National Science Foundation of China (No. 62072304), Shanghai Municipal Science and Technology Commission (No. 21511104700), the Shanghai East Talents Program, the Oceanic Interdisciplinary Program of Shanghai Jiao Tong University (No. SL2020MS032), and Zhejiang Aoxin Co. Ltd.

REFERENCES

- [1] Gediminas Adomavicius and Alexander Tuzhilin. 2005. Toward the Next Generation of Recommender Systems: A Survey of the State-of-the-Art and Possible Extensions. *IEEE Trans. Knowl. Data Eng.* 17, 6 (2005), 734–749.
- [2] Immanuel Bayer, Xiangnan He, Bhargav Kanagal, and Steffen Rendle. 2017. A generic coordinate descent framework for learning from implicit feedback. In *Proceedings of the 26th International Conference on World Wide Web*. 1341–1350.
- [3] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and Devon Hjelm. 2018. Mutual information neural estimation. In *International conference on machine learning*. PMLR, 531–540.
- [4] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*. 41–48.
- [5] Shuqing Bian, Wayne Xin Zhao, Kun Zhou, Jing Cai, Yancheng He, Cunxiang Yin, and Ji-Rong Wen. 2021. Contrastive curriculum learning for sequential user behavior modeling via data augmentation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 3737–3746.
- [6] Desheng Cai, Shengsheng Qian, Quan Fang, Jun Hu, Wenkui Ding, and Changsheng Xu. 2023. Heterogeneous Graph Contrastive Learning Network for Personalized Micro-Video Recommendation. *IEEE Transactions on Multimedia* 25 (2023), 2761–2773.
- [7] Xuheng Cai, Chao Huang, Lianghao Xia, and Xubin Ren. 2023. LightGCL: Simple Yet Effective Graph Contrastive Learning for Recommendation. In *The Eleventh International Conference on Learning Representations, ICLR 2023*. OpenReview.net.
- [8] Jiangxia Cao, Xin Cong, Jiawei Sheng, Tingwen Liu, and Bin Wang. 2022. Contrastive Cross-Domain Sequential Recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, Atlanta, GA, USA, October 17-21, 2022*, Mohammad Al Hasan and Li Xiong (Eds.). ACM, 138–147.
- [9] Jiangxia Cao, Xixun Lin, Shu Guo, Luchen Liu, Tingwen Liu, and Bin Wang. 2021. Bipartite graph embedding via mutual information maximization. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 635–643.
- [10] Gang Chen, Jiawei Chen, Fuli Feng, Sheng Zhou, and Xiangnan He. 2023. Unbiased Knowledge Distillation for Recommendation. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining, WSDM 2023*. ACM, 976–984.
- [11] Jiawei Chen, Hande Dong, Xiang Wang, Fuli Feng, Meng Wang, and Xiangnan He. 2023. Bias and Debias in Recommender System: A Survey and Future Directions. *ACM Trans. Inf. Syst.* 41, 3 (2023), 67:1–67:39.
- [12] Mengru Chen, Chao Huang, Lianghao Xia, Wei Wei, Yong Xu, and Ronghua Luo. 2023. Heterogeneous Graph Contrastive Learning for Recommendation. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 544–552.

- [13] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PMLR, 1597–1607.
- [14] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PMLR, 1597–1607.
- [15] Tong Chen, Hongzhi Yin, Jing Long, Quoc Viet Hung Nguyen, Yang Wang, and Meng Wang. 2022. Thinking inside The Box: Learning Hypercube Representations for Group Recommendation. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 1664–1673.
- [16] Yongjun Chen, Zhiwei Liu, Jia Li, Julian McAuley, and Caiming Xiong. 2022. Intent contrastive learning for sequential recommendation. In *Proceedings of the ACM Web Conference 2022*. 2172–2182.
- [17] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishu Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ipsir, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*. 7–10.
- [18] Mingyue Cheng, Fajie Yuan, Qi Liu, Xin Xin, and Enhong Chen. 2021. Learning Transferable User Representations with Sequential Behaviors via Contrastive Pre-training. In *IEEE International Conference on Data Mining, ICDM 2021*. 51–60.
- [19] Ching-Yao Chuang, Joshua Robinson, Yen-Chen Lin, Antonio Torralba, and Stefanie Jegelka. 2020. Debiased contrastive learning. *Advances in neural information processing systems* 33 (2020), 8765–8775.
- [20] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM conference on recommender systems*. 191–198.
- [21] Yizhou Dang, Enneng Yang, Guibing Guo, Linying Jiang, Xingwei Wang, Xiaoxiao Xu, Qinghui Sun, and Hong Liu. 2023. Uniform Sequence Better: Time Interval Aware Data Augmentation for Sequential Recommendation. (2023), 4225–4232.
- [22] Monroe D Donsker and SR Srinivasa Varadhan. 1983. Asymptotic evaluation of certain Markov process expectations for large time. IV. *Communications on pure and applied mathematics* 36, 2 (1983), 183–212.
- [23] Hanwen Du, Hui Shi, Pengpeng Zhao, Deqing Wang, Victor S Sheng, Yanchi Liu, Guanfeng Liu, and Lei Zhao. 2022. Contrastive Learning with Bidirectional Transformers for Sequential Recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 396–405.
- [24] Hanwen Du, Huanhuan Yuan, Pengpeng Zhao, Fuzhen Zhuang, Guanfeng Liu, Lei Zhao, Yanchi Liu, and Victor S. Sheng. 2023. Ensemble Modeling with Contrastive Knowledge Distillation for Sequential Recommendation. (2023), 58–67.
- [25] Ziwei Fan, Zhiwei Liu, Hao Peng, and Philip S Yu. 2023. Mutual Wasserstein Discrepancy Minimization for Sequential Recommendation. (2023), 1375–1385.
- [26] Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple Contrastive Learning of Sentence Embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021*.
- [27] Xumeng Gong, Cheng Yang, and Chuan Shi. 2023. MA-GCL: Model Augmentation Tricks for Graph Contrastive Learning. (2023), 4284–4292.
- [28] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. 2020. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems* 33 (2020), 21271–21284.
- [29] Shuyun Gu, Xiao Wang, Chuan Shi, and Ding Xiao. 2022. Self-supervised Graph Neural Networks for Multi-behavior Recommendation. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*, Luc De Raedt (Ed.). 2052–2058.
- [30] Wei Guo, Can Zhang, Zhicheng He, Jiarui Qin, Huifeng Guo, Bo Chen, Ruiming Tang, Xiuqiang He, and Rui Zhang. 2022. Miss: Multi-interest self-supervised learning framework for click-through rate prediction. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*. IEEE, 727–740.
- [31] Xiaobo Guo, Shaoshuai Li, Naicheng Guo, Jiangxia Cao, Xiaolei Liu, Qiongxiu Ma, Runsheng Gan, and Yunan Zhao. 2022. Disentangled Representations Learning for Multi-Target Cross-Domain Recommendation. *ACM Transactions on Information Systems* (2022).
- [32] Michael Gutmann and Aapo Hyvärinen. 2010. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 297–304.
- [33] Michael Gutmann and Aapo Hyvärinen. 2010. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 297–304.
- [34] Tengyue Han, Pengfei Wang, Shaozhang Niu, and Chenliang Li. 2022. Modality matches modality: Pretraining modality-disentangled item representations for recommendation. In *Proceedings of the ACM Web Conference 2022*. 2058–2066.

- [35] Bowen Hao, Hongzhi Yin, Jing Zhang, Cuiping Li, and Hong Chen. 2023. A Multi-Strategy-Based Pre-Training Method for Cold-Start Recommendation. *ACM Trans. Inf. Syst.* 41, 2, Article 31 (2023), 24 pages. <https://doi.org/10.1145/3544107>
- [36] Bowen Hao, Jing Zhang, Hongzhi Yin, Cuiping Li, and Hong Chen. 2021. Pre-training graph neural networks for cold-start users and items representation. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 265–273.
- [37] Yongjing Hao, Tingting Zhang, Pengpeng Zhao, Yanchi Liu, Victor S Sheng, Jiajie Xu, Guanfang Liu, and Xiaofang Zhou. 2023. Feature-Level Deeper Self-Attention Network With Contrastive Learning for Sequential Recommendation. *IEEE Trans. Knowl. Data Eng.* 35, 10 (2023), 10112–10124.
- [38] Ruining He and Julian McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *proceedings of the 25th international conference on world wide web*. 507–517.
- [39] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. 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*. 639–648.
- [40] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*. 173–182.
- [41] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. 2015. Distilling the Knowledge in a Neural Network. *CoRR* abs/1503.02531 (2015).
- [42] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. 2019. Learning deep representations by mutual information estimation and maximization. (2019).
- [43] Junjie Huang, Qi Cao, Ruobing Xie, Shaoliang Zhang, Feng Xia, Huawei Shen, and Xueqi Cheng. 2023. Adversarial Learning Data Augmentation for Graph Contrastive Learning in Recommendation. 13944 (2023), 373–388.
- [44] Ashish Jaiswal, Ashwin Ramesh Babu, Mohammad Zaki Zadeh, Debapriya Banerjee, and Fillia Makedon. 2020. A survey on contrastive self-supervised learning. *Technologies* 9, 1 (2020), 2.
- [45] Yangqin Jiang, Chao Huang, and Lianghao Xia. 2023. Adaptive Graph Contrastive Learning for Recommendation. (2023), 4252–4261.
- [46] Mengyuan Jing, Yanmin Zhu, Tianzi Zang, Jiadi Yu, and Feilong Tang. 2022. Graph Contrastive Learning with Adaptive Augmentation for Recommendation. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 590–605.
- [47] Yannis Kalantidis, Mert Bulent Sariyildiz, Noe Pion, Philippe Weinzaepfel, and Diane Larlus. 2020. Hard negative mixing for contrastive learning. *Advances in Neural Information Processing Systems* 33 (2020), 21798–21809.
- [48] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *2018 IEEE international conference on data mining (ICDM)*. IEEE, 197–206.
- [49] Adnan Khan, Sarah AlBarri, and Muhammad Arslan Manzoor. 2022. Contrastive self-supervised learning: a survey on different architectures. In *2022 2nd International Conference on Artificial Intelligence (ICAI)*. IEEE, 1–6.
- [50] Boyu Li, Ting Guo, Xingquan Zhu, Qian Li, Yang Wang, and Fang Chen. 2023. SGCCL: Siamese Graph Contrastive Consensus Learning for Personalized Recommendation. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 589–597.
- [51] Bohan Li, Yutai Hou, and Wanxiang Che. 2022. Data augmentation approaches in natural language processing: A survey. *AI Open* 3 (2022), 71–90.
- [52] Kang Li, Chang-Dong Wang, Jian-Huang Lai, and Huaqiang Yuan. 2023. Self-Supervised Group Graph Collaborative Filtering for Group Recommendation. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 69–77.
- [53] Xuwei Li, Aitong Sun, Mankun Zhao, Jian Yu, Kun Zhu, Di Jin, Mei Yu, and Ruiguo Yu. 2023. Multi-Intention Oriented Contrastive Learning for Sequential Recommendation. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 411–419.
- [54] Yicong Li, Hongxu Chen, Xiangguo Sun, Zhenchao Sun, Lin Li, Lizhen Cui, Philip S Yu, and Guandong Xu. 2021. Hyperbolic hypergraphs for sequential recommendation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 988–997.
- [55] Zihan Lin, Changxin Tian, Yupeng Hou, and Wayne Xin Zhao. 2022. Improving graph collaborative filtering with neighborhood-enriched contrastive learning. In *Proceedings of the ACM Web Conference 2022*. 2320–2329.
- [56] Xiao Liu, Fanjin Zhang, Zhenyu Hou, Li Mian, Zhaoyu Wang, Jing Zhang, and Jie Tang. 2021. Self-supervised learning: Generative or contrastive. *IEEE Transactions on Knowledge and Data Engineering* 35, 1 (2021), 857–876.
- [57] Yixin Liu, Ming Jin, Shirui Pan, Chuan Zhou, Yu Zheng, Feng Xia, and Philip Yu. 2022. Graph self-supervised learning: A survey. *IEEE Transactions on Knowledge and Data Engineering* (2022).
- [58] Zhiwei Liu, Yongjun Chen, Jia Li, Man Luo, S Yu Philip, and Caiming Xiong. 2021. Self-supervised Learning for Sequential Recommendation with Model Augmentation. (2021).

- [59] Zhiwei Liu, Yongjun Chen, Jia Li, Philip S Yu, Julian McAuley, and Caiming Xiong. 2021. Contrastive self-supervised sequential recommendation with robust augmentation. *CoRR abs/2108.06479* (2021).
- [60] Zhuang Liu, Yunpu Ma, Yuanxin Ouyang, and Zhang Xiong. 2021. Contrastive learning for recommender system. *arXiv preprint arXiv:2101.01317* (2021).
- [61] Lajanugen Logeswaran and Honglak Lee. 2018. An efficient framework for learning sentence representations. (2018).
- [62] Xiaoling Long, Chao Huang, Yong Xu, Huance Xu, Peng Dai, Lianghao Xia, and Liefeng Bo. 2021. Social Recommendation with Self-Supervised Metagraph Informax Network. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 1160–1169.
- [63] Yunshan Ma, Yingzhi He, An Zhang, Xiang Wang, and Tat-Seng Chua. 2022. CrossCBR: Cross-view Contrastive Learning for Bundle Recommendation. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 1233–1241.
- [64] Kelong Mao, Jieming Zhu, Jinpeng Wang, Quanyu Dai, Zhenhua Dong, Xi Xiao, and Xiuqiang He. 2021. SimpleX: A Simple and Strong Baseline for Collaborative Filtering. In *CIKM*. ACM, 1243–1252.
- [65] Albert W Marshall and Ingram Olkin. 1985. A family of bivariate distributions generated by the bivariate Bernoulli distribution. *J. Amer. Statist. Assoc.* 80, 390 (1985), 332–338.
- [66] Ping Nie, Yujie Lu, Shengyu Zhang, Ming Zhao, Ruobing Xie, William Yang Wang, and Yi Ren. 2022. MIC: Model-agnostic Integrated Cross-channel Recommender. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 3400–3409.
- [67] Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. 2016. f-gan: Training generative neural samplers using variational divergence minimization. 29 (2016), 271–279.
- [68] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).
- [69] Zhiqiang Pan, Fei Cai, Wanyu Chen, Chonghao Chen, and Honghui Chen. 2022. Collaborative graph learning for session-based recommendation. *ACM Transactions on Information Systems (TOIS)* 40, 4 (2022), 1–26.
- [70] Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. 2019. On variational bounds of mutual information. In *International Conference on Machine Learning*. PMLR, 5171–5180.
- [71] Xiuyuan Qin, Huanhuan Yuan, Pengpeng Zhao, Junhua Fang, Fuzhen Zhuang, Guanfeng Liu, Yanchi Liu, and Victor S. Sheng. 2023. Meta-optimized Contrastive Learning for Sequential Recommendation. *CoRR abs/2304.07763* (2023).
- [72] Yuqi Qin, Pengfei Wang, and Chenliang Li. 2021. The World Is Binary: Contrastive Learning for Denoising Next Basket Recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 859–868.
- [73] Ruihong Qiu, Zi Huang, Hongzhi Yin, and Zijian Wang. 2022. Contrastive Learning for Representation Degeneration Problem in Sequential Recommendation. In *WSDM '22: The Fifteenth ACM International Conference on Web Search and Data Mining, Virtual Event / Tempe, AZ, USA, February 21 - 25, 2022*. 813–823.
- [74] Xubin Ren, Lianghao Xia, Jiashu Zhao, Dawei Yin, and Chao Huang. 2023. Disentangled Contrastive Collaborative Filtering. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2023*. ACM, 1137–1146.
- [75] Aravind Sankar, Yanhong Wu, Yuhang Wu, Wei Zhang, Hao Yang, and Hari Sundaram. 2020. GroupIM: A Mutual Information Maximization Framework for Neural Group Recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020*. ACM, 1279–1288.
- [76] Jie Shuai, Kun Zhang, Le Wu, Peijie Sun, Richang Hong, Meng Wang, and Yong Li. 2022. A review-aware graph contrastive learning framework for recommendation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1283–1293.
- [77] Zihua Si, Zhongxiang Sun, Xiao Zhang, Jun Xu, Xiaoxue Zang, Yang Song, Kun Gai, and Ji-Rong Wen. 2023. When Search Meets Recommendation: Learning Disentangled Search Representation for Recommendation. (2023), 1313–1323.
- [78] Kihyuk Sohn. 2016. Improved deep metric learning with multi-class n-pair loss objective. (2016), 1849–1857.
- [79] Hwanjun Song, Minseok Kim, Dongmin Park, and Jae-Gil Lee. 2020. Learning from Noisy Labels with Deep Neural Networks: A Survey. *CoRR abs/2007.08199* (2020).
- [80] Guoqiang Sun, Yibin Shen, Sijin Zhou, Xiang Chen, Hongyan Liu, Chunming Wu, Chenyi Lei, Xianhui Wei, and Fei Fang. 2023. Self-Supervised Interest Transfer Network via Prototypical Contrastive Learning for Recommendation. (2023), 4614–4622.
- [81] Yueming Sun and Yi Zhang. 2018. Conversational recommender system. In *The 41st international acm sigir conference on research & development in information retrieval*. 235–244.
- [82] Hao Tang, Guoshuai Zhao, Yuxia Wu, and Xueming Qian. 2023. Multisample-Based Contrastive Loss for Top-K Recommendation. *IEEE Trans. Multim.* 25 (2023), 339–351.

- [83] Zhulin Tao, Xiaohao Liu, Yewei Xia, Xiang Wang, Lifang Yang, Xianglin Huang, and Tat-Seng Chua. 2022. Self-supervised learning for multimedia recommendation. *IEEE Transactions on Multimedia* (2022).
- [84] Changxin Tian, Zihan Lin, Shuqing Bian, Jinpeng Wang, and Wayne Xin Zhao. 2022. Temporal Contrastive Pre-Training for Sequential Recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 1925–1934.
- [85] Changxin Tian, Yuexiang Xie, Yaliang Li, Nan Yang, and Wayne Xin Zhao. 2022. Learning to Denoise Unreliable Interactions for Graph Collaborative Filtering. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 122–132.
- [86] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. 2020. What makes for good views for contrastive learning? *Advances in neural information processing systems* 33 (2020), 6827–6839.
- [87] Michael Tschannen, Josip Djolonga, Paul K Rubenstein, Sylvain Gelly, and Mario Lucic. 2020. On mutual information maximization for representation learning. (2020).
- [88] Petar Velickovic, William Fedus, William L Hamilton, Pietro Liò, Yoshua Bengio, and R Devon Hjelm. 2019. Deep graph infomax. *ICLR (Poster)* 2, 3 (2019), 4.
- [89] Chenyang Wang, Weizhi Ma, and Chong Chen. 2022. Sequential Recommendation with Multiple Contrast Signals. *ACM Trans. Inf. Syst.* (2022).
- [90] Fangye Wang, Yingxu Wang, Dongsheng Li, Hansu Gu, Tun Lu, Peng Zhang, and Ning Gu. 2023. CLACTR: A Contrastive Learning Framework for CTR Prediction. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 805–813.
- [91] Hao Wang, Yao Xu, Cheng Yang, Chuan Shi, Xin Li, Ning Guo, and Zhiyuan Liu. 2023. Knowledge-Adaptive Contrastive Learning for Recommendation. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 535–543.
- [92] Hui Wang, Kun Zhou, Xin Zhao, Jingyuan Wang, and Ji-Rong Wen. 2023. Curriculum Pre-Training Heterogeneous Subgraph Transformer for Top-N Recommendation. *ACM Transactions on Information Systems* 41, 1 (2023), 1–28.
- [93] Jingkun Wang, Yongtao Jiang, Haochen Li, and Wen Zhao. 2023. Improving News Recommendation with Channel-Wise Dynamic Representations and Contrastive User Modeling. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 562–570.
- [94] Lei Wang, Ee-Peng Lim, Zhiwei Liu, and Tianxiang Zhao. 2022. Explanation guided contrastive learning for sequential recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 2017–2027.
- [95] Tongzhou Wang and Phillip Isola. 2020. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*. PMLR, 9929–9939.
- [96] Xiao Wang and Guo-Jun Qi. 2023. Contrastive learning with stronger augmentations. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 5 (2023), 5549–5560.
- [97] Yu Wang, Hengrui Zhang, Zhiwei Liu, Liangwei Yang, and Philip S Yu. 2022. ContrastVAE: Contrastive Variational AutoEncoder for Sequential Recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 2056–2066.
- [98] Yuening Wang, Yingxue Zhang, and Mark Coates. 2021. Graph Structure Aware Contrastive Knowledge Distillation for Incremental Learning in Recommender Systems. In *CIKM*. ACM, 3518–3522.
- [99] Ziyang Wang, Huoyu Liu, Wei Wei, Yue Hu, Xian-Ling Mao, Shaojian He, Rui Fang, and Dangyang Chen. 2022. Multi-level Contrastive Learning Framework for Sequential Recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 2098–2107.
- [100] Chunyu Wei, Jian Liang, Di Liu, and Fei Wang. 2022. Contrastive Graph Structure Learning via Information Bottleneck for Recommendation. In *Advances in Neural Information Processing Systems*.
- [101] Wei Wei, Chao Huang, Lianghao Xia, Yong Xu, Jiashu Zhao, and Dawei Yin. 2022. Contrastive Meta Learning with Behavior Multiplicity for Recommendation. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 1120–1128.
- [102] Wei Wei, Chao Huang, Lianghao Xia, and Chuxu Zhang. 2023. Multi-Modal Self-Supervised Learning for Recommendation. (2023), 790–800.
- [103] Yinwei Wei, Xiang Wang, Qi Li, Liqiang Nie, Yan Li, Xuanping Li, and Tat-Seng Chua. 2021. Contrastive learning for cold-start recommendation. In *Proceedings of the 29th ACM International Conference on Multimedia*. 5382–5390.
- [104] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Self-Supervised Graph Learning for Recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 726–735.
- [105] Jiancan Wu, Xiang Wang, Xingyu Gao, Jiawei Chen, Hongcheng Fu, Tianyu Qiu, and Xiangnan He. 2022. On the Effectiveness of Sampled Softmax Loss for Item Recommendation. *CoRR* abs/2201.02327 (2022).

- [106] Lirong Wu, Haitao Lin, Cheng Tan, Zhangyang Gao, and Stan Z Li. 2021. Self-supervised learning on graphs: Contrastive, generative, or predictive. *IEEE Transactions on Knowledge and Data Engineering* (2021).
- [107] Yiqing Wu, Ruobing Xie, Yongchun Zhu, Xiang Ao, Xin Chen, Xu Zhang, Fuzhen Zhuang, Leyu Lin, and Qing He. 2022. Multi-view multi-behavior contrastive learning in recommendation. In *Database Systems for Advanced Applications: 27th International Conference, DASFAA 2022, Virtual Event, April 11–14, 2022, Proceedings, Part II*. Springer, 166–182.
- [108] Ying-Xin Wu, Xiang Wang, An Zhang, Xiangnan He, and Tat-Seng Chua. 2022. Discovering invariant rationales for graph neural networks. (2022).
- [109] Jun Xia, Lirong Wu, Jintao Chen, Bozhen Hu, and Stan Z Li. 2022. Simgrace: A simple framework for graph contrastive learning without data augmentation. In *Proceedings of the ACM Web Conference 2022*. 1070–1079.
- [110] Lianghao Xia, Chao Huang, Jiao Shi, and Yong Xu. 2023. Graph-less Collaborative Filtering. In *Proceedings of the ACM Web Conference 2023, WWW 2023*. ACM, 17–27.
- [111] Lianghao Xia, Chao Huang, Yong Xu, Jiashu Zhao, Dawei Yin, and Jimmy X. Huang. 2022. Hypergraph Contrastive Collaborative Filtering. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 70–79.
- [112] Xin Xia, Hongzhi Yin, Junliang Yu, Yingxia Shao, and Lizhen Cui. 2021. Self-supervised graph co-training for session-based recommendation. In *Proceedings of the 30th ACM International conference on information & knowledge management*. 2180–2190.
- [113] Xin Xia, Hongzhi Yin, Junliang Yu, Qinyong Wang, Lizhen Cui, and Xiangliang Zhang. 2021. Self-supervised hypergraph convolutional networks for session-based recommendation. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 35. 4503–4511.
- [114] Xin Xia, Hongzhi Yin, Junliang Yu, Qinyong Wang, Guandong Xu, and Quoc Viet Hung Nguyen. 2022. On-Device Next-Item Recommendation with Self-Supervised Knowledge Distillation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 546–555.
- [115] Ruobing Xie, Qi Liu, Liangdong Wang, Shukai Liu, Bo Zhang, and Leyu Lin. 2022. Contrastive Cross-domain Recommendation in Matching. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4226–4236.
- [116] Xu Xie, Fei Sun, Zhaoyang Liu, Shiwen Wu, Jinyang Gao, Jiandong Zhang, Bolin Ding, and Bin Cui. 2022. Contrastive Learning for Sequential Recommendation. In *38th IEEE International Conference on Data Engineering, ICDE 2022, Kuala Lumpur, Malaysia, May 9–12, 2022*. 1259–1273.
- [117] Yaochen Xie, Zhao Xu, Jingtun Zhang, Zhengyang Wang, and Shuiwang Ji. 2023. Self-supervised learning of graph neural networks: A unified review. *IEEE transactions on pattern analysis and machine intelligence* 45, 2 (2023), 2412–2429.
- [118] Caiming Xiong, Victor Zhong, and Richard Socher. 2017. Dynamic Coattention Networks For Question Answering. In *ICLR (Poster)*.
- [119] Yilun Xu, Shengjia Zhao, Jiaming Song, Russell Stewart, and Stefano Ermon. 2020. A theory of usable information under computational constraints. (2020).
- [120] Hongrui Xuan, Yi Liu, Bohan Li, and Hongzhi Yin. 2023. Knowledge Enhancement for Contrastive Multi-Behavior Recommendation. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 195–203.
- [121] Haoran Yang, Hongxu Chen, Lin Li, S Yu Philip, and Guandong Xu. 2021. Hyper meta-path contrastive learning for multi-behavior recommendation. In *2021 IEEE International Conference on Data Mining (ICDM)*. IEEE, 787–796.
- [122] Yuhao Yang, Chao Huang, Lianghao Xia, and Chunzhen Huang. 2023. Knowledge Graph Self-Supervised Rationalization for Recommendation. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2023*. ACM, 3046–3056.
- [123] Yuhao Yang, Chao Huang, Lianghao Xia, Chunzhen Huang, Da Luo, and Kangyi Lin. 2023. Debiased Contrastive Learning for Sequential Recommendation. In *WWW*. ACM, 1063–1073.
- [124] Yuhao Yang, Chao Huang, Lianghao Xia, and Chenliang Li. 2022. Knowledge Graph Contrastive Learning for Recommendation. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11 - 15, 2022*. ACM, 1434–1443.
- [125] Yonghui Yang, Le Wu, Richang Hong, Kun Zhang, and Meng Wang. 2021. Enhanced graph learning for collaborative filtering via mutual information maximization. In *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 71–80.
- [126] Yonghui Yang, Zhengwei Wu, Le Wu, Kun Zhang, Richang Hong, Zhiqiang Zhang, Jun Zhou, and Meng Wang. 2023. Generative-Contrastive Graph Learning for Recommendation. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2023*. ACM, 1117–1126.
- [127] Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. Graph Convolutional Networks for Text Classification. In *Proceedings of the AAAI conference on artificial intelligence*. 7370–7377.

- [128] Quanming Yao, Mengshuo Wang, Yuqiang Chen, Wenyuan Dai, Yu-Feng Li, Wei-Wei Tu, Qiang Yang, and Yang Yu. 2018. Taking human out of learning applications: A survey on automated machine learning. *arXiv preprint arXiv:1810.13306* (2018).
- [129] Tiansheng Yao, Xinyang Yi, Derek Zhiyuan Cheng, Felix X. Yu, Ting Chen, Aditya Krishna Menon, Lichan Hong, Ed H. Chi, Steve Tjoa, Jieqi (Jay) Kang, and Evan Ettinger. 2021. Self-supervised Learning for Large-scale Item Recommendations. In *CIKM '21: The 30th ACM International Conference on Information and Knowledge Management, Virtual Event, Queensland, Australia, November 1 - 5, 2021*. 4321–4330.
- [130] Haibo Ye, Xinjie Li, Yuan Yao, and Hanghang Tong. 2023. Towards Robust Neural Graph Collaborative Filtering via Structure Denoising and Embedding Perturbation. *ACM Trans. Inf. Syst.* 41, 3 (2023), 28 pages.
- [131] Junliang Yu, Xin Xia, Tong Chen, Lizhen Cui, Nguyen Quoc Viet Hung, and Hongzhi Yin. 2023. XSimGCL: Towards extremely simple graph contrastive learning for recommendation. *IEEE Transactions on Knowledge and Data Engineering* (2023).
- [132] Junliang Yu, Hongzhi Yin, Min Gao, Xin Xia, Xiangliang Zhang, and Nguyen Quoc Viet Hung. 2021. Socially-aware self-supervised tri-training for recommendation. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 2084–2092.
- [133] Junliang Yu, Hongzhi Yin, Jundong Li, Qinyong Wang, Nguyen Quoc Viet Hung, and Xiangliang Zhang. 2021. Self-Supervised Multi-Channel Hypergraph Convolutional Network for Social Recommendation. In *Proceedings of the Web Conference 2021*. 413–424.
- [134] Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Lizhen Cui, and Quoc Viet Hung Nguyen. 2022. Are Graph Augmentations Necessary?: Simple Graph Contrastive Learning for Recommendation. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1294–1303.
- [135] Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Jundong Li, and Zi Huang. 2023. Self-supervised learning for recommender systems: A survey. *IEEE Transactions on Knowledge and Data Engineering* (2023).
- [136] Xu Yuan, Hongshen Chen, Yonghao Song, Xiaofang Zhao, and Zhuoye Ding. 2021. Improving Sequential Recommendation Consistency with Self-Supervised Imitation. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021*. 3321–3327.
- [137] An Zhang, Wenchang Ma, Xiang Wang, and Tat-Seng Chua. 2022. Incorporating Bias-aware Margins into Contrastive Loss for Collaborative Filtering. In *NeurIPS*.
- [138] Junwei Zhang, Min Gao, Junliang Yu, Lei Guo, Jundong Li, and Hongzhi Yin. 2021. Double-scale self-supervised hypergraph learning for group recommendation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 2557–2567.
- [139] Lingzi Zhang, Yong Liu, Xin Zhou, Chunyan Miao, Guoxin Wang, and Haihong Tang. 2022. Diffusion-based graph contrastive learning for recommendation with implicit feedback. In *Database Systems for Advanced Applications: 27th International Conference, DASFAA 2022*. Springer, 232–247.
- [140] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. 2019. Deep learning based recommender system: A survey and new perspectives. *ACM computing surveys (CSUR)* 52, 1 (2019), 1–38.
- [141] Yixin Zhang, Yong Liu, Yonghui Xu, Hao Xiong, Chenyi Lei, Wei He, Lizhen Cui, and Chunyan Miao. 2022. Enhancing Sequential Recommendation with Graph Contrastive Learning. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022*. 2398–2405.
- [142] Yuanliang Zhang, Xiaofeng Wang, Jinxin Hu, Ke Gao, Chenyi Lei, and Fei Fang. 2022. Scenario-Adaptive and Self-Supervised Model for Multi-Scenario Personalized Recommendation. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 3674–3683.
- [143] Kun Zhou, Hui Wang, Wayne Xin Zhao, Yutao Zhu, Sirui Wang, Fuzheng Zhang, Zhongyuan Wang, and Ji-Rong Wen. 2020. S3-Rec: Self-Supervised Learning for Sequential Recommendation with Mutual Information Maximization. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 1893–1902.
- [144] Yuanhang Zhou, Kun Zhou, Wayne Xin Zhao, Cheng Wang, Peng Jiang, and He Hu. 2022. C²-CRS: Coarse-to-Fine Contrastive Learning for Conversational Recommender System. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 1488–1496.
- [145] Zhi-Hua Zhou and Ming Li. 2005. Tri-training: Exploiting unlabeled data using three classifiers. *IEEE Transactions on knowledge and Data Engineering* 17, 11 (2005), 1529–1541.
- [146] Ding Zou, Wei Wei, Xian-Ling Mao, Ziyang Wang, Minghui Qiu, Feida Zhu, and Xin Cao. 2022. Multi-level Cross-view Contrastive Learning for Knowledge-aware Recommender System. In *SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 1358–1368.
- [147] Ding Zou, Wei Wei, Ziyang Wang, Xian-Ling Mao, Feida Zhu, Rui Fang, and Danyang Chen. 2022. Improving knowledge-aware recommendation with multi-level interactive contrastive learning. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. 2817–2826.