
CLICK-THROUGH RATE PREDICTION WITH AUTO-QUANTIZED CONTRASTIVE LEARNING

Yuji Pan*

Shanghai Jiao Tong University
yujiepan@sjtu.edu.cn

Jiangchao Yao

DAMO Academy, Alibaba Group
jiangchao.yjc@alibaba-inc.com

Bo Han

Hong Kong Baptist University
bhanml@comp.hkbu.edu.hk

Kunyang Jia

DAMO Academy, Alibaba Group
kunyang.jky@alibaba-inc.com

Ya Zhang

Shanghai Jiao Tong University
ya_zhang@sjtu.edu.cn

Hongxia Yang

DAMO Academy, Alibaba Group
yang.yhx@alibaba-inc.com

ABSTRACT

Click-through rate (CTR) prediction becomes indispensable in ubiquitous web recommendation applications. Nevertheless, the current methods are struggling under the cold-start scenarios where the user interactions are extremely sparse. We consider this problem as an automatic identification about whether the user behaviors are rich enough to capture the interests for prediction, and propose an Auto-Quantized Contrastive Learning (AQCL) loss to regularize the model. Different from previous methods, AQCL explores both the instance-instance and the instance-cluster similarity to robustify the latent representation, and automatically reduces the information loss to the active users due to the quantization. The proposed framework is agnostic to different model architectures and can be trained in an end-to-end fashion. Extensive results show that it consistently improves the current state-of-the-art CTR models.

1 Introduction

Click-Through Rate (CTR) prediction in the recommender systems is indispensable, which helps rank the candidate items meticulously based on the user interests. In the past few years, several Deep-Learning-based methods *e.g.*, Wide&Deep [11], DeepFM [20] and DIN [67], have achieved impressive performance on this task. Nevertheless, these CTR models still suffer from the cold-start problem that breaks the assumption about the sufficient training samples available. Under the restricted user interactions in the cold-start scenarios, the model performance is dramatically limited [23]. This motivates a range of subsequent works [29, 34, 51, 69] to consider the cold-start recommendation.

One line of methods leverage the implicit regularization to prevent the CTR models from over-fitting [10]. The popular techniques *e.g.*, Dropout [48] and Early-stop [44], would be considered during the training. For example, DropoutNet [51] randomly disturbs the embedding of users or items to robustify the optimization procedure. Some other works explore to use the parameter or the embedding initialization to regularize the training of the CTR models [24]. For example, MeLU [29] learns to initialize the whole parameters of the models using MAML [18]. MAMO [17] extends MELU with two groups of memories to enhance the personalized initialization. MetaEmb [43], MWUF [69] and GME [40] explore to use side information, *e.g.*, the item attributes and user neighbors, to initialize the user and item id embedding. For the recommendation model itself, the supervision for each sample is not changed during training.

Another line of works explicitly construct the auxiliary task to help the training of the CTR models. For example, DeepMCP [41] uses the additional tasks that model the user-item and item-item relationship to improve the user embedding and the item embedding. DIEN [66] and DMR [35] encourage the historical state representation close to the next click item, to better capture the evolution of user interests. SSL4Rec [61] explores to use the unsupervised learning task, SimCLR [7], to enhance the generalization performance of the representation in the models. Similarly,

*Work done when Pan was an intern at Alibaba Group.

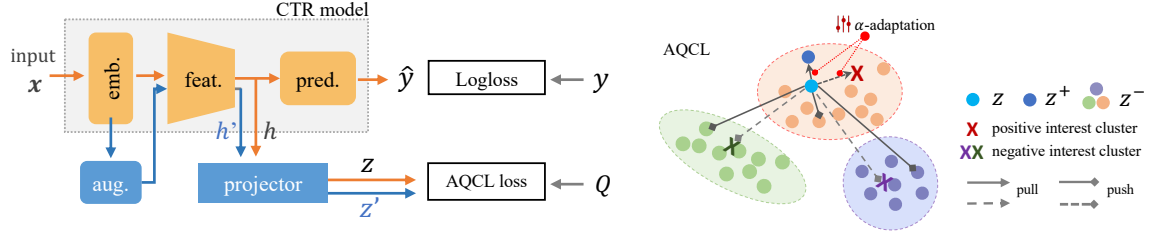


Figure 1: **Left:** Framework of Auto-quantized Contrastive Learning (AQCL) for the CTR model. For each input x , the embedding layer and feature interaction layer convert x into latent code h . The primary prediction task is guided by Logloss, using the output \hat{y} and ground truth label y . Besides, a simple MLP projector $g(\cdot)$ outputs the representation $z = g(h)$. An augmented input is also fed similarly to get z' . The proposed AQCL loss uses the representations z, z' and the interest cluster Q . During training, both Logloss and AQCL Loss are applied. **Right:** Motivation of AQCL. It achieves: (1) instance-instance similarity, where the representation z should be close to z^+ from an augmented input and far from z^- from negative sampling; and (2) interest clusters support, where z should be close to the positive interest cluster and far from negative interest clusters among Q ; (3) automatic balance of the instance-instance and the instance-cluster similarities for non-active/active users by α -adaptation.

CLCRec [56] implements this goal by maximizing the mutual dependencies between item content and the collaborative signals. Empirically, these methods essentially incorporate the prior knowledge on the feature representation to help the training of the CTR models, which have achieved the state-of-the-art performance.

This work follows the second line and explores to leverage the structure of the representation space to automatically constrain the similarity among representations from different users. Specifically, we design an Auto-Quantized Contrastive Learning (AQCL) loss to regularize the training of the CTR models. Unlike the traditional contrastive learning approaches [7–9, 21] and some attempts on the recommendation tasks [56, 58, 61, 68] that focus only on instance-level discrimination, AQCL encourages both the instance-instance similarity and the instance-cluster similarity to automatically contribute to the modeling of the user interests. Figure 1 illustrates the framework and the intuition of AQCL. We conduct a range of experiments on three sparse datasets and the results show AQCL consistently improves the CTR performance under the cold-start scenarios.

In total, our contributions can be summarized as follows:

- We introduce an auxiliary AQCL loss that automatically leverages the instance-instance similarity and the instance-cluster similarity to regularize the representations in the CTR models under the cold-start scenarios.
- The interest clusters used in AQCL are learned together with the loss of the primary task in an end-to-end manner. Simultaneously, an α -adaptation strategy is searched to automatically control the geometric balance for the representation of the active users and the non-active users.
- Extensive experiments on three datasets prove the effectiveness of our proposed method. Besides, AQCL is compatible with the common DNN-based CTR models like W&D, DeepFM and DIN, which improves the ranking performance on both non-active users and active users.

2 Related Works

2.1 CTR in Cold-Start Recommendations

Recommender systems have been well studied in the past decades [4, 11, 13, 15, 23, 26, 32, 46, 49, 60, 62, 64], while the cold-start problem is a long-standing challenge in recommendation tasks. As mentioned before, many previous works can be considered as implicit or explicit regularization on the model. The former tries to interfere with the optimization without extra tasks. For example, DropoutNet [51] applies a data-augmentation on input to encourage robust user or item representations. Many meta-learning based works, *e.g.*, MeLU [29], MetaEmb [43], MAMO [17], MetaHIN [34], PAML [53], GME [40] and MWUF [69], explore to initialize model parameters or embeddings with user and item side information. Some other works train the recommendation model with auxiliary task as explicit regularization. DeepMCP [41] is an early attempt to explore representation learning by designing matching subnet and correlation subnet. SSL4Rec [61] and CLCRec [56] applies traditional contrastive learning loss on either users or item representations. Our method also introduces an auxiliary task, but more comprehensively explores the representation space.

2.2 Contrastive Learning

Self-supervised Learning (SSL) is an unsupervised approach in learning data representations [33] and has shown success in computer vision [30, 52], audio [3, 45], natural language processing [16, 28] and many cross-modality tasks [1, 42, 63]. Contrastive learning (CL) is one representative line of works, including CPC [39], MoCo [21], SimCLR [7] and PIRL [38]. CL maximizes a lower bound on the mutual information between two or among more “views” of an instance [57]. By identifying the positive sample pairs among other negative pairs, it succeeds in capturing the intrinsic features from individual instances in the latent space. Several attempts have used contrastive learning in sequential recommendations to learn either better item-level features [56, 61, 68] or user representations [58] individually, while our method considers the composed representations of the user and the item. Some works have also extended the traditional contrastive learning with more positive pairs. For example, SupCon [27] utilizes extra labels and make each instance close to others with the same class. PCL [31] uses EM to conduct unsupervised clustering and contrastive learning together. Similar to our AQCL method, these works also explore instance representations with other neighbors or clusters. However, it ignores the negative effect on the representation of active users who actually require the sufficient details in representation regarding the recommendation.

3 Preliminary

3.1 CTR Prediction

The CTR prediction as a binary classification problem is to find a map $f(\mathbf{x}_j) \rightarrow y_j$ for each pair $(\mathbf{x}_j, y_j) \in \mathcal{D}$. Generally, for each input \mathbf{x}_j , it at least contains the user id u_j and the candidate item id i_j . Besides, the user historical clicks $\mathbf{s}_j = [i_{j,1}, i_{j,2}, \dots, i_{j,L_j}]$ are often considered, where L_j is the length of the click sequence. Combining the user id u_j , the item id i_j , the historical clicks \mathbf{s}_j and the other features o_j , we have $\mathbf{x}_j = (u_j, i_j, \mathbf{s}_j, o_j)$. The corresponding target label is a binary scalar $y_j \in \{0, 1\}$ meaning whether the user u_j clicks on the candidate item i_j . A typical deep CTR model f consists of the following parts [65]:

- *Embedding layer.* It transforms the sparse categorical features into dense-valued vectors *i.e.*, embedding. Features like the item id are projected as the fixed-length embedding. For the historical sequence, we correspondingly acquire a sequence of embedding for the interacted items.
- *Feature interaction layer.* The transformed embeddings are then fed into the interaction layers to produce a compact representation \mathbf{h}_j for input instance \mathbf{x}_j . This component has diverse designs such as Multi-Layer Perception (MLP) [20], Cross Network [54] and Multi-Head Self-Attention [47].
- *Prediction layer.* Finally, a simple prediction layer (usually a logistic regression module) produces the final score $\hat{y}_j = f(\mathbf{x}_j) \in [0, 1]$ for \mathbf{x}_j on the representation \mathbf{h}_j .

With model output $f(\mathbf{x}_j)$ and ground truth y_j , the CTR model is trained on dataset \mathcal{D} with the Logloss \mathcal{L}_C :

$$\mathcal{L}_C = -\frac{1}{|\mathcal{D}|} \sum_{j=1}^{|\mathcal{D}|} (y_j \log f(\mathbf{x}_j) + (1 - y_j) \log (1 - f(\mathbf{x}_j))). \quad (1)$$

3.2 Cold-start Problem in CTR

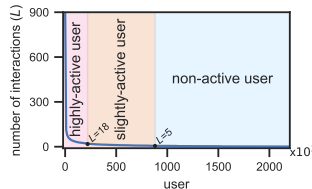


Figure 2: The user interactions (descending sort) on a cold-start industrial scenarios.

Users diverse a lot when it comes to the activeness in the cold-start scenarios¹. We divide users into three groups, non-active, slightly-active and highly-active users based on the length of user click sequence \mathbf{s}_j . Figure 2 plots the

¹The cold-start scenarios here we mean are the early stage of the recommendation feed applications that have many slightly-active or non-active users but also have a small fraction of active users.

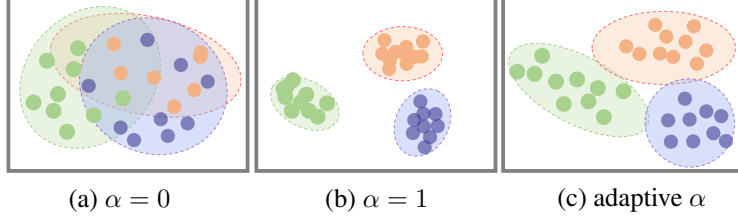


Figure 3: Illustration of AQCL loss with different α in latent space. With $\alpha = 0$, AQCL approximately degrades to ICL, which imposes each sample to be different from others. Ideally the representations are evenly distributed [55]. With $\alpha = 1$, AQCL acts like quantized contrastive learning (QCL), which focuses on the interest clustering to help non-active users. The representations might lose the rich details for the CTR primary task. With an adaptive control of α , AQCL combines the advantages of both ICL and QCL. It encourages the sample representations to be assigned to certain interest clusters, while the distances between samples are still guaranteed to maintain enough information for the CTR prediction.

curve of the user sample number on one cold-start industrial dataset, and categorizes three groups roughly of 60%, 30% and 10% users. We can see that there are only a few of active users and a large proportion of users produces very limited interactions, making the CTR prediction task challenging.

4 Auto-Quantized Contrastive Learning

4.1 Self-supervision Framework for CTR

In the cold-start scenarios, the click signal is usually scarce for training. In this case, the recent self-supervised learning (SSL), *e.g.*, the contrastive learning loss SimCLR [7], is a straightforward choice as an auxiliary task to regularize the CTR model. In SimCLR, for each training instance \mathbf{x} and its randomly augmented version \mathbf{x}^+ , they go through the same feature interaction layer to get corresponding representations \mathbf{z} and \mathbf{z}^+ after projector $g(\cdot)$. The classical loss focuses on the instance-level contrastive learning (ICL) to maximize the similarity between \mathbf{z} and \mathbf{z}^+ ,

$$\mathcal{L}_{\text{ICL}}(\mathbf{z}) = -\log \frac{\exp(\text{sim}(\mathbf{z}, \mathbf{z}^+) / \tau)}{\sum_{\mathbf{z}' \in \{\mathbf{z}^+\} \cup \mathbf{Z}^-} \exp(\text{sim}(\mathbf{z}, \mathbf{z}') / \tau)}, \quad (2)$$

where sim means cosine similarity, τ is a temperature hyper-parameter and \mathbf{Z}^- are samples from negative sampling. Regarding the CTR task, we applying the following input augmentation: first, we randomly mask some items from the user click history \mathbf{s}_j ; and then we randomly set some embedding bits into zero, like dropout operation. Similar to [7], we transform the latent code \mathbf{h} by a small non-linear projection module $g(\cdot)$ to get the actual representation \mathbf{z} for SSL task, *i.e.*, $\mathbf{z}_i = g(\mathbf{h}_i)$. In practice, $g(\cdot)$ is implemented as a 3-layer MLP with the leaky-ReLU [37] activation. The final training loss in the framework is a combination of both the primary CTR prediction task and the self-supervised auxiliary task. Note that, the SSL module only participates in the training stage.

4.2 Auto-Quantized Contrastive Learning

However, ICL loss only explores the instance-level similarity and fails to capture the relationship between neighbors. For non-active users, it is reasonable to get benefits from the neighbors with the rich behaviors in the latent space. This motivates us to model the structure information in the latent representation space, which can be also considered as the user interest clusters. Here, we define the codeword as interest in the representation codebook, borrowed from the concept of vector quantization [19], and propose an Auto-Quantized counterpart of contrastive learning. Figure 3 shows the difference between ICL and AQCL. Formally, for T interests $Q = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_T]$ and a certain \mathbf{z} , we find the top- K closest codewords Q^+ , *i.e.*,

$$Q^+ = \arg \max_{\mathbf{q}^1, \mathbf{q}^2, \dots, \mathbf{q}^K \in Q} \sum_{k=1}^K \text{sim}(\mathbf{z}, \mathbf{q}^k).$$

Note that, sim is the cosine similarity in unit-sphere space as in Eqn. (2), since it is much less vulnerable to mode collapse [36] and is widely used in prototype-based methods [31]. The proposed Auto-Quantized contrastive learning

loss is formulated as a dynamic combination of instance-level and cluster-level contrastive learning:

$$\begin{aligned} \mathcal{L}_{\text{AQCL}}(\mathbf{z}) &= -\log \frac{[d_1(\mathbf{z}, \mathbf{z}^+)]^{1-\alpha} \left[\sum_{\mathbf{q}^+ \in Q^+} d_2(\mathbf{z}, \mathbf{q}^+) \right]^\alpha}{\sum_{\mathbf{z}' \in \{\mathbf{z}^+\} \cup \mathbf{Z}^-} d_1(\mathbf{z}, \mathbf{z}') + \sum_{\mathbf{q}' \in Q} d_2(\mathbf{z}, \mathbf{q}')}}, \\ d_1(\mathbf{z}, \mathbf{z}') &= \exp\left(\text{sim}(\mathbf{z}, \mathbf{z}') / \tau_1\right), \\ d_2(\mathbf{z}, \mathbf{q}') &= \exp\left(\text{sim}(\mathbf{z}, \mathbf{q}') / \tau_2\right), \end{aligned} \quad (3)$$

where τ_1 and τ_2 are temperature hyper-parameters for instances and clusters, and α controls the geometric mean of instance-instance and instance-cluster similarities. The design of α is for automatic balance for the representation personalization of non-active users and active users, which will be explained later. AQCL extends ICL by introducing positive support from both its augmented version \mathbf{z}' and a set of K closest codewords. We allow K to be equal or larger than 1, because a sample representation may contain several interests, and using multiple codewords can be more stable to the probably incomplete interest clustering. Note that, the codewords are built based on the whole dataset, while the positive and negative pairs are from the same batch.

4.2.1 Building the codebook Q

In AQCL, a good codebook should try to cover all the sample representations with relatively small distances. Considering the possibly large-scale training dataset, we leverage an online method [5] to learn the codebook along with AQCL training. In detail, for a codebook $Q = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_T]$ and a batch of representation $Z = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_B]$ with batch size B , we would acquire the corresponding assignment code matrix $A = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_B] \in \mathbb{R}_+^{T \times B}$. Each column \mathbf{a}_b of A denotes the probability of assigning \mathbf{z}_b into the totally T codewords. Similar to [2, 5], the objective is

$$\max_{A \in \mathcal{A}} \text{Tr}(A^\top Q^\top Z) + \epsilon H(A), \quad (4)$$

where H is the entropy function serving as a regularization with a small weight ϵ . We define the constraints of A by

$$\mathcal{A} = \left\{ A \in \mathbb{R}_+^{T \times B} \mid A \mathbf{1}_B = \frac{1}{T} \mathbf{1}_T, Q^\top \mathbf{1}_T = \frac{1}{B} \mathbf{1}_B \right\},$$

where $\mathbf{1}_T$ denotes the vector of ones in dimension T . The constraint ensures that each codeword is roughly assigned evenly. This can be considered as an optimal transport problem [2] and solved by the iterative Sinkhorn-Knopp [14] algorithm with the small computation cost. Then, we convert the continuous solution A^* into its discrete, one-hot version using $\arg \max$ operation. For each representation \mathbf{z}_b , we encourage it to be close to only one of interest codewords. In summary, the loss to build the codebook is transformed as

$$\begin{aligned} \mathcal{L}_{\text{codebook}} &= -\sum_{b=1}^B \sum_{t=1}^T \mathbf{a}_b^{(t)} \log \mathbf{p}_b^{(t)}, \\ \mathbf{p}_b^{(l)} &= \frac{\exp(\text{sim}(\text{sg}(\mathbf{z}_b), \mathbf{q}_l) / \tau_3)}{\sum_{\mathbf{q} \in Q} \exp(\text{sim}(\text{sg}(\mathbf{z}_b), \mathbf{q}) / \tau_3)}, \end{aligned} \quad (5)$$

where $\text{sg}(\cdot)$ means the stop-gradient operation, and τ_3 is a temperature hyper-parameter. Note that, Eqn. (5) is only to learn the codebook and not to update the model parameters.

4.2.2 Auto-Quantization via α -adaptation

Intuitively, users with the scarce clicks are more uncertain and need support from their neighbors, and conversely, the performance of active users might suffer from over-quantization, since their representation maintains more details for prediction. Therefore, the representation learning should account for the different user activeness in the cold-start scenarios. To achieve this, we search α to automatically balance the importance of instance-instance measure and instance-cluster measure in Eqn. (3). Specially, we design a weight control module for α as a function of variable L_j for user u_j , i.e., $\alpha_j = R(L_j)$. With L_j increasing, we know better about the user history, and empirically we shall not enforce conformity and make α smaller. However, it is challenging and labor-consuming to design the weight vs. activeness curve for each agnostic cold-start scenario. Therefore, we resort to AutoML [25] to search the appropriate function $R(L_j)$ for Eqn. (3) in the following.

Algorithm 1 Algorithm for AQCL training

Input: Training samples $\{\mathbf{x}_j\}$ and their history length $\{L_j\}$,
parameters w_1 and w_2 for α -adaptation

Output: the CTR model

- 1: Initiate CTR model f and user interest codebook Q .
- 2: **while** not early stop **do**
- 3: Fetch a batch of $\{\mathbf{x}_j\}$ and the click length $\{L_j\}$.
- 4: Get output $\{\hat{y}_j\}$ and projected representation $\{\mathbf{z}_j\}$.
- 5: Get top- K positive interests with codebook Q .
- 6: Update f by Eqn. (8) with α computed by Eqn. (6).
- 7: Get the discrete assignment matrix A with Sinkhorn and update the user interest codebook Q by Eqn. (5).
- 8: **end while**
- 9: **return** model f

Search space. First, the search space about $R(L_j)$ should take the following two intuitions into account: (1) when L_j increases, α should be reduced; (2) α value should be in the range $[0, 1]$. In this paper, we design the search space as

$$\mathcal{F} = \left\{ R(L_j) = e^{-w_1 \cdot (L_j/L)^{w_2}} : w_1 > 0, w_2 > 0 \right\}, \quad (6)$$

where L is the mean length of the click history for all users. The exact choice of basis function is not important. Figure 4 illustrates some possible search results.

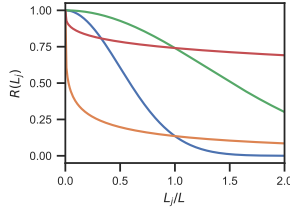


Figure 4: Possible search results for $R(L_j)$.

Search objective. For the problem in this section, we need a subset \mathcal{D}_{val} (partitioned from the training set) to help auto-search. Given θ as the parameters of CTR model f , we target to search for the proper α function such that the model trained on training set $\mathcal{D}_{\text{train}}$ has the best performance on validation set \mathcal{D}_{val} . Concretely, the objective is defined as

$$\begin{aligned} \{w_1^*, w_2^*\} &= \arg \min_{R(\cdot) \in \mathcal{F}} \mathcal{L}_{\text{val}}(f(\theta^*; R), \mathcal{D}_{\text{val}}), \\ \theta^* &= \arg \min_{\theta} \mathcal{L}_{\text{val}}(f(\theta; R), \mathcal{D}_{\text{val}}). \end{aligned} \quad (7)$$

where \mathcal{L}_{val} is the Logloss \mathcal{L}_c on validation set \mathcal{D}_{val} . Based on Eqn. (7), we can find the optimal function in Eqn. (6), which achieves the dynamic balance in representation learning.

4.3 AQCL Algorithm

The AQCL is implemented as an auxiliary loss to the primary CTR task. Therefore, the overall loss is

$$\mathcal{L} = \mathcal{L}_c + w\mathcal{L}_{\text{AQCL}}, \quad (8)$$

where w is the weight for the auxiliary task. Note that, the codebook Q is learned together with Eqn. (8) in an end-to-end manner, using the loss function in Eqn. (5). The training procedure is summarized in Algorithm 1. AutoML eases the search of hyper-parameters w_1, w_2 in Algorithm 1 by using the objective (7). Once the proper hyper-parameters are found by AutoML, we get the final CTR model. During the test phase, all components in the auxiliary task are omitted.

5 Experiments

In this section, we will evaluate the proposed AQCL framework. Specifically, we would like to answer the questions:

- **RQ1.** Compared with other methods, how does AQCL perform with the CTR model on different group of users?

- **RQ2.** Is AQCL as an auxiliary task generally compatible with the different CTR models?
- **RQ3.** How does AQCL work and how do the hyperparameters make the effect?

5.1 Datasets

We conduct experiments on three datasets with relatively severe data sparsity. The statistics are listed in Table 1.

Table 1: Statistics for experiment datasets.

dataset	#users	#items	#samples
Amazon	22,363	12,101	198,502
Ta Feng	32,266	23,812	817,741
Oncold	908,400	21,078	53,754,238

- **Amazon**². This dataset [22] is composed of product reviews from the Amazon website. We follow [6, 12] and use the subset of Beauty to verify AQCL. The task is defined to predict whether a user will comment about a certain item.
- **Ta Feng**³. This is a sparse grocery shopping dataset released by ACM RecSys. It covers products from food, office supplies to furniture. The dataset consists of user transactions from November 2000 to February 2001. We predict whether a user will buy a certain item.
- **Oncold**. This is an industrial dataset of the real-world online cold-start recommendation feeds, which is collected from May to July, 2021. The dataset is extremely sparse as most of users only have clicked a few of items.

Like [36, 67], we sort the user behaviors by timestamp. For Amazon (Ta Feng, respectively), we use the last interaction (day) as the test, the second last interaction (day) as the validation, and the rest as the training data. For Oncold dataset, we split the data of last 20 days equally as the test set and the validation set, and the rest clicks are used as the training data.

5.2 Experiment Settings

5.2.1 Baselines

For RQ1, we compare with some representative methods of two research lines, *i.e.*, the implicit regularization and the explicit regularization. For fair comparison, the following methods and AQCL all use DIN as the backbone model. Besides, we only refer to the design of regularization and omit their modifications of the CTR model architectures in these works, *e.g.*, the positional encoding.

- **DropoutNet** [51] is a training strategy that randomly masks user or item embeddings to handle the cold-start problem. It encourages the CTR model to make full use of the side information.
- **DeepMCP** [41] is a representation-learning-aided model. Except the Logloss, it uses the matching subnet to capture the user-item relation, and the correlation subnet to explore the item-history relation.
- **DMR** [35] proposes an auxiliary matching loss to measure the correspondence between the user preference and the target item in the embedding space.
- **ICL** [7] is the vanilla instance-level contrastive learning. We here use ICL during the training stage rather than pre-training.

For RQ2, we verify AQCL with the following backbones.

- **W&D** [11] is a classical method that uses the feature-cross to help the model capture the high-order relationship hidden in the data for the better prediction.
- **DeepFM** [20] is a successful attempt to combine the power of factorization machines in recommendation and deep learning in the feature learning.
- **DIN** [67] uses the attention mechanism to learn the representation of the user click history given the candidate item, to better explore the user interests.

²<https://jmcauley.ucsd.edu/data/amazon>

³http://recsyswiki.com/wiki/Grocery_shopping_datasets

Table 2: The average result of 5 trials on Amazon, Ta Feng and Oncold datasets. DIN is the base without regularization.

Dataset	Model	Overall		Non-active user		Slightly-active user		Highly-active user	
		AUC	RelaImpr	AUC	RelaImpr	AUC	RelaImpr	AUC	RelaImpr
Amazon	DIN	0.6956	0.00%	0.6719	0.00%	0.7064	0.00%	0.7998	0.00%
	DropoutNet	0.6929	-1.38%	0.6691	-1.63%	0.7034	-1.45%	0.7975	-0.77%
	DeepMCP	0.7053	4.96%	0.6797	4.54%	0.7191	6.15%	0.8097	3.30%
	DMR	0.6982	1.33%	0.6736	0.99%	0.7087	1.11%	0.8064	2.20%
	ICL	0.7016	3.07%	0.6767	2.79%	0.7138	3.59%	0.8059	2.03%
	AQCL	0.7078	6.24%	0.6826	6.22%	0.7231	8.09%	0.8105	3.57%
Ta Feng	DIN	0.6865	0.00%	0.6783	0.00%	0.6929	0.00%	0.7009	0.00%
	DropoutNet	0.6812	-2.84%	0.6762	-1.18%	0.6843	-4.46%	0.6927	-4.08%
	DeepMCP	0.6881	0.86%	0.6791	0.45%	0.6947	0.93%	0.7057	2.39%
	DMR	0.6880	0.80%	0.6806	1.29%	0.6930	0.05%	0.7031	1.10%
	ICL	0.6914	2.63%	0.6812	1.63%	0.6931	0.10%	0.7032	1.14%
	AQCL	0.6931	3.54%	0.6832	2.75%	0.6969	2.07%	0.7098	4.43%
Oncold	DIN	0.7601	0.00%	0.7361	0.00%	0.7608	0.00%	0.7787	0.00%
	DropoutNet	0.7615	0.54%	0.7387	1.10%	0.7617	0.35%	0.7756	-1.11%
	DeepMCP	0.7675	2.85%	0.7409	2.03%	0.7728	4.60%	0.7914	4.56%
	DMR	0.7598	-0.12%	0.7365	0.17%	0.7570	-1.46%	0.7778	-0.32%
	ICL	0.7640	1.50%	0.7367	0.25%	0.7706	3.76%	0.7895	3.88%
	AQCL	0.7691	3.46%	0.7396	1.48%	0.7823	8.24%	0.7968	6.49%

5.2.2 Implementations

All experiments are implemented in PyTorch and run on NVIDIA Tesla V100 GPUs. We use Adam optimizer and the learning rate is 0.001 during training. The dropout rate for DNNs is set as 0.2. We search the L2 regularization weight for the embedding among $\{10^{-5}, 10^{-4}, \dots, 10^{-1}\}$ in all models. For AQCL, we define the weight w of auxiliary task as $\{0.01, 0.05, 0.1\}$, and the temperatures τ_1, τ_2, τ_3 are all set as 0.1. We set the codebook capacity T as 128 for each dataset, and K as 5.

5.2.3 Evaluation metrics

We adopt AUC and RelaImpr to evaluate the performance of the CTR models. AUC measures the probability that a randomly chosen positive sample is ranked higher than a randomly chosen negative sample. Higher AUC indicates better performance. RelaImpr, used in many works like [59, 69], shows the relative improvement of the target compared with the base. Here, we define the base as the backbone used in our experiments, *e.g.*, DIN. RelaImpr is thus defined as

$$\text{RelaImpr} = \left(\frac{\text{AUC}(\text{target}) - 0.5}{\text{AUC}(\text{base}) - 0.5} - 1 \right) \times 100\%. \quad (9)$$

We calculate the metrics on non-active, slightly-active and active users individually to monitor the model performance on the users with the different activeness.

5.3 Results and Discussion

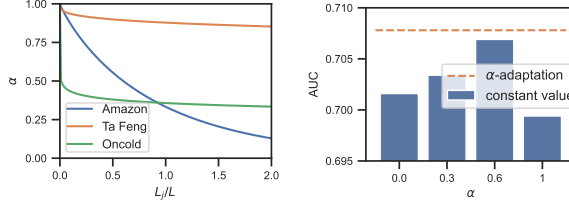
5.3.1 CTR prediction

To answer RQ1, we conduct experiments on the baselines and our method. From Table 2, we can find AQCL achieves consistent improvements on three datasets.

- In general, we observe the overall performance improvement with the implicit/explicit regularization. However, there is a slight performance drop for DropoutNet on Ta Feng and DMR on Oncold. Considering that DropoutNet emphasizes the importance of the side information, the performance may get hurt if there is no enough auxiliary information available. For DMR, the representation correspondence constraint between the user and the target item might be harmful to the active users. For other cases, the positive effects are shown on both non-active users and active users. This means that the proper regularization can help the CTR model in cold-start scenarios.

Table 3: Performance of AQCL with different backbones.

Model	Amazon	Ta Feng	Oncold
W&D	0.6803	0.6716	0.7547
AQCL _{W&D}	0.6841	0.6732	0.7578
DeepFM	0.6810	0.6721	0.7633
AQCL _{DeepFM}	0.6848	0.6768	0.7659
DIN	0.6956	0.6865	0.7601
AQCL _{DIN}	0.7078	0.6931	0.7691

Figure 5: **Left:** The curves of α in AQCL for three datasets; **Right:** The experiments on Amazon with the constant α .

- Our method outperforms the baselines in most cases. The advantage of AQCL is that it can automatically capture the interest clusters to support the non-active users and weaken the information loss for the representation of the active users via α adaptation. For Amazon dataset, we can see the significant improvement on non-active users, which shows the effectiveness by using neighbour information to alleviate the sparsity issue. For Ta Feng and Oncold datasets, there is more gain on the highly-active users. This can be attributed to that the proper relaxation to incorporate the instance-instance similarity in the AQCL loss yields a more robust representation.

5.3.2 Different backbones

To answer RQ2, we adapt AQCL with other backbones, *i.e.*, W&D and DeepFM to verify its effectiveness. For both W&D and DeepFM, we do not modify their wide or FM part, and apply AQCL to the output before the last linear layer on the deep side. Table 3 summarizes the results of AQCL with different backbones. According to the table, AQCL consistently improves the backbones, which demonstrates its compatibility with the popular CTR models.

5.4 Visualization and Ablation Study

To answer RQ3, we conduct a range of visualization and ablation study about AQCL in the following.

5.4.1 α -adaptation

The search result of α in AQCL for each dataset is plot in the left panel of Figure 5. The curvatures are significantly different on three datasets. Amazon requires the relatively small instance-cluster regularization on active users. In comparison, Ta Feng emphasizes more dependency on the clusters and Oncold is similar but gives large α to all users. Besides, we explore to replace α -adaptation with the constant value in $\{0, 0.3, 0.6, 1.0\}$. According to the right panel of Figure 5, there is a performance drop compared with α -adaptation, which demonstrates the effectiveness of AQCL to avoid the human labor.

5.4.2 Representation h

To visualize the learned representation in the latent space, we randomly choose a subset from Oncold dataset, and project the representations $\{h_j\}$ into the 2D space via t-SNE [50] in Figure 6. For a better view, we assign each point with the color representing the interest cluster deduced by AQCL. We find that AQCL can divide the sample representations roughly into several groups, while the vanilla DIN does not. This confirms the motivation of AQCL regarding the interest clusters, which might be useful to the non-active users.

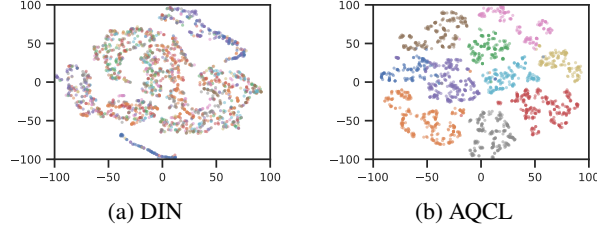


Figure 6: t-SNE visualization of the latent representations respectively learned by DIN and AQCL on Oncold dataset.

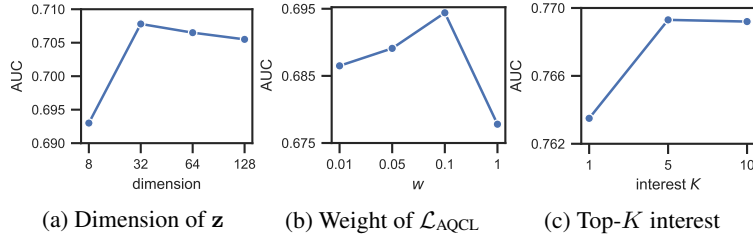


Figure 7: (a) The effect of the dimension of z on Amazon; (b) The effect of the auxiliary weight w on Ta Feng. (c) The effect of the positive interest number K on Oncold.

5.4.3 Dimension of z

As mentioned before, we project the hidden latent code \mathbf{h}_i to another space by MLP $g(\cdot)$. The resulting vector \mathbf{z}_i can have different dimensions. In Figure 7a, we show the dimension of z is also important to the final performance. Specifically, we find that the dimension ≥ 32 can keep the effectiveness of AQCL, while too small dimensions hinder the projector from collecting the enough information for the auxiliary task, and thus causes negative effects.

5.4.4 Weight of \mathcal{L}_{AQCL}

In this part, we conduct experiments with the different auxiliary task weight w on Ta Feng dataset to verify its effect. As shown in Figure 7b, we empirically find that w should not be very large. This might be because AQCL here serves as a data-driven regularization, and too large w may result in the little attention to the primary task. Besides, w should not be too small, since in this case, it will degrade to the vanilla CTR task without the auxiliary gain.

5.4.5 Interest number K

In AQCL, it allows each sample to be assigned into several interest clusters by adjusting K . Figure 7c shows the effect of changing K as $\{1, 5, 10\}$. We see that the performance decreased when $K = 1$. This implies the user representation consists of multiple interests.

6 Conclusion

This paper aims at handling the cold-start scenarios to help the CTR model by designing an auxiliary task. We propose an Auto-quantized Contrastive Learning (AQCL) loss to encourage the model to leverage the possible interest clusters to help the non-active users and maintain the generalization ability to the active users. By training the CTR models with AQCL, we demonstrate our method consistently improve the current models, especially in the face of scarce interactions. The proposed framework is compatible with different model architectures and can be trained in an end-to-end fashion. We hope our work can inspire more explore to improve the CTR models with self-supervised representation learning.

References

- [1] Humam Alwassel, Dhruv Mahajan, Bruno Korbar, Lorenzo Torresani, Bernard Ghanem, and Du Tran. Self-supervised learning by cross-modal audio-video clustering. *Advances in Neural Information Processing Systems*, 33, 2020.

- [2] Yuki Markus Asano, Christian Rupprecht, and Andrea Vedaldi. Self-labelling via simultaneous clustering and representation learning. In *International Conference on Learning Representations (ICLR)*, 2020.
- [3] Alexei Baevski, Steffen Schneider, and Michael Auli. vq-wav2vec: Self-supervised learning of discrete speech representations. In *International Conference on Learning Representations (ICLR)*, 2020.
- [4] Oren Barkan and Noam Koenigstein. Item2vec: neural item embedding for collaborative filtering. In *2016 IEEE 26th International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6. IEEE, 2016.
- [5] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [6] Chong Chen, Min Zhang, Yiqun Liu, and Shaoping Ma. Neural attentional rating regression with review-level explanations. In *Proceedings of the World Wide Web Conference (WWW)*, pages 1583–1592, 2018.
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- [8] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big self-supervised models are strong semi-supervised learners. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pages 22243–22255, 2020.
- [9] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *International Conference on Computer Vision (ICCV)*, 2021.
- [10] Yihong Chen, Bei Chen, Xiangnan He, Chen Gao, Yong Li, Jian-Guang Lou, and Yue Wang. λ opt: Learn to regularize recommender models in finer levels. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, pages 978–986, 2019.
- [11] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on Deep Learning for Recommender Systems*, 2016.
- [12] Zhiyong Cheng, Ying Ding, Xiangnan He, Lei Zhu, Xuemeng Song, and Mohan Kankanhalli. A³nfc: An adaptive aspect attention model for rating prediction. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, pages 3748–3754, 7 2018.
- [13] Paul Covington, Jay Adams, and Emre Sargin. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM conference on recommender systems*, pages 191–198, 2016.
- [14] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems (NeurIPS)*, 26:2292–2300, 2013.
- [15] Mukund Deshpande and George Karypis. Item-based top-n recommendation algorithms. *ACM Transactions on Information Systems (TOIS)*, 22(1):143–177, 2004.
- [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 4171–4186, 2019.
- [17] Manqing Dong, Feng Yuan, Lina Yao, Xiwei Xu, and Liming Zhu. Mamo: Memory-augmented meta-optimization for cold-start recommendation. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, pages 688–697, 2020.
- [18] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1126–1135, 2017.
- [19] Robert Gray. Vector quantization. *IEEE ASSP Magazine*, 1(2):4–29, 1984.
- [20] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. Deepfm: a factorization-machine based neural network for ctr prediction. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, 2017.
- [21] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9729–9738, 2020.
- [22] Ruining He and Julian McAuley. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *Proceedings of the International Conference on World Wide Web (WWW)*, 2016.

- [23] Xinran He, Junfeng Pan, Ou Jin, Tianbing Xu, Bo Liu, Tao Xu, Yanxin Shi, Antoine Atallah, Ralf Herbrich, Stuart Bowers, et al. Practical lessons from predicting clicks on ads at facebook. In *Proceedings of the Eighth International Workshop on Data Mining for Online Advertising*, pages 1–9, 2014.
- [24] T. M. Hospedales, A. Antoniou, P. Micaelli, and A. J. Storkey. Meta-learning in neural networks: A survey. *IEEE Transactions on Pattern Analysis & Machine Intelligence (TPAMI)*, 2021.
- [25] Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren. *Automated machine learning: methods, systems, challenges*. Springer Nature, 2019.
- [26] Wang-Cheng Kang and Julian McAuley. Self-attentive sequential recommendation. In *IEEE International Conference on Data Mining (ICDM)*, pages 197–206. IEEE, 2018.
- [27] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 33, 2020.
- [28] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations (ICLR)*, 2020.
- [29] Hyeop Lee, Jinbae Im, Seongwon Jang, Hyunsouk Cho, and Sehee Chung. Melu: Meta-learned user preference estimator for cold-start recommendation. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, pages 1073–1082, 2019.
- [30] Hongming Li and Yong Fan. Non-rigid image registration using self-supervised fully convolutional networks without training data. In *IEEE International Symposium on Biomedical Imaging (ISBI)*, pages 1075–1078. IEEE, 2018.
- [31] Junnan Li, Pan Zhou, Caiming Xiong, and Steven CH Hoi. Prototypical contrastive learning of unsupervised representations. In *International Conference on Learning Representations (ICLR)*, 2020.
- [32] Greg Linden, Brent Smith, and Jeremy York. Amazon. com recommendations: Item-to-item collaborative filtering. *IEEE Internet computing*, 7(1):76–80, 2003.
- [33] Xiao Liu, Fanjin Zhang, Zhenyu Hou, Zhaoyu Wang, Li Mian, Jing Zhang, and Jie Tang. Self-supervised learning: Generative or contrastive. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 2021.
- [34] Yuanfu Lu, Yuan Fang, and Chuan Shi. Meta-learning on heterogeneous information networks for cold-start recommendation. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, pages 1563–1573, 2020.
- [35] Ze Lyu, Yu Dong, Chengfu Huo, and Weijun Ren. Deep match to rank model for personalized click-through rate prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 34, pages 156–163, 2020.
- [36] Jianxin Ma, Chang Zhou, Peng Cui, Hongxia Yang, and Wenwu Zhu. Learning disentangled representations for recommendation. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 5711–5722, 2019.
- [37] Andrew L Maas, Awni Y Hannun, Andrew Y Ng, et al. Rectifier nonlinearities improve neural network acoustic models. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2013.
- [38] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6707–6717, 2020.
- [39] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. In *International Conference on Learning Representations (ICLR)*, 2019.
- [40] Wentao Ouyang, Xiuwu Zhang, Shukui Ren, Li Li, Kun Zhang, Jinmei Luo, Zhaojie Liu, and Yanlong Du. Learning graph meta embeddings for cold-start ads in click-through rate prediction. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, 2021.
- [41] Wentao Ouyang, Xiuwu Zhang, Shukui Ren, Chao Qi, Zhaojie Liu, and Yanlong Du. Representation learning-assisted click-through rate prediction. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, 2019.
- [42] Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 631–648, 2018.
- [43] Feiyang Pan, Shuokai Li, Xiang Ao, Pingzhong Tang, and Qing He. Warm up cold-start advertisements: Improving ctr predictions via learning to learn id embeddings. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 695–704, 2019.

- [44] Garvesh Raskutti, Martin J Wainwright, and Bin Yu. Early stopping and non-parametric regression: an optimal data-dependent stopping rule. *The Journal of Machine Learning Research*, 15(1):335–366, 2014.
- [45] Mirco Ravanelli, Jianyuan Zhong, Santiago Pascual, Pawel Swietojanski, Joao Monteiro, Jan Trmal, and Yoshua Bengio. Multi-task self-supervised learning for robust speech recognition. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6989–6993. IEEE, 2020.
- [46] Steffen Rendle. Factorization machines. In *Proceedings of the IEEE International Conference on Data Mining (ICDM)*, 2010.
- [47] Weiping Song, Chence Shi, Zhiping Xiao, Zhijian Duan, Yewen Xu, Ming Zhang, and Jian Tang. Autoint: Automatic feature interaction learning via self-attentive neural networks. In *Proceedings of the ACM International Conference on Information and Knowledge Management (CIKM)*, 2019.
- [48] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [49] Qiaoyu Tan, Jianwei Zhang, Jiangchao Yao, Ninghao Liu, Jingren Zhou, Hongxia Yang, and Xia Hu. Sparse-interest network for sequential recommendation. In *Proceedings of the ACM International Conference on Web Search and Data Mining (WSDM)*, pages 598–606, 2021.
- [50] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008.
- [51] Maksims Volkovs, Guang Wei Yu, and Tomi Poutanen. Dropoutnet: Addressing cold start in recommender systems. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 4957–4966, 2017.
- [52] Jiangliu Wang, Jianbo Jiao, Linchao Bao, Shengfeng He, Yunhui Liu, and Wei Liu. Self-supervised spatio-temporal representation learning for videos by predicting motion and appearance statistics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4006–4015, 2019.
- [53] Li Wang, Binbin Jin, Zhenya Huang, Hongke Zhao, Defu Lian, Qi Liu, and Enhong Chen. Preference-adaptive meta-learning for cold-start recommendation. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, pages 1607–1614, 2021.
- [54] Ruoxi Wang, Bin Fu, Gang Fu, and Mingliang Wang. Deep & cross network for ad click predictions. In *Proceedings of the ADKDD’17*, 2017.
- [55] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning (ICML)*, 2020.
- [56] Yinwei Wei, Xiang Wang, Qi Li, Liqiang Nie, Yan Li, Xuanping Li, and Tat-Seng Chua. Contrastive learning for cold-start recommendation. In *The ACM International Conference on Multimedia (MM)*, 2021.
- [57] Mike Wu, Chengxu Zhuang, Milan Mosse, Daniel Yamins, and Noah D. Goodman. On mutual information in contrastive learning for visual representations. *arXiv preprint arXiv:2005.13149*, 2020.
- [58] Xu Xie, Fei Sun, Zhaoyang Liu, Shiwen Wu, Jinyang Gao, Bolin Ding, and Bin Cui. Contrastive learning for sequential recommendation. *arXiv preprint arXiv:2010.14395*, 2020.
- [59] Ling Yan, Wu-Jun Li, Gui-Rong Xue, and Dingyi Han. Coupled group lasso for web-scale ctr prediction in display advertising. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 32, pages 802–810, 2014.
- [60] Jiangchao Yao, Feng Wang, Kunyang Jia, Bo Han, Jingren Zhou, and Hongxia Yang. Device-cloud collaborative learning for recommendation. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, pages 3865–3874, 2021.
- [61] Tiansheng Yao, Xinyang Yi, Derek Zhiyuan Cheng, Felix Yu, Ting Chen, Aditya Menon, Lichan Hong, Ed H Chi, Steve Tjoa, Jieqi Kang, et al. Self-supervised learning for large-scale item recommendations. *arXiv preprint arXiv:2007.12865*, 2020.
- [62] Xinyang Yi, Ji Yang, Lichan Hong, Derek Zhiyuan Cheng, Lukasz Heldt, Aditee Kumthekar, Zhe Zhao, Li Wei, and Ed Chi. Sampling-bias-corrected neural modeling for large corpus item recommendations. In *Proceedings of the 13th ACM Conference on Recommender Systems*, pages 269–277, 2019.
- [63] Shengyu Zhang, Tan Jiang, Tan Wang, Kun Kuang, Zhou Zhao, Jianke Zhu, Jin Yu, Hongxia Yang, and Fei Wu. DevIbert: Learning deconfounded visio-linguistic representations. In *Proceedings of the ACM International Conference on Multimedia (MM)*, pages 4373–4382, 2020.

- [64] Shengyu Zhang, Dong Yao, Zhou Zhao, Tat-Seng Chua, and Fei Wu. Causerec: Counterfactual user sequence synthesis for sequential recommendation. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, 2021.
- [65] Weinan Zhang, Jiarui Qin, Wei Guo, Ruiming Tang, and Xiuqiang He. Deep learning for click-through rate estimation. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, pages 4695–4703, 8 2021. Survey Track.
- [66] Guorui Zhou, Na Mou, Ying Fan, Qi Pi, Weijie Bian, Chang Zhou, Xiaoqiang Zhu, and Kun Gai. Deep interest evolution network for click-through rate prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2019.
- [67] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. Deep interest network for click-through rate prediction. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*, 2018.
- [68] Kun Zhou, Hui Wang, Wayne Xin Zhao, Yutao Zhu, Sirui Wang, Fuzheng Zhang, Zhongyuan Wang, and Ji-Rong Wen. S3-rec: Self-supervised learning for sequential recommendation with mutual information maximization. In *Proceedings of the ACM International Conference on Information & Knowledge Management (CIKM)*, pages 1893–1902, 2020.
- [69] Yongchun Zhu, Ruobing Xie, Fuzhen Zhuang, Kaikai Ge, Ying Sun, Xu Zhang, Leyu Lin, and Juan Cao. Learning to warm up cold item embeddings for cold-start recommendation with meta scaling and shifting networks. In *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, 2021.