

# UniMatch: A Unified User-Item Matching Framework for the Multi-purpose Merchant Marketing

Qifang Zhao, Tianyu Li, Meng Du, Yu Jiang, Qinghui Sun, Zhongyao Wang, Hong Liu, Huan Xu

Alibaba Group, Hangzhou, China

{james.zqf, qianchuan.lty, dmmeng.dm, jiangyu.jiangyu}@alibaba-inc.com

{yuyang.sqh, zhongyao.wangzy, liuhong.liu, huan.xu}@alibaba-inc.com

**Abstract**—When doing private domain marketing with cloud services, the merchants usually have to purchase different machine learning models for the multiple marketing purposes, leading to a very high cost. We present a unified user-item matching framework to simultaneously conduct item recommendation and user targeting with just one model. We empirically demonstrate that the above concurrent modeling is viable via modeling the user-item interaction matrix with the multinomial distribution, and propose a bidirectional bias-corrected NCE loss for the implementation. The proposed loss function guides the model to learn the user-item joint probability  $p(u, i)$  instead of the conditional probability  $p(i|u)$  or  $p(u|i)$  through correcting both the users and items’ biases caused by the in-batch negative sampling. In addition, our framework is model-agnostic enabling a flexible adaptation of different model architectures. Extensive experiments demonstrate that our framework results in significant performance gains in comparison with the state-of-the-art methods, with greatly reduced cost on computing resources and daily maintenance.

**Index Terms**—Marketing, Matching, Item Recommendation, User Targeting, Bias Correction.

## I. INTRODUCTION

Nowadays, merchants commonly sell their products in multiple channels, such as the public platforms like Amazon, Alibaba, and the private channels like their own websites, offline shops, and exclusive customer groups on social medias like Wechat, etc. The marketing on those public platforms, managed by the e-commerce companies, has reached a limit in recent years. As a result, merchants are paying more attention to operate their businesses via the private channels, *i.e.*, conducting the private domain marketing. In order to manage businesses more effectively, merchants utilize the cloud services like Amazon Web Services and Alibaba Cloud, to link all the private channels.

The cloud services not only manage data for merchants, but also provide machine learning techniques for the intelligent marketing. There are two common marketing directions of the merchants: the item recommendation (IR) [33] and the user targeting (UT). To be more specific, merchants try to keep their high-value users active and loyal by periodically sending them messages or emails with recommended items. Meanwhile, merchants always look forward to discovering the potential buyers for certain items, *e.g.*, new releases or popular products, etc. Then, they can send personalized promotion

content to those targeted users. Owing to the machine learning techniques, both item recommendation and user targeting contribute to the profit of merchants significantly.

However, merchants have to purchase a handful of machine learning models for different marketing purposes. First, the item recommendation usually requires one model. Then, the user targeting usually requires more than one model, because practitioners need to create multiple targeting lists according to different promotion subjects, *e.g.*, popular products or bundles of items. It takes great efforts to conduct feature engineering, model training and inference for each model. These practices push up the cost dramatically.

This paper proposes a unified user-item matching framework, named *UniMatch*, to serve for the item recommendation and user targeting with one model only. The previous recommendation algorithms utilize the conditional probability  $p(i|u)$  as the modeling objective [8], [23], while the user targeting models are commonly optimized via the objective  $p(u|i)$ . In our UniMatch framework, the modeling objective is the joint probability  $p(u, i)$ , which is implemented with a bidirectional bias-corrected NCE loss, named *bbcNCE*. When applied for the item recommendation,  $p(u, i) = p(i|u)p(u)$  will produce a similar item list compared to  $p(i|u)$  given a specific user. The same logic holds for the user targeting as well. Thus, our framework is able to reduce the cost of computing resources and data storage, and relieve the burden of the daily maintenance as well.

Different from the online recommendation on the e-commerce platforms, merchants usually apply these intelligent marketing models less frequently when doing private domain marketing. For instance, they send promotion emails or personalized messages weekly or even longer. To adapt for this specific scenario, both the potential user and recommended item lists are produced under a next- $n$ -day prediction setting. Conventionally, both the item recommendation and user targeting tasks are solved by modeling the Bernoulli or multinomial distribution on the user-item interaction matrix. In this paper, we first theoretically prove that it is equivalent to model with the Bernoulli and multinomial distribution since they converge to the same optima in practice. Then, we uncover that modeling with the multinomial distribution has better efficiency in terms of the data preparation and model

convergence. Therefore, we follow our discovery and propose a bidirectional NCE loss with bias correction to model the user-item joint probability  $p(u, i)$ . Additionally, our framework adopts a classical two-tower architecture which enables a flexible utilization of different models.

Our framework has been implemented in the Alibaba cloud product, *QuickAudience(QA)*<sup>1</sup>, for the intelligent marketing of the merchants. Our contributions are summarized as follows:

- We present a unified user-item matching framework, *UniMatch*, which trains only one model to serve both the item recommendation and user targeting simultaneously. To the best of our knowledge, this is the first work on the topic.
- We theoretically prove the equivalence between modeling the user-item interaction matrix with the Bernoulli and multinomial distributions, and empirically demonstrate that modeling with the multinomial distribution yields more robust results with much less resources.
- We propose a bidirectional bias-corrected NCE loss, *bbcNCE*, and train models with the joint probability of  $u$  and  $i$  being the learning objective in theory. Also, we empirically show that the *bbcNCE* loss will guide the model to learn the joint distribution.
- Extensive experiments on two public datasets and two real-world datasets demonstrate that the proposed framework consistently yields improved performance, in comparison with the state-of-the-art methods on both item recommendation and user targeting tasks. In addition, our framework saves up to 94%+ of the total cost compared to previous practices.

## II. PRELIMINARIES

In this section, we first describe the characteristics of the private domain marketing for merchants, and then introduce how the previous research solves the tasks of IR and UT via modeling the user-item interaction matrix separately with the Bernoulli or multinomial distributions.

### A. Problem Definition

For the private domain marketing, the merchants generally carry out the IR and UT for marketing periodically via their private channels. They send out messages or emails to their active users or potential customers, and then expect them to take actions, *e.g.*, visiting offline shops, making inquiries online, and placing an order, etc, some time later. It relatively takes longer time for the merchants to achieve the private domain marketing results.

To fit for this application scenario, we formulate both the IR and UT as a next- $n$ -day prediction problem. When a user  $u$  purchases an item  $i$  at time  $t$ , a record  $(u, i, t)$  is logged. Given the raw logs  $\{(u, i, t)\}$ , the following data-processing method is applied for the next- $n$ -day prediction: we create a dataset  $\mathcal{D} = \{(x_{u,t}, y_{u,t}) : t \in \{1, 2, \dots, T_u\} | u \in \{1, 2, \dots, N\}\}$ , where  $x_{u,t}$  represents user  $u$ 's purchases prior to day  $t$ , *i.e.*, a

sequence of purchased items, and  $y_{u,t}$  is any item purchased during the next  $n$  days  $[t, t + n)$ , and  $T_u = T$  is the number of days.

The dataset  $\mathcal{D}$  forms a user-item interaction matrix  $S_{ui}$ , as shown in Fig. 1, where the rows and columns represent  $x_{u,t}$  and  $y_{u,t}$ , respectively. We call  $x_{u,t}$  as the pseudo-user, and all possible sequences of purchases form the pseudo-user set. Without loss of generality, we use  $u$  to represent the pseudo-user  $x_{u,t}$ , and  $i$  to represent  $y_{u,t}$  in the rest of the paper.

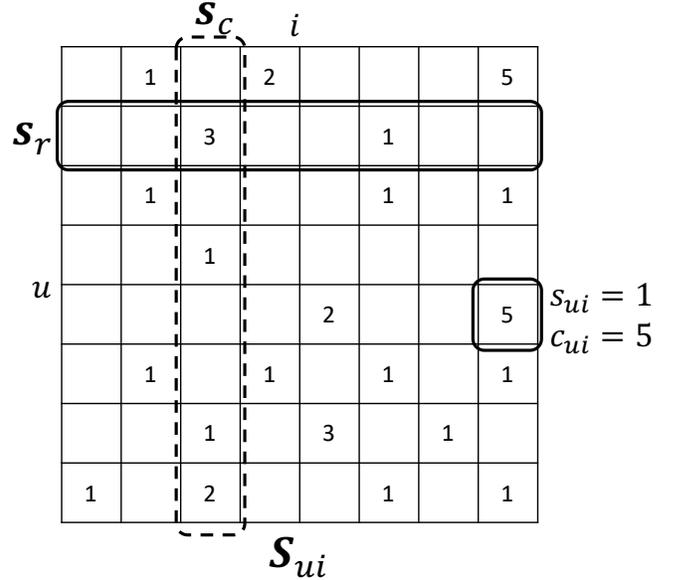


Fig. 1. An illustration of the user-item interaction matrix  $S_{ui}$ . We use  $s \in \{0, 1\}$  to denote whether there is an interaction between a user and an item, and  $c \in \mathbb{N}$  to count the interaction times, *e.g.*,  $c_{ui}$  denotes the interaction times between  $u$  and  $i$ . The row vector  $s_r$  records the interactions between one user and all the items, and the column vector  $s_c$  contains the interactions between one item and all the users. The ultimate goal is to solve the unknown entries in the matrix  $S_{ui}$ , but it is hard to model the whole matrix directly. So conventionally we either model the vectors  $s_r$  or  $s_c$  with the multinomial distribution or model the scalar  $s$  with the Bernoulli distribution.

In the matrix  $S_{ui}$ , the entries are either the counts of  $(u, i)$  interactions  $c_{ui}$  or unknown. We have all the users form the set  $\mathbb{U} = \{u_1, u_2, \dots, u_M\}$ , and all the items form the set  $\mathbb{I} = \{i_1, i_2, \dots, i_K\}$ . In IR, given a user  $u \in \mathbb{U}$ , we generate matched items from the item pool  $\mathbb{I}$ . In UT, we shall find out the potential users out of all the users  $\mathbb{U}$  given an item  $i \in \mathbb{I}$ . See Fig. 1 for a detailed illustration.

### B. Modeling with the Bernoulli and Multinomial Distributions

Both the tasks of IR and UT try to estimate the value of the unknown entries in  $S_{ui}$  with the probability. Traditionally, they are modeled with either the Bernoulli distribution on the scalar  $s$  [14], [18], [32], or the multinomial distribution on the vectors  $s_r$  [24] or  $s_c$ , as depicted in Fig. 1.

1) *The Bernoulli Distribution*: To conduct the IR and UT, we can model  $s$  as a binary random variable drawn from a Bernoulli distribution. Then, we have  $s \sim \mathbf{B}(\sigma(\phi_\theta(u, i)))$ , which means  $p(s = s_{u,i} | u, i) = \sigma(\phi_\theta(u, i))$ , where  $\sigma(\cdot)$  is the

<sup>1</sup>[https://help.aliyun.com/document\\_detail/136924.html](https://help.aliyun.com/document_detail/136924.html)

sigmoid function,  $\phi_\theta(u, i)$  is the scalar output of the model parameterized by  $\theta$ , and it will be further illustrated in Sec. III-B1.

Conventionally, the training dataset is constructed with the positive and negative samples with  $s \in \{0, 1\}$ . The likelihood of the dataset is the product of probabilities of all the single data point, *i.e.*,  $\prod_{u \in \mathbb{U}, i \in \mathbb{I}, (u, i) \in \mathcal{D}_b} p(s = s_{u, i} | u, i)$ . By maximizing the log-likelihood, we obtain the binary cross-entropy (BCE) loss:

$$l = -\frac{1}{|\mathcal{D}_b|} \sum_{u \in \mathbb{U}, i \in \mathbb{I}, (u, i) \in \mathcal{D}_b} s_{u, i} \log \sigma(\phi_\theta(u, i)) + (1 - s_{u, i}) \log(1 - \sigma(\phi_\theta(u, i))), \quad (1)$$

where  $s_{u, i} \in \{0, 1\}$ , and  $s_{u, i} = 1$  are the positive pairs in  $\mathcal{S}_{ui}$ , and  $s_{u, i} = 0$  are the negative samples randomly sampled from  $\mathcal{S}_{ui}$  with the probability  $p_n(u, i)$ , and  $\mathcal{D}_b$  is the training dataset consisting of the positive and negative samples.

Many studies have employed the Bernoulli distribution for the IR and achieved state-of-the-art (SOTA) results [14], [32]. Theoretically, it should also work for the UT, and it will be further depicted in Sec. III-A1.

2) *The Multinomial Distribution*: Within the multinomial distribution scope, the modeling objective can be either  $s_r$  and  $s_c$ . Traditionally only  $s_r$  is studied in the IR and UT research area. To the best of our knowledge, we do not find any research in the IR or UT that models  $s_c$ .

When modeling  $s_r$ , we assume that it is drawn from a multinomial distribution  $\mathbf{Mult}(N_u, \mathbf{p}_u)$  for a given user  $u$ . Here the total number of interactions  $N_u = \sum_{i \in \mathbb{I}} c_{u, i}$  of  $u$ ,  $\mathbf{p}_u = (p_{u1}, p_{u2}, \dots, p_{uK})^T$  is a  $K$ -dimensional probability vector summing to 1 [24]. The  $p_{uk} := p(i = k | u)$  is the conditional probability modeled as

$$p_{uk} = \frac{\exp \phi_\theta(u, k)}{\sum_{j \in \mathbb{I}} \exp \phi_\theta(u, j)}, \quad (2)$$

where  $\phi_\theta(u, i)$  is the scalar measuring the similarity between  $u$  and  $i$ , output of the model parameterized by  $\theta$  as in Fig. 2. Although not explicitly discussed, many research works build upon the modeling with the multinomial distribution [2], [8], [23], [38].

When the modeling objective is  $s_c$ , it is assumed to follow the multinomial distribution  $\mathbf{Mult}(N_i, \mathbf{p}_i)$  for a given item  $i$ . Here the total number of interactions  $N_i = \sum_{u \in \mathbb{U}} c_{u, i}$  of  $i$ ,  $\mathbf{p}_i = (p_{1i}, p_{2i}, \dots, p_{Mi})^T$  is a  $M$ -dimensional probability vector summing to 1, and  $p_{mi} := p(u = m | i)$  is the conditional probability.

By maximizing the multinomial loglikelihood [24], we have the losses in Eqs. 3 and 4 for them respectively:

$$l = -\frac{1}{|\mathcal{D}_m|} \sum_{u \in \mathbb{U}, i \in \mathbb{I}, (u, i) \in \mathcal{D}_m} \log \frac{\exp \phi_\theta(u, i)}{\sum_{i' \in \mathbb{I}} \exp \phi_\theta(u, i')}, \quad (3)$$

$$l = -\frac{1}{|\mathcal{D}_m|} \sum_{u \in \mathbb{U}, i \in \mathbb{I}, (u, i) \in \mathcal{D}_m} \log \frac{\exp \phi_\theta(u, i)}{\sum_{u' \in \mathbb{U}} \exp \phi_\theta(u', i)}, \quad (4)$$

TABLE I  
THE BCE LOSS WITH DIFFERENT NEGATIVE SAMPLING PROBABILITIES  $p_n(u, i)$  LEAD TO DIFFERENT OPTIMA.

$p_n(u, i) \propto$	$\phi_\theta(u, i) \sim$
$\hat{p}_{\text{data}}(u)$	$\log \hat{p}_{\text{data}}(i   u)$
$\hat{p}_{\text{data}}(i)$	$\log \hat{p}_{\text{data}}(u   i)$
$\hat{p}_{\text{data}}(u) \cdot \hat{p}_{\text{data}}(i)$	$\log \frac{\hat{p}_{\text{data}}(u, i)}{\hat{p}_{\text{data}}(u) \hat{p}_{\text{data}}(i)}$
1	$\log \hat{p}_{\text{data}}(u, i)$

where  $\mathcal{D}_m$  is the training dataset consisting of only the positive user-item pairs.

### III. A UNIFIED USER-ITEM MATCHING FRAMEWORK

In this section, we first show that modeling with the Bernoulli and multinomial distributions are theoretically equivalent. We prove that in theory they can converge to the same optima by properly selecting the negative sampling methods for the Bernoulli distribution, and setting up the configurations for the multinomial distribution.

Then, we elaborate on the proposed framework, *UniMatch*, which consists of a bidirectional bias-corrected NCE loss, *bbcNCE*, that models the  $s_r$  and  $s_c$  concurrently with the multinomial distribution leading the model convergence to  $\hat{p}_{\text{data}}(u, i)$ , a two-tower architecture that can incorporate various models, and an incremental training procedure that is tailored to the application of the private domain marketing.

The incremental training mechanism enables the model training to avoid learning from highly biased user-item distributions. The *bbcNCE* loss allows us to train only one model and then infer one set of user and item embeddings, which can be used for both the IR and UT. We choose the *bbcNCE* loss originating from the multinomial distribution over the equivalent loss setup from the Bernoulli distribution, because we empirically unveil the discovery that the former produces better, more robust results, and saves training costs dramatically as in Sec. IV.

#### A. The equivalence between modeling with the Bernoulli and multinomial distributions

When modeling with the Bernoulli distribution, we deal with  $p(s | u, i)$  directly, while we study  $p(i | u)$  or  $p(u | i)$  with the multinomial distribution. In order to bridge the gap between modeling with these two distributions, we uncover the theoretical connection between these two modeling strategies, and prove that they can converge to the same optima.

1) *Optima of modeling with the Bernoulli distribution*: We derive the optima of modeling with the Bernoulli distribution for various negative sampling methods. Inspired by Noise Contrastive Estimation (NCE) [12], we assume the positive samples of the training dataset form the set  $\mathbb{X}$ , and the negative samples form the set  $\mathbb{Y}$ . The negative samples are randomly sampled based on a certain distribution  $p_n(u, i)$ .

Assume  $\mathbb{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L\}$  contains  $L$  samples, and  $\mathbb{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_F\}$  contains  $F$  samples, and  $\mathbb{Z} = \mathbb{X} \cup \mathbb{Y} = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_{L+F}\}$  contains all the  $L + F$  samples. Here

TABLE II  
THE OPTIMA OF THE SSM LOSS AND LOSSES WITH DIFFERENT SETTINGS OF EQ. 10. WE PROPOSE TO USE THE BBCNCE AS THE LOSS OF OUR UNIMATCH FRAMEWORK FOR BOTH ITEM RECOMMENDATION AND USER TARGETING.

Settings	Objective	$\phi_\theta(u, i) \sim$	Loss
N/A	$\mathbf{s}_r$	$\log \hat{p}_{\text{data}}(i u)$	SSM [17]
$\alpha = 1, \delta_\alpha = \beta = \delta_\beta = 0$ $\alpha = \beta = 1, \delta_\alpha = \delta_\beta = 0$	$\mathbf{s}_r$ $\mathbf{s}_r, \mathbf{s}_c$	$\log \frac{\hat{p}_{\text{data}}(u, i)}{\hat{p}_{\text{data}}(u)\hat{p}_{\text{data}}(i)}$	InfoNCE [29] SimCLR [5]
$\alpha = \delta_\alpha = 1, \beta = \delta_\beta = 0$ $\alpha = \delta_\alpha = 0, \beta = \delta_\beta = 1$ $\alpha = \delta_\alpha = \beta = \delta_\beta = 1$	$\mathbf{s}_r$ $\mathbf{s}_c$ $\mathbf{s}_r, \mathbf{s}_c$	$\log \hat{p}_{\text{data}}(i u)$ $\log \hat{p}_{\text{data}}(u i)$ $\log \hat{p}_{\text{data}}(u, i)$	row-bcNCE col-bcNCE <b>bbcNCE</b>

$\mathbf{x}_l := (u, i)$ ,  $\mathbf{y}_f := (u, i)$  and  $\mathbf{z}_j := (u, i)$ , where  $u \in \mathbb{U}$  and  $i \in \mathbb{I}$ . We assign each sample  $\mathbf{z}_j$  a binary class  $C_j$ :  $C_j = 1$  if  $\mathbf{z}_j$  comes from  $\mathbb{X}$ , and  $C_j = 0$  if  $\mathbf{z}_j$  is from  $\mathbb{Y}$ .

We assume the joint probability of the positive samples in  $\mathbb{X}$  is parameterized by  $\tilde{\theta}$  as  $p_{\text{model}}(u, i; \tilde{\theta}) = p_{\text{model}}(\mathbf{z}; \tilde{\theta})$ . So we have the conditional probabilities:

$$p(\mathbf{z}|C = 1; \tilde{\theta}) = p_{\text{model}}(\mathbf{z}; \tilde{\theta}) \quad p(\mathbf{z}|C = 0; \tilde{\theta}) = p_n(\mathbf{z}). \quad (5)$$

The posterior probabilities are:

$$\begin{aligned} p(C = 1|\mathbf{z}; \tilde{\theta}) &= \frac{p_{\text{model}}(\mathbf{z}; \tilde{\theta})P(C = 1)}{p_{\text{model}}(\mathbf{z}; \tilde{\theta})P(C = 1) + p_n(\mathbf{z})P(C = 0)} \\ &= \frac{1}{1 + \exp(-G(\mathbf{z}; \tilde{\theta}))} = \sigma(G(\mathbf{z}; \tilde{\theta})), \end{aligned} \quad (6)$$

where

$$G(\mathbf{z}; \tilde{\theta}) = \log \frac{p_{\text{model}}(\mathbf{z}; \tilde{\theta})P(C = 1)}{p_n(\mathbf{z})P(C = 0)}, \quad (7)$$

and we also have

$$p(C = 0|\mathbf{z}; \tilde{\theta}) = \frac{1}{1 + \exp(G(\mathbf{z}; \tilde{\theta}))}.$$

So the log-likelihood is:

$$\begin{aligned} l(\tilde{\theta}) &= \sum_{j=1}^{L+F} C_j \log P(C_j = 1|\mathbf{z}_j; \tilde{\theta}) \\ &\quad + (1 - C_j) \log P(C_j = 0|\mathbf{z}_j; \tilde{\theta}). \end{aligned}$$

Through optimizing the binary classification using the samples in  $\mathbb{Z}$  and the corresponding binary classes  $C$ , we are actually recovering the modeling of  $s$  with the Bernoulli distribution as in Eq. 1. The dot product  $\phi_\theta(u, i)$  of Eq. 1 is  $G(\mathbf{z}; \theta)$  in Eq. 7:

$$\phi_\theta(u, i) = \log \frac{p_{\text{model}}(\mathbf{z}; \tilde{\theta})P(C = 1)}{p_n(\mathbf{z})P(C = 0)}. \quad (8)$$

It is shown that  $p_{\text{model}}(\mathbf{z}; \tilde{\theta})$  will converge to the empirical distribution  $\hat{p}_{\text{data}}(\mathbf{z})$  in [12]. As  $\mathbf{z} := (u, i)$ , we will have the following equation:

$$\phi_\theta(u, i) \approx \log \frac{\hat{p}_{\text{data}}(u, i)}{p_n(u, i)} + C', \quad (9)$$

where  $C'$  denotes some constant.

Different negative sampling strategies  $p_n(u, i)$  will lead to very different optimal  $\phi_\theta(u, i)$ . For example, if we randomly sample  $\tilde{n}$  items for each positive  $(u, i)$  pair to form the negative samples with the user  $u$ , then we have  $p_n(u, i) = \hat{p}_{\text{data}}(u) \cdot 1/K$ , where  $\hat{p}_{\text{data}}(u)$  is the empirical marginal probability of the  $u$ . Substitute  $p_n(u, i)$  into Eq. 9, we have  $\phi_\theta(u, i) \approx \hat{p}_{\text{data}}(i|u)$ . Similarly, we can derive the results in Tab. I for other negative sampling methods. Specifically, when sampling with the uniform probability, we would have  $\phi_\theta(u, i) \sim \log \hat{p}_{\text{data}}(u, i)$ , which could be used for both the IR and UT.

2) *Optima of modeling with the multinomial distribution:* When modeling  $\mathbf{s}_r$  or  $\mathbf{s}_c$  with the multinomial distribution, we have the losses in Eqs. 3 and 4. The vocabularies of the user set  $\mathbb{U}$  and item set  $\mathbb{I}$  are very large, so the calculation of the partition functions in the losses are very time-consuming and memory-exhausting, causing problems during the optimization [17]. In practice, it can be solved by implementing with the sampled softmax loss (SSM) [8], [17]. The InfoNCE loss [29] provides an alternative implementation with the sampling bias attached as in CLRec [39].

In our applications of the IR and UT, we propose to model  $\mathbf{s}_r$  and  $\mathbf{s}_c$  concurrently by combining the two losses into one, and implement it with the bias-corrected NCE loss inspired by the InfoNCE loss. The resulting loss is shown in Eq. 10:

$$\begin{aligned} l &= -\frac{1}{|\mathcal{S}_{u,i}|} \sum_{u \in \mathbb{U}, i \in \mathbb{I}, s_{u,i}=1} \alpha \cdot \log \frac{h(u, i)}{h(u, i) + \sum_{i' \in \mathbb{I}_u} h(u, i')} \\ &\quad + \beta \cdot \log \frac{o(u, i)}{o(u, i) + \sum_{u' \in \mathbb{U}_i} o(u, i')}, \end{aligned} \quad (10)$$

where  $h(u, i) = \exp(\phi_\theta(u, i) - \delta_\alpha \log \hat{p}_{\text{data}}(i))$ , and  $o(u, i) = \exp(\phi_\theta(u, i) - \delta_\beta \log \hat{p}_{\text{data}}(u))$ ,  $\mathbb{I}_u \subset \mathbb{I}$  and  $\mathbb{U}_i \subset \mathbb{U}$  contain hundreds or thousands of in-batch negative samples as in Tab. IV, and  $\hat{p}_{\text{data}}(i)$  and  $\hat{p}_{\text{data}}(u)$  are empirical marginal distributions calculated using the training data.  $\alpha$ ,  $\beta$ ,  $\delta_\alpha$  and  $\delta_\beta$  are binary numbers.  $\delta_\alpha \log \hat{p}_{\text{data}}(i)$  and  $\delta_\beta \log \hat{p}_{\text{data}}(u)$  are the bias correction terms, which ‘correct’ the biases caused by the in-batch negative sampling.

We call the first part row loss and the second part column loss (See Fig. 1). As shown in [29], the row loss can be regarded as an approximation of the loss in Eq. 3, and the column loss as an approximation of the loss in Eq. 4.

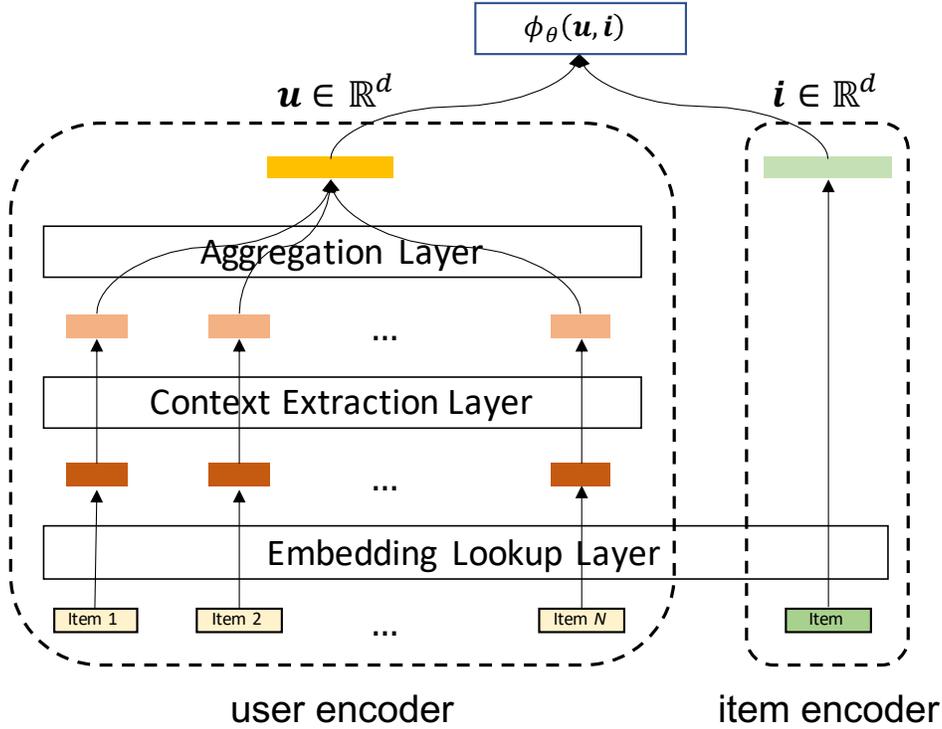


Fig. 2. The model architecture of the *UniMatch* framework. Users’ behavior sequences and items’ features go through the encoders separately, and output the  $d$ -dimensional representation vectors  $\mathbf{u}$  and  $\mathbf{i}$ . The two encoders share the same item embedding lookup table. Their  $l_2$ -normalized dot product is rescaled by the temperature hyperparameter  $\tau$  to obtain  $\phi_\theta(\mathbf{u}, \mathbf{i})$  as in Eq. 13, which is then passed to the loss functions, e.g., Eq. 1 and 10.

The InfoNCE and SimCLR losses are the special cases of Eq. 10 when the bias correction terms are omitted, i.e.,  $\delta_\alpha = \delta_\beta = 0$ . Then we have the InfoNCE loss by setting  $\alpha = 1, \beta = 0$ , and the SimCLR loss with  $\alpha = 1 = \beta = 0$  as in Tab. II.

Different settings will lead to different optima of  $\phi_\theta(\mathbf{u}, \mathbf{i})$ . It has been shown that the setting for InfoNCE will have the optimum  $\exp(\phi_\theta(\mathbf{u}, \mathbf{i})) \propto \frac{\hat{p}_{\text{data}}(\mathbf{i}|\mathbf{u})}{\hat{p}_{\text{data}}(\mathbf{i})}$  in [29]. It is straightforward that the SimCLR loss has the same optimum.

In the cases that the bias terms are retained, i.e.,  $\delta_\alpha = \delta_\beta = 1$ , the optimum of  $\phi_\theta(\mathbf{u}, \mathbf{i})$  would be different. It results in an NCE loss with the bias-correction. Different from the InfoNCE loss that converges at the point where the mutual information between users and items is maximized, the bias-corrected NCE losses (bcNCE) converge when the conditional probability is fitted by  $\phi_\theta(\mathbf{u}, \mathbf{i})$ . For example, when  $\alpha = 1, \beta = 0$ , the loss degenerates to the row loss in Eq. 10, and the optimum is  $\frac{\exp(\phi_\theta(\mathbf{u}, \mathbf{i}))}{\hat{p}_{\text{data}}(\mathbf{i})} \propto \frac{\hat{p}_{\text{data}}(\mathbf{i}|\mathbf{u})}{\hat{p}_{\text{data}}(\mathbf{i})}$ . Then we have  $\phi_\theta(\mathbf{u}, \mathbf{i}) \propto \log \hat{p}_{\text{data}}(\mathbf{i}|\mathbf{u})$ . Analogously, we have  $\phi_\theta(\mathbf{u}, \mathbf{i}) \propto \log \hat{p}_{\text{data}}(\mathbf{u}|\mathbf{i})$  for the setting  $\alpha = 0, \beta = 1$  as in Tab. II.

When  $\alpha = \beta = \delta_\alpha = \delta_\beta = 1$ , the optimum is not trivial to derive. The proof is given as follows:

In this setting, we have  $\phi_\theta(\mathbf{u}, \mathbf{i}) \propto \log \hat{p}_{\text{data}}(\mathbf{i}|\mathbf{u})$  for the row loss of Eq. 10 at its optimum, so we can assume

$$\begin{aligned} \phi_\theta(\mathbf{u}, \mathbf{i}) &= \log \hat{p}_{\text{data}}(\mathbf{i}|\mathbf{u}) + f(\mathbf{u}) \\ &= \log \hat{p}_{\text{data}}(\mathbf{u}, \mathbf{i}) - \log \hat{p}_{\text{data}}(\mathbf{u}) + f(\mathbf{u}), \end{aligned} \quad (11)$$

for a given  $\mathbf{u}$  and any  $\mathbf{i} \in \mathbb{I}$ , where  $f(\cdot)$  is an arbitrary function depending on  $\mathbf{u}$  only. Similarly, in the optimum of the column loss, we have

$$\phi_\theta(\mathbf{u}, \mathbf{i}) = \log \hat{p}_{\text{data}}(\mathbf{u}, \mathbf{i}) - \log \hat{p}_{\text{data}}(\mathbf{i}) + g(\mathbf{i}) \quad (12)$$

for a given  $\mathbf{i}$  and any  $\mathbf{u} \in \mathbb{U}$ , and  $g(\cdot)$  is an arbitrary function depending on  $\mathbf{i}$  only.

Thus, from the equivalence of Eq. 11 and 12, we have the equation that always holds for any  $\mathbf{u}$  and  $\mathbf{i}$ :  $-\log \hat{p}_{\text{data}}(\mathbf{u}) + f(\mathbf{u}) \equiv -\log \hat{p}_{\text{data}}(\mathbf{i}) + g(\mathbf{i})$ , so it must be some constant. Then we have  $\phi_\theta(\mathbf{u}, \mathbf{i}) = \log \hat{p}_{\text{data}}(\mathbf{u}, \mathbf{i}) + C'$ , where  $C'$  is some constant that is independent of  $\mathbf{u}$  and  $\mathbf{i}$ . When  $\phi_\theta(\mathbf{u}, \mathbf{i})$  converges to  $\log \hat{p}_{\text{data}}(\mathbf{u}, \mathbf{i})$ , both parts in Eq. 10 can reach their optima, so it is at least one of the solutions of the whole loss. We name the resulting loss **bbcNCE**, short for bidirectional-bias-corrected NCE. It is employed in our framework for the IR and UT, as further illustrated in Sec. III-B.

As proved in the above sections, the optima of different settings of modeling with the Bernoulli and multinomial distributions are listed in Tab. I and II. We can see that they can guide the parameterized models to converge at the same optima. Therefore, we can conclude that they are equivalent on modeling the user-item interactions depicted in Fig. 1.

## B. The UniMatch Framework

We propose to model the IR and UT jointly in the *UniMatch* framework. In details, the *UniMatch* consists of the classical two-tower architecture, the bbcNCE loss proposed in the previous section and the incremental training procedure.

1) *Architecture*: We choose this architecture for two reasons. The first is that the users and items can be processed equivalently, while the other one is that there is no feature crossing occurring before the final logits  $\phi_\theta(u, i)$ , as shown in Fig. 2. As a result, users’ and items’ embeddings can be inferred separately, and then the approximate nearest neighbor (ANN) search algorithm can be applied during serving [25].

The output of two towers are  $d$ -dimensional vectors  $\mathbf{u} = f_\theta(u) \in \mathbb{R}^d$  and  $\mathbf{i} = g_\theta(i) \in \mathbb{R}^d$ , where  $\theta$  is the model parameter. The dot product  $\langle \mathbf{u} | \mathbf{i} \rangle$ , or the function of it is used as the sufficient statistics of the probability distributions defined in Sec. III-A. We find that l2-normalizing  $\mathbf{u}$  and  $\mathbf{i}$  and then rescaling the dot product by the temperature  $\tau$  lead to better and robust results:

$$\phi_\theta(u, i) = \frac{1}{\tau} \frac{\langle \mathbf{u} | \mathbf{i} \rangle}{\|\mathbf{u}\|_2 \|\mathbf{i}\|_2}. \quad (13)$$

- **User Encoder.** In this work, users’ behavior sequences are used as the features of the user encoder, so any sequential model can be adopted here. For example, the CNN [21], [22], RNN (GRU [7], LSTM [10]), Transformer [37] and their enormous variants can be used here. In fact, they have been widely applied in the item recommendation, such as Caser [35], GRU4Rec [15] SASRec [19], and etc.

We abstract the user encoder into 3 parts, embedding lookup layer, context extraction layer and aggregation layer. Through the lookup layer, item-ids are turned into vectors. In the context extraction layer, we extract and fuse the contextual information for each item in the behavior sequence with CNN/RNN/Transformer. Finally, we aggregate all the sequential item embeddings with max/mean/last/attention pooling methods. Here the last pooling means picking the last embedding in the sequence, and the attention pooling means summing over all the embeddings with learned attention weights.

- **Item Encoder.** The item encoder takes item features as the input and outputs a representative vector. In this work, we obtain the item vectors directly from the lookup table.

Our framework is model agnostic, and in case that other formats of data is taken as the input, the corresponding models can be used to replace the sequential models.

2) *Loss Function*: As illustrated in Sec. II, we can choose to model  $S_{ui}$  with either the Bernoulli or multinomial distributions to do the IR and UT simultaneously. And we have proved that they could lead to the same optima in Sec. III-A. This implies that we can employ either the BCE loss in Eq. 1 or the NCE loss in Eq. 10 in our framework.

Both the BCE loss with the uniform negative sampling and the bbcNCE converge at the joint probability  $\hat{p}_{\text{data}}(u, i)$ . Theoretically, the two losses with the corresponding settings should perform well in both IR and UT. Our experiments show that the bbcNCE loss yields better and more robust results across all 4 datasets. In addition, bbcNCE requires 1/10 ~ 1/5 of the training time of the BCE, thus reduce the cost very much. So we choose the bbcNCE as the loss of our *UniMatch* framework.

3) *Incremental training*: Incremental training feeds the training data sequentially based on the absolute time. It has two advantages compare to feeding all the training data randomly instead. First, in this setting, we train the model every month from the saved checkpoint using the latest 1-month training data. This will save lots of cost compare to training with all the data in the past dozens of months that are shuffled randomly. The other is that the results are much better. When trained with the latest 1-month data, the model parameters will shift to fit the updated user-item distribution, and thus boost the results on predicting the near future.

## IV. EXPERIMENTS

We verify whether the proposed framework is able to yield the SOTA results for the IR and UT tasks on two public datasets and two real-world datasets. Comprehensive experiments are designed to compare the two modeling strategies of the Bernoulli and multinomial distributions, and different losses within the multinomial distribution scope are evaluated. In addition, we show that our *UniMatch* framework is model agnostic. It can adopt different models and produce better results consistently. Finally, we show that the incremental training is necessary for our applications.

### A. Experimental Setup

1) *Datasets*: We use two public datasets, the Amazon books and electronics data<sup>2</sup> as well as two real-world datasets

<sup>2</sup><http://jmcauley.ucsd.edu/data/amazon/index.html>

TABLE III  
THE STATISTICS OF THE EXPERIMENTAL DATASETS, INCLUDING TWO OPEN DATASETS AMAZON BOOKS AND ELECTRONICS AS WELL AS TWO REAL-WORLD DATASETS FROM OUR QA CLIENTS.

Data	# users	# items	# interaction	time-span	avg. #actions/user	avg. #actions/item
Books	536,409	338,739	6,132,506	31 months	11.4	18.1
Electronics	3,142,438	382,246	5,566,859	31 months	1.8	14.6
QA_e_comp	237,052	15,168	1,350,566	47 months	5.7	89.0
QA_w_comp	867,107	507	2,762,870	24 months	3.2	5449.4

TABLE IV

THE TRAINING DATA SAMPLES OF THE LOSSES DERIVED FROM MODELING WITH MULTINOMIAL DISTRIBUTIONS, *e.g.*, SSM, INFONCE, BBCNCE AND ETC.

user_id	item_seq	item_id	$\log(p(u))$	$\log(p(i))$
406690	27886 755 4609 1319	19273	-11.83447	-9.34957
357729	8926 42571 9499	39415	-9.63725	-10.91818
392972	14172 6887 177888	21632	-11.42901	-11.83447
354500	85014 850 16291	176520	-11.42901	-10.73586
15839	10528 690 173 17	272267	-11.42901	-12.52762

TABLE V

THE TRAINING DATA SAMPLES OF THE BCE LOSS DERIVED FROM MODELING WITH BERNOULLI DISTRIBUTIONS.

user_id	item_seq	item_id	label
406690	27886 755 4609 1319	19273	1
357729	8926 42571 9499	39415	1
394560	60076 5568 186 11 7 274408	16751	0
392972	14172 6887 177888	21632	1
391953	70 20167 171	6493	0

collected from two QuickAudience clients. The statistics of the data is shown in Tab. III.

- **Amazon.** We choose two commonly used datasets, Amazon books and electronics from *Amazon.com*. The datasets span from May 1996 to July 2014. We utilize the data from January 2012 to July 2014 in our experiments. For Amazon books and electronics, each sample’s behavior sequence is truncated at the length of 20 and 36, respectively.
- **QuickAudience clients.** We use *e\_comp* and *w\_comp* to represent the two merchant clients, and these two datasets span 47 and 24 months. Compared to Amazon datasets, they have comparable number of users and interactions, but the number of items are significantly less, so they are less sparse as in Tab. III. We truncate the samples at the length of 29 and 18 for *e\_comp* and *w\_comp*, respectively.

In our experiments, the next- $n$ -day prediction is set to predict for the next month. With the whole dataset spanning  $T$  months, we split the the data into train, validation and test data as  $(0, T - 1]$ ,  $(T - 2, T - 1]$  and  $(T - 1, T]$ .

In the train/validation/test data, we filter out the users/items who interact with less than 3 items/users. To train the model, we adopt the incremental training method, and consume the data consecutively according to the date  $t$ . In other words, we feed data of  $t = 1$  first and train for some epochs, and then followed by  $t = 2, 3, \dots, T - 1$ .

The losses derived from the multinomial and Bernoulli distributions require different input data formats. The differences are shown in Tab. IV and V with data samples. To be more specific, the losses like *bbcNCE* require the bias-correction terms pre-calculated from the empirical distributions of users and items in the training data, as illustrated in Eq. 10. Each record is the positive user-item pair, and the in-batch negative sampling use the users or items in the same batch to form

the negative user-item pairs. On the other hand, for the BCE loss derived from modeling with the Bernoulli distribution, the records with label 1 are positive samples, and label 0 are negative samples that are sampled with certain distributions  $p_n(u, i)$  as in Tab. I. The ratio between positive and negative samples is 1 : 1.

The test data of the IR and UT are prepared separately. The statistics of the 4 experimental datasets after splitting to train and test are shown in Tab. VI. We describe the table using the Amazon Books dataset as an example. Its train data contains 2,985,163 records as the positive samples. For the IR, the number of test users is 43,867. Each user has 1 positive item and 99 negative items that is sampled randomly from the item pool of 67967 items in total. Our experimental models predict top 10 items out of the 100 candidates for each user, and the results are evaluated using the metrics Recall@10 and NDCG@10 depicted in the following section. The same logic applies in the UT.

2) *Hyperparameters:* We have three hyperparameters to be tuned in the experiments, *i.e.*, temperature  $\tau$ , batch-size and the number of epochs, and choose the hyperparameters based on the validation data through grid search. Different datasets modeled with different distributions have their own specific hyperparameters, and the grid search results are listed in Tab. VII. The dimension  $d = 16$  is adopted for all the datasets in the experiments.

3) *Evaluation Metrics:* We employ two commonly used metrics, Recall and Normalized Discounted Cumulative Gain (NDCG), and report Recall@ $N$  and NDCG@ $N$  to evaluate the top- $N$  ranked items/users in the IR and UT. Another popular metric HitRate@ $N$  is the same as Recall@ $N$  when there is only 1 positive in the candidate pool, so we will not repeat its results here.

For the IR, the two metrics are defined as follows:

$$\text{Recall@}N = \frac{1}{|\mathbb{U}|} \sum_{u \in \mathbb{U}} \frac{|\hat{\mathcal{I}}_{u,N} \cap \mathcal{I}_u|}{\text{Min}(|\mathcal{I}_u|, N)}, \quad (14)$$

$$\text{NDCG@}N = \frac{1}{|\mathbb{U}|} \sum_{u \in \mathbb{U}} \frac{1}{Z_u} \sum_{n=1}^N \frac{\delta(\hat{i}_{u,n} \in \mathcal{I}_u)}{\log_2(n+1)}, \quad (15)$$

where  $\hat{\mathcal{I}}_{u,N}$  denotes the top- $N$  ranked items for  $u$ ,  $\hat{i}_{u,n}$  is  $n$ -th recommended item,  $\mathcal{I}_u$  is the set of ground-truth items, and  $\delta(\cdot)$  is the indicator function.  $Z_u$  is the normalization constant denoting the best possible discounted cumulative gain (DCG) for the user  $u$ , which means that all the ground-truth items are ranked at the top. For the UT the metrics are defined symmetrically.

In the experiments, we use Recall/NDCG@10 for the IR and UT across all the datasets except for QA *w\_comp*. We measure QA *w\_comp* with Recall/NDCG@5 due to its small number of items as in Tab. III. The NDCG measures the ranking status of the recommended items or targeted users, and contains more subtle information of the results, therefore we use it to select the hyperparameters as well.

TABLE VI  
THE STATISTICS OF THE 4 EXPERIMENTAL DATASETS AFTER SPLITTING INTO TRAIN AND TEST.

		Amazon Books	Amazon Electronics	QA e_comp	QA w_comp
train data		2,985,163	451,283	504,500	328,770
IR	# test users	43,867	7,916	4,685	29,168
	# item pool	67,967	14,118	1,943	221
	# top- $n$ items	10	10	10	5
	# negatives	99	99	99	49
UT	# test items	27,541	4,708	1,324	203
	# user pool	317,667	207,060	30,439	171,354
	# top- $n$ users	10	10	10	5
	# negatives	99	99	99	49

TABLE VII  
THE HYPERPARAMETERS OF ALL THE DATASETS MODELED WITH BERNOULLI OR THE MULTINOMIAL DISTRIBUTIONS.

Hyperparameters	Amazon books		Amazon Electronics		QA e_comp		QA w_comp	
	Bernoulli	Multinomial	Bernoulli	Multinomial	Bernoulli	Multinomial	Bernoulli	Multinomial
Batch-size	128	64	256	64	128	64	128	64
Temperature	0.1667	0.1667	0.5	0.5	0.25	0.125	0.125	0.1
Epochs	8	3	6	2	6	2	10	2

4) *Experimental Comparisons*: First, we compare the results between modeling with the multinomial and Bernoulli distributions, Then, the losses derived from modeling with the multinomial distribution are compared, and finally we study distinct model architectures. The detailed comparisons are listed as follows:

- **The bbcNCE versus the BCE.** Our proposed bbcNCE loss (Eq. 10) and the BCE loss with the specific negative sampling are the practical realizations of modeling the user-item interaction matrix with the multinomial and Bernoulli distributions, respectively. We experiment and evaluate their performance in both IR and UT of all the four datasets.
- **The bbcNCE versus other losses in the multinomial distribution scope.** When modeling the interaction matrix with the multinomial distribution, there are various implementations using different losses. We compare the bbcNCE loss with other well-applied losses like SSM to show its effectiveness.
- **The model-agnostic characteristic of the UniMatch framework.** We demonstrate that our framework is capable of integrating various models, including Youtube-

DNN [8], CNN used in Caser [35], RNN employed in GRU4Rec [15] and Transformer utilized in SASRec [19], and show that they produce consistent results in the IR and UT tasks.

- **The effectiveness of the incremental training.** We setup the training procedure as the incremental training month by month. By this way, the model can adapt to the latest distribution of user-item interactions, and produces much better results, particularly in the case that the item trends and users' interests shift quickly.
- **Cost Saving.** We summarize how the concurrent modeling of the IR and UT tasks by our framework saves the total cost up to 94%+.

5) *Implementations*: The code of our experiments is implemented with TensorFlow 1.12 [1] in Python 2.7, running on Nvidia GPU Tesla T4. We use the existing CNN and RNN modules in TensorFlow, and implement the Transformer based on this github repository<sup>3</sup>.

<sup>3</sup><https://github.com/Kyubyong/transformer>

TABLE VIII  
RESULTS OF IR AND UT OBTAINED BY THE BCE LOSS WITH DIFFERENT NEGATIVE SAMPLING STRATEGIES VERSUS THE BBCNCE. THE METRIC IS NDCG@10 FOR AMAZON BOOKS, AMAZON ELECTRONICS, e\_comp AND NDCG@5 FOR w\_comp. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD, THE SECOND BEST IS UNDERLINED. THE % IS OMITTED.

losses	NS: $p_n(u, i)$	Amazon Books			Amazon Electronics			QA e_comp			QA w_comp		
		IR	UT	AVG	IR	UT	AVG	IR	UT	AVG	IR	UT	AVG
BCE	$\hat{p}_{data}(u)$	<u>53.07</u>	41.95	47.51	<b>24.43</b>	10.50	<u>17.46</u>	<u>36.99</u>	4.98	20.98	<u>35.73</u>	20.59	28.16
BCE	$\hat{p}_{data}(i)$	42.85	<u>44.77</u>	43.81	13.68	11.47	12.58	6.44	<u>7.35</u>	6.90	24.29	22.24	23.27
BCE	$\hat{p}_{data}(u)\hat{p}_{data}(i)$	44.67	43.76	44.21	13.66	<u>11.81</u>	12.73	6.08	6.41	6.25	24.78	21.57	23.17
BCE	$1/MK$	52.79	42.46	<u>47.63</u>	24.34	9.81	17.08	36.51	6.70	<u>21.60</u>	35.55	<u>23.42</u>	<u>29.48</u>
bbcNCE	-	<b>57.20</b>	<b>47.67</b>	<b>52.44</b>	<u>24.39</u>	<b>12.77</b>	<b>18.58</b>	<b>37.65</b>	<b>8.25</b>	<b>22.95</b>	<b>36.48</b>	<b>24.30</b>	<b>30.39</b>

TABLE IX

RESULTS OF IR AND UT OF bbcNCE LOSS VERSUS OTHER LOSSES MODELING  $s_r$  OR  $s_c$  WITH THE MULTINOMIAL DISTRIBUTION ON AMAZON DATASETS. THE % IS OMITTED. THE DETAILS OF THE LOSSES CAN BE FOUND IN TAB. II. ‘SSM W. N’ IS THE SSM LOSS WITH THE USERS’ AND ITEMS’ REPRESENTATIONS BEING L2-NORMALIZED. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD, THE SECOND BEST RESULT IS UNDERLINED.

loss	Amazon Books						Amazon Electronics					
	IR		UT		AVG		IR		UT		AVG	
SSM w. n.	77.07	56.88	58.01	35.85	67.54	46.36	<u>48.78</u>	<u>25.82</u>	13.14	6.06	30.96	15.94
InfoNCE	70.73	48.12	68.78	46.67	69.76	47.39	28.19	15.83	<b>24.54</b>	<b>13.35</b>	26.36	14.59
SimCLR	71.53	49.35	<u>69.24</u>	47.64	70.38	48.50	27.97	15.83	<u>22.96</u>	12.68	25.46	14.26
row-bcNCE	<b>78.56</b>	<b>58.71</b>	66.71	44.44	<u>72.64</u>	<u>51.58</u>	<b>49.54</b>	<b>28.88</b>	20.13	11.00	<b>34.83</b>	<b>19.94</b>
col-bcNCE	68.23	46.03	<b>71.12</b>	<b>50.42</b>	69.67	48.23	25.68	14.33	21.57	11.95	23.63	13.14
bbcNCE	<u>77.43</u>	<u>57.20</u>	69.18	<u>47.67</u>	<b>73.31</b>	<b>52.44</b>	41.89	24.39	22.55	<u>12.77</u>	<u>32.22</u>	<u>18.58</u>

## B. Experimental Results

1) *The bbcNCE versus the BCE*: The backbone model is the Youtube-DNN with mean pooling for all the losses. As shown in Tab. VIII, we have these observations:

*i*). We compare different negative sampling methods with the BCE loss. Negative sampling with  $p_n(u, i) \propto \hat{p}_{\text{data}}(u)$  gives consistent good results in the IR, while  $p_n(u, i) \propto \hat{p}_{\text{data}}(i)$  performs well in the UT. The uniform sampling with  $p_n(u, i) = 1/MK$  gives equally good results for both the IR and UT.

Particularly, the negative sampling with  $\hat{p}_{\text{data}}(u)$  outperforms  $\hat{p}_{\text{data}}(i)$  by 51.2% on average for the IR task on the Amazon datasets. On our QA datasets, the difference is more significant. The NDCG@10 is almost 5 times higher on the QA e\_comp dataset. In contrast, the results obtained from the negative sampling with  $\hat{p}_{\text{data}}(i)$  surpasses  $\hat{p}_{\text{data}}(u)$  for the UT task. On the QA e\_comp dataset, the metric result is about 48% relatively higher.

The above comparisons show that the IR and UT tasks require different sampling methods for the BCE loss to achieve competing results.

*ii*). The proposed bbcNCE loss obtains the best or second best results for both IR and UT across all the datasets. This verifies that in theory modeling the user-item interaction matrix  $S_{ui}$  with Bernoulli and multinomial distributions makes no difference, but in practice the bbcNCE can reach better and robust results.

We hypothesize that it is due to the comparison mechanism rooted in the loss. In the BCE loss, the model parameter  $\theta$  is optimized to push the sample’s sigmoid towards 1 or 0. With the bbcNCE loss, the positive items/users are forced to compare with other items/users in the same batch as detailed in Tab. IV, and the model is optimized to allow the positive item/user to surpass all the others.

*iii*). The bbcNCE loss costs much less than the BCE loss during training. As illustrated in Tab. VII, the losses derived from the multinomial modeling paradigm requires much less training epochs to converge. For the Amazon books dataset, the BCE loss reaches the best results with 8 epochs, while bbcNCE converges in 3 epochs. In addition, the BCE loss processes 2 times of data due to the 1 : 1 sampling of

negatives per epoch as illustrated in Tab. IV and V. Therefore, the computation cost of training is about 5 times. For other datasets, the computation cost is about 6 ~ 10 times, which is calculated from the epochs in Tab. VII.

We conjecture that the comparison mechanism stated above applies here: in each epoch, a sample in the BCE loss does not provide additional information except for its divergence from the true label 1 or 0. According to the information theory, it provides no more than 1 bit information. On the contrary, a sample in the bbcNCE loss is employed to differentiate one positive item/user from the rest of the items/users in the same batch as in Tab. IV, so it can offer at most  $\log_2 64 = 6$  bits of information if the batch-size is 64. To conclude, the bbcNCE can utilize much more information per-epoch during the training, and thus speeds up the convergence and reduce the cost by a large amount.

2) *The bbcNCE versus other losses in the multinomial distribution scope*: In comparison with other losses that modeling either  $s_r$  or  $s_c$  or both, the bbcNCE guides the model to learn the joint probability  $p(u, i)$  from the empirical distribution  $\hat{p}_{\text{data}}(u, i)$ . The experimental results in Tab. IX and X are in alignment with our theoretical proof. We analyze the experiments from the following three aspects.

*i*). We can use only one model trained with the bbcNCE loss to serve for both IR and UT in our applications. For IR, the bbcNCE is on par with the losses that model  $s_r$  and lead to the convergence at  $\hat{p}_{\text{data}}(i|u)$ , *i.e.*, SSM and row-bcNCE. For UT, its results match with the col-bcNCE loss that models  $s_c$  and converges at  $\hat{p}_{\text{data}}(u|i)$  as in Tab. II. For both IR and UT, the bbcNCE loss produces the best or second best results robustly, which makes it the best choice for our QA applications.

*ii*). The bias correction plays the key role on lifting the results in IR. As shown in Tab. IX and X, the bbcNCE and row-bcNCE losses always produce much better results than the InfoNCE and SimCLR losses that have no bias correction. The SSM loss is implemented with negative sampling and bias correction so that it converges to  $\hat{p}_{\text{data}}(i|u)$  in theory. However, its performance is usually inferior than the bbcNCE and row-bcNCE. The SSM loss draws negative samples from the whole item vocabulary, while in contrast the bbcNCE samples negatives from in-batch training data. Therefore, in

TABLE X  
RESULTS OF IR AND UT OF THE bbcNCE LOSS VERSUS OTHER LOSSES MODELING  $s_r$  OR  $vs_u$  WITH THE MULTINOMIAL DISTRIBUTION ON QUICKAUDIENCE DATASETS. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD, THE SECOND BEST IS UNDERLINED. THE % IS OMITTED.

loss	IR		QA e_comp UT		AVG		IR		QA w_comp UT		AVG	
	Recall	NDCG	Recall	NDCG	Recall	NDCG	Recall	NDCG	Recall	NDCG	Recall	NDCG
SSM w. n.	58.22	35.57	15.37	6.85	36.79	21.21	49.17	<u>36.54</u>	34.95	23.95	<u>42.06</u>	<u>30.25</u>
InfoNCE	15.82	7.33	14.62	7.29	15.22	7.31	38.14	28.63	27.59	18.09	32.87	23.36
SimCLR	16.25	7.19	15.30	7.18	15.77	7.18	36.56	27.26	35.36	24.10	35.96	25.68
row-bcNCE	<b>62.37</b>	<b>38.49</b>	15.98	7.27	39.17	<u>22.88</u>	<b>50.17</b>	<b>37.10</b>	31.53	20.92	40.85	29.01
col-bcNCE	25.62	12.61	<b>18.09</b>	<b>8.35</b>	21.86	10.48	34.32	24.24	<b>38.42</b>	<b>24.87</b>	36.37	24.56
bbcNCE	<u>61.35</u>	<u>37.65</u>	<u>17.63</u>	<u>8.25</u>	<b>39.49</b>	<b>22.95</b>	<u>49.54</u>	36.48	<u>35.47</u>	<u>24.30</u>	<b>42.50</b>	<b>30.39</b>

TABLE XI  
STATISTICS OF THE POPULARITY/ACTIVENESS OF ITEMS/USERS RETRIEVED BY DIFFERENT LOSSES. WE MEASURE THE MEDIAN AND AVERAGE OF ITEMS' POPULARITY AND USERS' ACTIVENESS. THE POPULARITY AND ACTIVENESS ARE DEFINED AS THE NUMBER OF INTERACTIONS OCCURRED IN THE PAST ONE YEAR. HERE 'MED' AND 'AVG' STANDS FOR MEDIAN AND AVERAGE RESPECTIVELY.

losses	Amazon Books				Amazon Electronics				QA e_comp				QA w_comp			
	IR		UT		IR		UT		IR		UT		IR		UT	
	med	avg	med	avg	med	avg	med	avg	med	avg	med	avg	med	avg	med	avg
SSM w. n.	25	72	6	13.6	232	491	4	4.8	94	187	4	6.4	11969	15176	5	7.1
InfoNCE	16	46	6	13.6	34	139	4	5.1	52	104	4	6.3	2332	5815	4	6.3
SimCLR	16	47	6	13.3	33	125	4	5.1	55	113	4	6.5	2294	5961	5	7.1
row-bcNCE	27	78	6	12.9	236	496	4	4.9	138	245	4	6.6	11969	15189	5	6.7
col-bcNCE	16	47	6	14.7	39	150	4	5.2	96	173	4	7.2	3063	6456	5	7.8
bbcNCE	23	69	6	14.1	160	400	4	5.1	138	246	5	7.4	12837	15320	5	7.1

our monthly incremental training setting, the bbcNCE only encounters negative items within the current month, and the SSM compares the positive items to the negative items sampled from the whole vocabulary. This brings the advantage of the bbcNCE on fitting on the latest data distribution, thus makes the results better.

In UT, the performance of bbcNCE and col-bcNCE does not always surpass other losses by a large margin. This implies that the bias correction from the users' empirical distribution is not that effective anymore when the data is sparse. We speculate that the marginal distribution calculated from the sparse dataset is not reliable. Because most users only have very few interactions, the  $\hat{p}_{data}(u)$  computed will not be statistically significant. On the contrary, if the user behavior data is relatively rich, the col-bcNCE and bbcNCE that correct the users' distribution bias produce much better results compared to the InfoNCE and SimCLR loss, as shown on the Amazon Books, QA e\_comp and w\_comp datasets.

iii). The experiment results in Tab. IX and X show that the SimCLR and InfoNCE losses yield very close results. This is in align with our claim that they both converges at  $\hat{p}_{data}(u, i)/(\hat{p}_{data}(u)\hat{p}_{data}(i))$  as in Tab. II. As shown in [29], the InfoNCE loss converges when the mutual information between the two matching variables is maximized. When it is applied in IR, it will tend to recommend users with unpopular items, because those items usually have high point-wise mutual information with the users.

We use the number of historical interactions to measure

the popularity/activeness of items/users as in Tab. XI. For example, if an item is purchased 100 times in the past 1 year, then its popularity is 100. For all the top- $n$  items/users retrieved by the model, we calculate the median and average popularity/activeness. Taking the Amazon Books dataset as an example, the median item popularity from the InfoNCE and SimCLR losses is 16, which means that 50% of the recommended items have less or equal than 16 interactions in the past 1 year. In contrast, the median popularity of the bbcNCE, row-bcNCE and SSM losses ranges from 23 to 27. In all the datasets in Tab. XI, the InfoNCE and SimCLR losses always tend to recommend unpopular items. In UT, the two losses prefer inactive users in general, but it is not that evident because most users do not have many interactions.

3) *The model-agnostic characteristic of the UniMatch framework*: We implement 5 types of context extractors, *i.e.*, Youtube-DNN<sup>4</sup>, 1-layer CNN, GRU, LSTM, and 1-layer Transformer, and 4 types of aggregators, *i.e.*, mean pooling, last pooling, max pooling and attention pooling. So we have 20 models in total. We report the results of QA w\_comp dataset in Tab. XII. The results of Max pooling are always worse than other aggregators, so we omit them.

In general, the results from different models under the same loss does not differ notably. This suggests that the complicated context extraction methods that thrive in processing the text sequence are not superior when dealing with users' sequential

<sup>4</sup>Youtube-DNN represents no context extractor here, *i.e.*, the lookup item embeddings are directly passed to the aggregation layer.

TABLE XII

RESULTS OF IR AND UT OF bbcNCE LOSS VERSUS OTHER LOSSES EVALUATED USING VARIOUS MODELS ON OUR QA w\_comp DATASETS. THE METRIC IS NDCG@5. EACH COLUMN REPRESENTS ONE KIND OF MODEL, AND ‘CEX’ MEANS THE CONTEXT EXTRACTOR OF THE MODEL, AND ‘AGG’ IS THE SEQUENCE AGGREGATOR. ‘MEAN’, ‘LAST’, ‘ATTN’ ARE MEAN, LAST AND ATTENTION POOLING RESPECTIVELY. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD, THE SECOND BEST IS UNDERLINED. THE % IS OMITTED.

	CEX AGG	Youtube-DNN			CNN-1			GRU			LSTM			Transformer-1		
		mean	last	attn	mean	last	attn									
IR	SSM w. n.	<u>36.54</u>	27.51	<u>36.62</u>	35.67	30.61	36.05	<u>36.85</u>	35.21	<b>36.40</b>	35.93	<u>35.92</u>	35.59	35.27	<b>34.89</b>	<b>37.06</b>
	InfoNCE	28.63	19.84	27.58	27.82	23.38	27.12	28.13	27.58	27.12	27.12	27.98	28.39	27.91	25.67	26.22
	SimCLR	27.26	20.45	27.17	27.48	23.96	26.59	27.46	27.32	26.95	26.47	26.89	27.49	26.02	24.92	25.87
	row-bcNCE	<b>37.10</b>	<u>28.35</u>	<u>36.61</u>	<b>36.86</b>	<u>31.08</u>	<u>37.01</u>	36.27	<u>36.22</u>	<u>35.87</u>	<b>36.69</b>	<b>36.47</b>	<b>36.53</b>	<b>36.81</b>	33.95	<u>36.79</u>
	col-bcNCE	24.24	9.10	<u>26.25</u>	24.76	11.66	<u>23.52</u>	19.93	19.02	18.80	16.90	12.72	16.58	15.89	17.73	24.99
	bbcNCE	<u>36.48</u>	<b>28.52</b>	<b>36.77</b>	<u>36.26</u>	<b>31.39</b>	<b>37.12</b>	<b>37.05</b>	<b>36.31</b>	35.70	<u>36.11</u>	35.85	<u>35.98</u>	<u>36.02</u>	<u>34.78</u>	36.43
UT	SSM w. n.	23.95	13.61	23.88	25.17	14.97	23.65	22.01	24.49	25.61	24.49	<u>26.02</u>	24.02	20.53	20.15	23.25
	InfoNCE	18.09	<u>13.91</u>	19.92	24.21	13.35	22.85	23.91	23.78	24.08	26.24	24.07	24.31	19.21	19.66	20.24
	SimCLR	24.10	<u>13.73</u>	23.57	23.89	14.86	25.12	25.72	26.23	25.43	<u>26.41</u>	24.54	26.01	<b>25.68</b>	20.69	23.40
	row-bcNCE	20.92	13.63	20.61	23.04	16.86	22.42	24.66	24.39	24.96	23.89	25.09	25.09	22.33	20.23	21.09
	col-bcNCE	<b>24.87</b>	<b>14.64</b>	<b>26.54</b>	<b>26.53</b>	<b>18.65</b>	<u>26.72</u>	<u>26.82</u>	<u>28.50</u>	<u>25.81</u>	25.79	24.44	<u>28.02</u>	24.03	<b>24.02</b>	<b>25.93</b>
	bbcNCE	<u>24.30</u>	<u>13.33</u>	<u>24.29</u>	<u>25.67</u>	<u>18.35</u>	<b>26.79</b>	<b>27.64</b>	<b>29.47</b>	<b>26.63</b>	<b>27.85</b>	<b>27.93</b>	<b>28.88</b>	<u>25.44</u>	<u>22.64</u>	<u>24.06</u>

behavior data, *e.g.*, the Amazon and QA datasets. It further implies that the contextual information is not very important on understanding users’ single interaction with an item. In contrast, a word’s exact meaning is defined by its context in natural language processing (NLP). This finding motivates us to adopt the Youtube-DNN with mean pooling as our default model for QA to save the computation cost in practice.

Across all the models, the proposed bbcNCE loss gives best or close to the best results in both tasks. For IR, the bbcNCE and row-bcNCE are the top two in almost all the models, while the bbcNCE and col-bcNCE are the best for UT. The results further verify our statements on the losses’ optima as in Tab. II, regardless of the choice of models.

4) *The effectiveness of the incremental training*: In our *UniMatch* framework, we train the model incrementally month by month. We adopt this training setup for our specific applications in QA, because IR and UT campaigns are usually organized monthly in private domain marketing. In addition, the incremental training can save cost, and improve the prediction results as depicted in Fig. 3.

We observe that the gain of the incremental training is crucial on the Amazon Books and QA e\_comp datasets. The NDCG metric of the model trained till 1 months before the test date is much higher than trained till 2 or 3 months before. We speculate that their items are very sensitive to time. For example, users prefer to buy new and popular books, which may vary quickly. So we have to keep training the model with the latest data to adapt for the trend. On the other hand, the results of the Amazon Electronics and QA w\_comp datasets are relatively stable during the incremental training. This implies their items are more stable over time, and their users’ interests are relatively static.

5) *Cost Saving*: We demonstrate that the flexibility of our framework and the proposed bbcNCE loss enable a significant cost saving in practice.

*i*). We choose the bbcNCE instead of the BCE loss to

reduce the data consumption during the training as in Sec. IV-B1. Thus we reduce the cost to  $1/10 \sim 1/5$ . Our theoretical analysis and experimental results show that the performance is on par or better.

*ii*). We propose the bbcNCE so that we can train only one model to do both IR and UT without performance decline as depicted in Sec. IV-B2. We can reduce both the training and prediction cost to  $1/2$ . At the same time, the underlying management cost of multiple models is also saved.

*iii*). We analyse various popular models and choose the simplest Youtube-DNN with mean pooling as our default model. Experiments in Sec. IV-B3 show that it can generate SOTA comparable results on all the datasets. Therefore, we do not have to choose too complex models and relieve ourselves from the large computation burden.

*iv*). We adopt the incremental training setup as illustrated in Sec. IV-B4. With the conventional training strategy, we use past 1 year data to train from scratch monthly. Using the incremental training, we can just utilize the past 1 month data and train from the latest model checkpoint. By this way, we can reduce the training cost to  $1/12$ .

To conclude, in the applications of QA, our *UniMatch* framework can reduce the training cost up to  $1/240 \sim 1/120$  while keeping the SOTA comparable performance. The prediction cost is reduced to  $1/2$  and eliminate the management cost of multiple models. The training cost is usually about 90% of total cost, so we can save up to 94%+.

## V. RELATED WORK

### A. Item Recommendation

Item recommendation (IR) have been studied in both academia and industry for decades. The collaborative filtering (CF) and its variants are widely adopted in the early years [34]. Later its descendant, the matrix factorization (MF) [9], is proposed to solve the problem more elegantly with higher accuracy. Then, the Probabilistic Matrix Factorization (PMF)

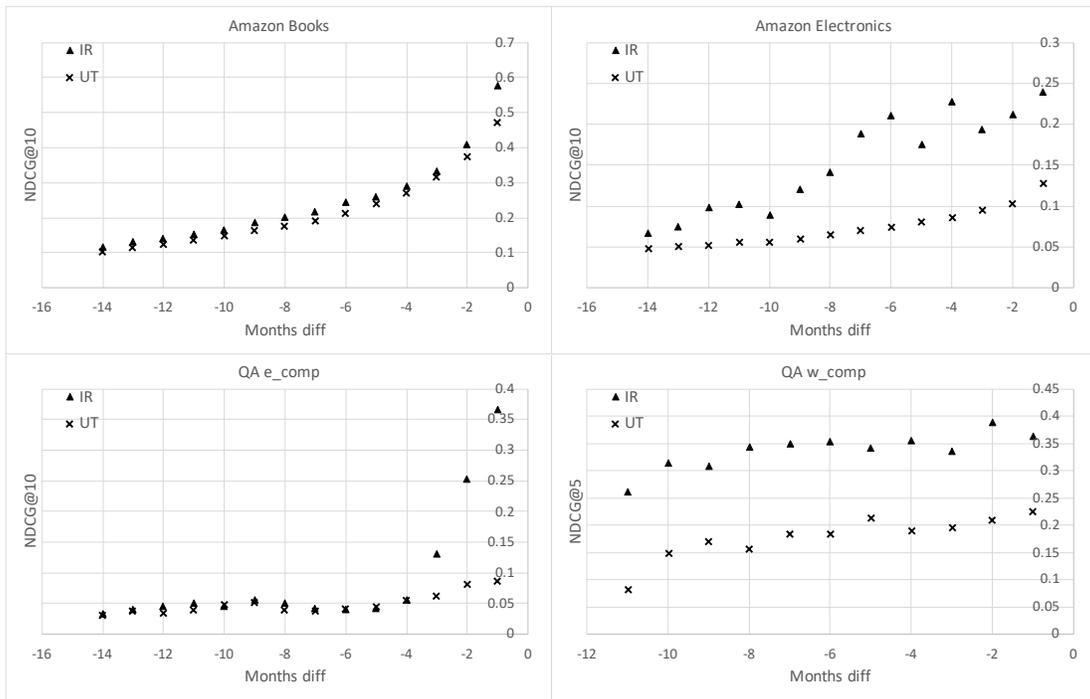


Fig. 3. The intermediate results from the incremental training setup for the Amazon and QA datasets. The  $y$ -axis is the NDCG metric, and  $x$ -axis represents the number of months ahead of the date of the test data. The metric increases steadily as newer data is feed into training.

[27] builds a solid theoretic foundation for the MF models based on the probability theory, *i.e.*, PMF models  $s$  with Gaussian distributions. Later, the Bernoulli distribution is shown to be superior in modeling  $s$  [14].

In recent years, the neural networks have become a significant component for the recommendation algorithms [14], [32], and contributed greatly for the recommender systems in industry. There are two common stages in a large-scale industrial recommendation application, *i.e.*, the candidate generation stage and the ranking stage. The former stage is usually formulated as a multi-class classification problem to quickly select a small set of item candidates from a vast number of items [2], [8], [23]. In the ranking stage, the problem is formulated as a binary classification to rank all the selected candidates [6], [28], [40]. Although not directly declared in many of these research, the underlying probability theory for the candidate generation stage is to model  $s_r$  with the multinomial distribution [24], and is to model  $s$  with the Bernoulli distribution for the ranking stage [14].

In the candidate generation stage, the huge number of items causes problems on calculating the partition function of the loss during the optimisation (as in Eq. 3). The sampled softmax (SSM) loss [17] is widely employed to solve the problem [8], [23]. Recently, the InfoNCE loss is exploited in item recommendation to suppress popular items during candidate generation [39].

## B. User Targeting

User targeting (UT) mines the potential users for given items. The *item* could be anything that users can interact with, *e.g.*, an insurance product [20], a company/business [26], [30], [31], a specific message (*e.g.*, tweets) on social medias [11], [36] and even another user [13], etc. UT is usually formulated as a binary classification problem, and solved with models like SVM, LR and neural networks, etc [3], [4], [16].

In an e-commerce company, the *item* could be a product, a brand, a product category, and a merchant, etc. The number of the *items* ranges from thousands to hundreds of millions. It is impractical to model each item respectively, so we commonly model the items all together via binary classification like [32].

In the above applications, researchers implicitly model  $s$  with the Bernoulli distribution, and the negative samples are generated with probability  $p_n(u, i) \propto \hat{p}_{\text{data}}(i)$ .

## VI. CONCLUSIONS

In this work, we propose the *UniMatch* framework that can be applied in both IR and UT to help merchants conduct the private domain marketing. Our framework is model agnostic and consists of a two-tower architecture as well as the incremental training setup and the proposed bbcNCE loss. Through comprehensive comparisons with different losses, models and training procedures, we show that our framework can generate SOTA comparable results theoretically and experimentally. Our framework reduces more than 94% of the total cost, making it affordable for merchants to do private domain marketing with the SOTA performance.

## REFERENCES

- [1] Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al.: Tensorflow: a system for large-scale machine learning. In: OSDI. vol. 16, pp. 265–283 (2016)
- [2] Cen, Y., Zhang, J., Zou, X., Zhou, C., Yang, H., Tang, J.: Controllable multi-interest framework for recommendation. In: Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. pp. 2942–2951 (2020)
- [3] Chang, C.C., Lin, C.J.: Libsvm: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)* **2**(3), 1–27 (2011)
- [4] Chen, T., Guestrin, C.: Xgboost: A scalable tree boosting system. In: Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining. pp. 785–794. ACM (2016)
- [5] Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: International conference on machine learning. pp. 1597–1607. PMLR (2020)
- [6] Cheng, H.T., Koc, L., Harmsen, J., Shaked, T., Chandra, T., Aradhye, H., Anderson, G., Corrado, G., Chai, W., Ispir, M., et al.: Wide & deep learning for recommender systems. In: Proceedings of the 1st Workshop on Deep Learning for Recommender Systems. pp. 7–10. ACM (2016)
- [7] Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y.: Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078* (2014)
- [8] Covington, P., Adams, J., Sargin, E.: Deep neural networks for youtube recommendations. In: Proceedings of the 10th ACM conference on recommender systems. pp. 191–198 (2016)
- [9] Funk, S.: Netflix update: Try this at home, <http://sifter.org/~simon/journal/20061211.html>
- [10] Gers, F.A., Schmidhuber, J., Cummins, F.: Learning to forget: Continual prediction with lstm. *Neural computation* **12**(10), 2451–2471 (2000)
- [11] Gui, T., Liu, P., Zhang, Q., Zhu, L., Peng, M., Zhou, Y., Huang, X.: Mention recommendation in twitter with cooperative multi-agent reinforcement learning. In: Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 535–544 (2019)
- [12] Gutmann, M., Hyvärinen, A.: Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In: Proceedings of the thirteenth international conference on artificial intelligence and statistics. pp. 297–304. JMLR Workshop and Conference Proceedings (2010)
- [13] Guy, I.: People recommendation on social media. In: Social information access, pp. 570–623. Springer (2018)
- [14] He, X., Liao, L., Zhang, H., Nie, L., Hu, X., Chua, T.S.: Neural collaborative filtering. In: Proceedings of the 26th international conference on world wide web. pp. 173–182 (2017)
- [15] Hidasi, B., Karatzoglou, A., Baltrunas, L., Tikk, D.: Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939* (2015)
- [16] Hosmer Jr, D.W., Lemeshow, S., Sturdivant, R.X.: Applied logistic regression, vol. 398. John Wiley & Sons (2013)
- [17] Jean, S., Cho, K., Memisevic, R., Bengio, Y.: On using very large target vocabulary for neural machine translation. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 1–10 (2015)
- [18] Johnson, C.C.: Logistic matrix factorization for implicit feedback data. *Advances in Neural Information Processing Systems* **27**(78), 1–9 (2014)
- [19] Kang, W.C., McAuley, J.: Self-attentive sequential recommendation. In: 2018 IEEE international conference on data mining (ICDM). pp. 197–206. IEEE (2018)
- [20] Kim, Y., Street, W.N.: An intelligent system for customer targeting: a data mining approach. *Decision Support Systems* **37**(2), 215–228 (2004)
- [21] Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Pereira, F., Burges, C.J.C., Bottou, L., Weinberger, K.Q. (eds.) *Advances in Neural Information Processing Systems* 25, pp. 1097–1105. Curran Associates, Inc. (2012)
- [22] LeCun, Y., Boser, B.E., Denker, J.S., Henderson, D., Howard, R.E., Hubbard, W.E., Jackel, L.D.: Handwritten digit recognition with a back-propagation network. In: Touretzky, D.S. (ed.) *Advances in Neural Information Processing Systems* 2, pp. 396–404. Morgan-Kaufmann (1990)
- [23] Li, C., Liu, Z., Wu, M., Xu, Y., Zhao, H., Huang, P., Kang, G., Chen, Q., Li, W., Lee, D.L.: Multi-interest network with dynamic routing for recommendation at tmall. In: Proceedings of the 28th ACM International Conference on Information and Knowledge Management. pp. 2615–2623 (2019)
- [24] Liang, D., Krishnan, R.G., Hoffman, M.D., Jebara, T.: Variational autoencoders for collaborative filtering. In: Proceedings of the 2018 world wide web conference. pp. 689–698 (2018)
- [25] Liu, T., Moore, A.W., Gray, A.G., Yang, K.: An investigation of practical approximate nearest neighbor algorithms. In: NIPS. vol. 12, p. 2004. Citeseer (2004)
- [26] Lo, S.L., Chiong, R., Cornforth, D.: Ranking of high-value social audiences on twitter. *Decision Support Systems* **85**, 34–48 (2016)
- [27] Mnih, A., Salakhutdinov, R.R.: Probabilistic matrix factorization. In: *Advances in neural information processing systems*. pp. 1257–1264 (2008)
- [28] Ni, Y., Ou, D., Liu, S., Li, X., Ou, W., Zeng, A., Si, L.: Perceive your users in depth: Learning universal user representations from multiple e-commerce tasks. *arXiv preprint arXiv:1805.10727* (2018)
- [29] Oord, A.v.d., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018)
- [30] Pang, G., Jiang, S., Chen, D.: A simple integration of social relationship and text data for identifying potential customers in microblogging. In: *International Conference on Advanced Data Mining and Applications*. pp. 397–409. Springer (2013)
- [31] Pennacchiotti, M., Popescu, A.M.: Democrats, republicans and starbucks aficionados: user classification in twitter. In: Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining. pp. 430–438 (2011)
- [32] Rendle, S., Krichene, W., Zhang, L., Anderson, J.: Neural collaborative filtering vs. matrix factorization revisited. In: *Fourteenth ACM Conference on Recommender Systems*. pp. 240–248 (2020)
- [33] Ricci, F., Rokach, L., Shapira, B.: Introduction to recommender systems handbook. In: *Recommender systems handbook*, pp. 1–35. Springer (2011)
- [34] Su, X., Khoshgoftaar, T.M.: A survey of collaborative filtering techniques. *Advances in artificial intelligence* **2009** (2009)
- [35] Tang, J., Wang, K.: Personalized top-n sequential recommendation via convolutional sequence embedding. In: Proceedings of the eleventh ACM international conference on web search and data mining. pp. 565–573 (2018)
- [36] Tang, L., Ni, Z., Xiong, H., Zhu, H.: Locating targets through mention in twitter. *World Wide Web* **18**(4), 1019–1049 (2015)
- [37] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: *Advances in Neural Information Processing Systems*. pp. 5998–6008 (2017)
- [38] Yi, X., Yang, J., Hong, L., Cheng, D.Z., Heldt, L., Kumthekar, A., Zhao, Z., Wei, L., Chi, E.: Sampling-bias-corrected neural modeling for large corpus item recommendations. In: Proceedings of the 13th ACM Conference on Recommender Systems. pp. 269–277 (2019)
- [39] Zhou, C., Ma, J., Zhang, J., Zhou, J., Yang, H.: Contrastive learning for debiased candidate generation in large-scale recommender systems. In: Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining. pp. 3985–3995 (2021)
- [40] Zhou, G., Zhu, X., Song, C., Fan, Y., Zhu, H., Ma, X., Yan, Y., Jin, J., Li, H., Gai, K.: Deep interest network for click-through rate prediction. In: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. pp. 1059–1068. ACM (2018)