

CauseRec: Counterfactual User Sequence Synthesis for Sequential Recommendation

Shengyu Zhang*

Zhejiang University, China

sy_zhang@zju.edu.cn

Dong Yao*

Zhejiang University, China

yaodongai@zju.edu.cn

Zhou Zhao†

Zhejiang University, China

zhaozhou@zju.edu.cn

Tat-seng Chua

National University of Singapore,

Singapore

dcscts@nus.edu.sg

Fei Wu†

Zhejiang University, China

wufei@zju.edu.cn

ABSTRACT

Learning user representations based on historical behaviors lies at the core of modern recommender systems. Recent advances in sequential recommenders have convincingly demonstrated high capability in extracting effective user representations from the given behavior sequences. Despite significant progress, we argue that solely modeling the observational behaviors sequences may end up with a brittle and unstable system due to the noisy and sparse nature of user interactions logged. In this paper, we propose to learn accurate and robust user representations, which are required to be less sensitive to (attack on) noisy behaviors and trust more on the indispensable ones, by modeling counterfactual data distribution. Specifically, given an observed behavior sequence, the proposed CauseRec framework identifies dispensable and indispensable concepts at both the fine-grained item level and the abstract interest level. CauseRec conditionally samples user concept sequences from the counterfactual data distributions by replacing dispensable and indispensable concepts within the original concept sequence. With user representations obtained from the synthesized user sequences, CauseRec performs contrastive user representation learning by contrasting the counterfactual with the observational. We conduct extensive experiments on real-world public recommendation benchmarks and justify the effectiveness of CauseRec with multi-aspects model analysis. The results demonstrate that the proposed CauseRec outperforms state-of-the-art sequential recommenders by learning accurate and robust user representations.

CCS CONCEPTS

- Information systems → Recommender systems.

*These authors contributed equally to this work.

†Corresponding Authors.

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

SIGIR '21, July 11–15, 2021, Virtual Event, Canada

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8037-9/21/07...\$15.00

<https://doi.org/10.1145/3404835.3462908>

KEYWORDS

Sequential Recommendation, User Modeling, Contrastive Learning, Counterfactual Representation

ACM Reference Format:

Shengyu Zhang, Dong Yao, Zhou Zhao, Tat-seng Chua, Fei Wu. 2021. CauseRec: Counterfactual User Sequence Synthesis for Sequential Recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '21), July 11–15, 2021, Virtual Event, Canada*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3404835.3462908>

1 INTRODUCTION

Due to the overwhelming data that people are facing on the Internet, personalized recommendation has become vital for retrieving information and discovering content. Accurately characterizing and representing users plays a vital role in a successful recommendation framework. Since users' historical interactions are sequentially dependent and by nature time-evolving, recent advances [40, 42, 43, 48, 57–60, 63, 71, 74, 82] pay attention to sequential recommendation, which captures the current and recent preference by exploiting the sequentially modeled user-item interactions.

A sequential recommender aims to predict the next item a user might interact with based on the historical interactions. The challenging and open-ended nature of sequence modeling lends itself to a variety of diverse models. Traditional methods mainly exploit Markov chains [15] and factorization machines [22, 49] to capture lower-order sequential dependencies. Following these works, the higher-order Markov Chain and RNN (Recurrent Neural Network) [18, 21, 66] are proposed to model the complex high-order sequential dependencies. More recently, MIND is proposed to transform the historical interactions into multiple interest vectors using the capsule network [51]. ComiRec [5] differs from MIND by leveraging the attention mechanism and introducing a factor to control the balance of recommendation accuracy and diversity.

Despite significant progress made with these frameworks, there are some challenges demanding further explorations. A vital challenge comes from the noisy nature of implicit feedback. Due to the ubiquitous distractions that may affect the users' first impressions (such as caption bias [39], position bias [24], and sales promotions), there are inconsistencies between users' interest and their clicking behaviors, known as the natural noise [45]. Another challenge

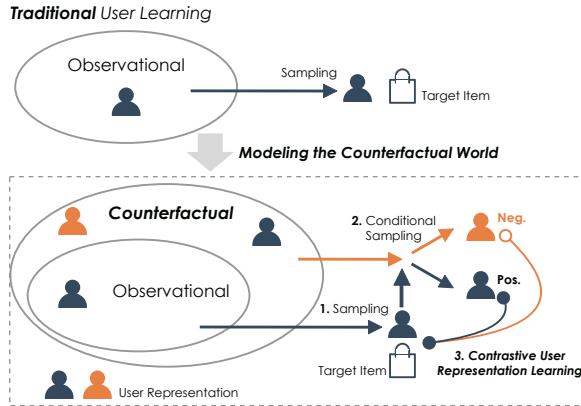


Figure 1: An illustration of the proposed contrastive user representation learning by modeling the counterfactual world (below), compared with most traditional approaches that solely model the observational user sequences (above).

relates to the deficiency of existing methods in confronting data sparsity problem in recommender systems where users in general only interact with a limited number of items compared with the item gallery which can easily reach 100 million in large live systems. Therefore, solely modeling the observational behavior sequences that can be both sparse and noisy may end up with a brittle system that is less satisfactory. To this end, learning **accurate** and **robust** users' user representations is essential for recommender systems.

In this paper, we propose **Counterfactual User Sequence Synthesis** for Sequential Recommendation, abbreviated as **CauseRec**. The essence of CauseRec in confronting the data sparsity problem is to model the counterfactual data distribution rather than the observational sparse data distribution where the latter can be a subset of the former one, as shown in Figure 1. We mainly aim to answer the counterfactual question, "what the user representation would be if we intervene on the observed behavior sequence?". Specifically, given the observed behavior sequence, we identify indispensable/dispensable concepts at both the fine-grained item level and the abstract interest level. A concept indicates a certain aspect of user interest/preference. We perform counterfactual transformations on both the item-level and the interest-level user concept sequences. We obtain counterfactually positive user representation by modifying dispensable concepts, and counterfactually negative user representation by replacing indispensable concepts. To learn **accurate** and **robust** user representations, we propose to conduct contrastive learning between: 1) the observational and the counterfactual user representations; and 2) the user representations and the target items. Contrast with such out-of-distribution hard negatives potentially makes the learned representations **robust** since they are less sensitive to dispensable/noisy concepts. Contrast with such out-of-distribution positives potentially makes the learned representations **accurate** since they will trust more on the indispensable concepts that are better representing user's interest.

We conduct in-depth experiments to validate the effectiveness of the proposed CauseRec architectures on various public recommendation datasets. With a naive deep candidate generation (or matching) architecture as the baseline method, CauseRec outperforms

SOTA sequential recommenders for deep candidate generation. We conduct comprehensive model analysis to uncover how different building blocks and hyper-parameters affect the performance of CauseRec. Case studies further demonstrate that CauseRec can help learn accurate user representations. To summarize, this paper makes the following key contributions:

- We propose to model the counterfactual data distribution (besides the observational data distribution) to confront the data sparsity problem for recommendation.
- We devise the CauseRec framework which learns *accurate* and *robust* user representations with counterfactual transformations on both fine-grained item-level and abstract interest-level, and with various contrastive objectives.
- We conduct extensive experiments and show that with a naive deep candidate generation architecture as the baseline, CauseRec can outperform SOTA sequential recommenders.

2 RELATED WORKS

2.1 Sequential Recommendation

Sequential recommendation can be traced back to leveraging Markov-chain [14, 15] and factorization machines [22, 49]. To capture long-term and multi-level cascading dependencies, deep learning based techniques (e.g., RNNs [10, 21, 47, 66] and CNNs [55, 75]) are incorporated into sequential modeling. DNNs are known to have enticing representation capability and have the natural strength to capture comprehensive relations [76] over different entities (e.g., items, users, interactions). Recently, there are works that explore advanced techniques, e.g., memory networks [53], attention mechanisms [56, 79], and graph neural networks [9, 26, 31, 36, 81] for sequential recommendation [6, 23, 29, 54, 61, 67, 72]. Typically, MIND [32] adopts the dynamic routing mechanism to aggregate users' behaviors into multiple interest vectors. ComiRec [5] differs from MIND by leveraging the attention mechanism for user representations and proposes a factor for the trade-off between recommendation diversity and accuracy. Different from the above works that solely model the observational user sequences, we take a step further to model the counterfactual data distributions. By contrasting the user representations of the observation with the counterfactual, we aim to learn user encoders that can better confront out-of-distribution user sequences and learn accurate and robust user representations.

2.2 Contrastive Learning for Recommendation

A growing number of attempts have been made to exploit the complementary power of self-supervised learning (e.g., contrastive learning) and deep learning, with domains varying from computer vision [12, 17, 68], natural language generation [7, 77], to graph embedding [46]. However, how to consolidate the merits of contrastive learning into recommendation remains largely unexplored in the literature. Recently, Sun *et al.*[33] adopt noise contrastive estimation [16] to transfer the knowledge from a large natural language corpus to recommendation-specific content that is sparse on long-tail publishers and thus learning effective word representations. CLRec [83] bridges the theoretical gap between contrastive learning objective and traditional recommendation objective, e.g., sampled softmax loss, as well as more advanced inverse propensity weighted

(IPW) loss. They show that directly performing contrastive learning can help to reduce exposure bias. CP4Rec [69] and S³-Rec [84] integrates Bert structure and contrastive learning objective for user pretraining, which require a fine-tuning stage. Compared with these works, we design model-agnostic and non-intrusive frameworks that help any baseline model learn more effective user representations in an end-to-end manner. Such representations are more accurate and robust by contrasting the original user representation with counterfactually positive samples and counterfactually negative samples.

2.3 Counterfactual for Recommendation

Causality and counterfactual reasoning have attracted great attentions in various domains [11, 78, 80]. Previous counterfactual frameworks in recommendation focus on debiasing the learning-to-rank problems. A rigorous counterfactual learning framework, *i.e.*, PropDCG [1], is proposed to overcome the distorting effect of presentation bias. The position bias and the clickbait issue are investigated in [2, 64] and [62], respectively. The Inverse Propensity Score [52, 70] method obtains unbiased estimation by sample re-weighting based on the likelihood of being logged. Another line of works encapsulates the uniform data into recommendation by learning imputation models [73], computing propensity [52], using knowledge distillation [13, 37], and directly modeling the uniform data [4, 28, 38, 50]. Different from these works, we focus on denoising user representation learning and considers the retrospect question, *i.e.*, ‘what the user representation would be if we intervene on the observed behavior sequence?’. Technically, we propose several counterfactual transformations based on the identification of indispensable/dispensable concepts and devise several contrasting objectives for learning accurate and robust user representations.

3 METHODS

3.1 Problem Formulation

In the view of sequential recommendation, datasets can be formulated as $\mathcal{D} = \{(x_{u,t}, y_{u,t})\}_{u=1,2,\dots,N, t=1,2,\dots,T_u}$, where $x_{u,t} = \{y_{u,1:(t-1)}\}$ denotes a user’s historical behaviors prior to the t th behavior $y_{u,t}$ and arranged in a chronological order, and T_u denotes the number of behaviors for the user u . The goal of sequential recommendation is to predict the next item $y_{u,t}$ given the historical behaviors $x_{u,t}$, which can be formulated as modeling the probability of all possible items:

$$p(y_{u,t} = y|x_{u,t}), \quad (1)$$

We will drop the sub-scripts occasionally and write (x, y) in place of $(x_{u,t}, y_{u,t})$ for simplicity. Let \mathcal{X} denote the set of all possible click sequences, *i.e.* $x \in \mathcal{X}$ and each $y \in \mathcal{Y}$ represent a clicked item, while \mathcal{Y} is the set of all possible items.

Since the number of items $|\mathcal{Y}|$ can easily reach 100 million, industrial recommender systems consist typically of two phases, *i.e.*, the matching phase and the ranking phase, due to concerns on system latency. The matching (also called deep candidate generation) phase focuses on retrieving Top N candidates for each user, while the ranking phase further sorts the N candidates by typically considering more fine-grained user/item features and incorporating complex modeling architectures. In this paper, we mainly conduct

experiments in the matching stage (*e.g.*, comparing with SOTA matching models).

3.2 A Naive Matching Baseline

The paradigm of a matching model includes a user encoder $f_\theta(x) \in \mathbb{R}^d$, which takes the user’s historical behavior sequence as input and output one (or more) dense vector representing the user’s interests, and an item encoder $g_\theta(y) \in \mathbb{R}^d$, that represents the items in the same vector space as the user encoder. We denote all the trainable parameters in the system as θ , which includes the parameters in f_θ and g_θ . With the learned encoders and the extracted item vectors, *i.e.*, $\{g_\theta(y)\}_{y \in \mathcal{Y}}$, a k-nearest-neighbor search service, *e.g.*, Faiss [27], will be deployed for Top-N recommendation. Specifically, at serving time, an arbitrary user behavior sequence x will be transformed into a vectorial representation $f_\theta(x)$ and top N items with the largest matching scores will be retrieved as Top-N candidates. Such matching scores are typically computed as inner product $\phi_\theta(x, y) = \langle f_\theta(x), g_\theta(y) \rangle$ or cosine similarity. In a nutshell, the learning procedure can be formulated as the following maximum likelihood estimation:

$$\arg \min_{\theta} \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} -\log p_\theta(y | x), \quad (2)$$

$$\text{where } p_\theta(y | x) = \frac{\exp \phi_\theta(x, y)}{\sum_{y' \in \mathcal{Y}} \exp \phi_\theta(x, y')}. \quad (3)$$

In the matching phase, it can be infeasible to sum over all possible items y' as in the denominator. Here we adopt the sample softmax objective [3, 25]. To demonstrate the effectiveness of the proposed CauseRec architecture, we utilize a naive framework as the baseline. Specifically, the *item* encoder $g_\theta(y)$ is a plain lookup embedding matrix where the n th vector represents the item embedding with item id n . The *user* encoder $f_\theta(x)$ aggregates the embeddings of historically interacted items using global average pooling and then transforms the aggregated embedding into the same embedding space as item embeddings using multi-layer perceptrons (MLP):

$$f_\theta(x) = \text{MLP}\left(\frac{1}{t-1} \sum_{i=1}^{t-1} g_\theta(y_i)\right). \quad (4)$$

3.3 The CauseRec Architecture

In this section, we give a brief illustration on the intuition and overall schema/pipeline of the CauseRec architecture, which is depicted in Figure 2, and introduce the building blocks in detail.

3.3.1 Overall Schema. The essence of CauseRec is to answer the retrospect question, ‘what the user representation would be if we intervene on the observed behavior sequence?’ The counterfactual transformation in CauseRec relates to the ‘intervention on the observed behavior sequence.’ For answering ‘what the user representation would be,’ we introduce an important inductive bias that makes the intervention work as expected. Specifically, we first identify indispensable/dispensable concepts in the historical behavior sequence. An indispensable concept indicates a subset of one behavior sequence that can jointly represent a meaningful aspect of the user’s interest. A dispensable concept indicates a noisy subset that is less meaningful/important in representing an aspect of interest. We introduce the details in Section 3.3.2.

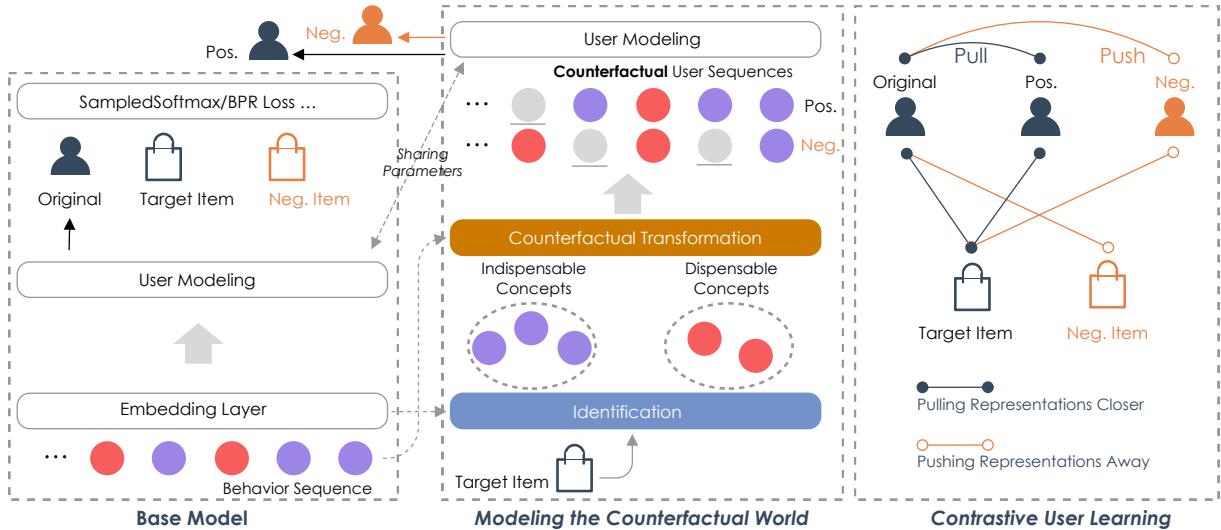


Figure 2: Schematic of the proposed CauseRec-Item framework.

Given the identified concepts, a representative counterfactual transformation is designed to build out-of-distribution counterfactual user sequences. Here comes the inductive bias, *i.e.*, counterfactual sequences constructed by replacing the dispensable concepts in the original user sequence should still have similar semantics to the original one. Here *semantics* refer to the characteristics of user interests/preferences. Therefore, replacing indispensable concepts in the original user sequence should incur a preference deviation from the resulted user representation to the original user representation. We denote these resulted user representations as counterfactually *negative* user representations. We note that such negatives are hard negatives where *hard* refers to that other dispensable concepts stay the same as the original user sequence and *negatives* means that the semantics of the user sequence should be different. In contrast, replacing dispensable concepts in the original user sequence should incur no preference change in representations. We denote these resulted user representations as counterfactually *positive* user representations. Different contrastive learning objectives are proposed to learn accurate and robust user representations that are less sensitive to the (attack on the) dispensable/noisy concepts and that trust more on the indispensable concepts that better represent the user's interest. Details can be found in Section 3.3.5.

3.3.2 Identification of Indispensable/Disposable concepts. To identify indispensable/disposable concepts, we propose to first extract concept proposals and compute the proposal scores.

Item-level Concepts. Inspired by instance discrimination [17], a straightforward while workable solution is to treat each item in the behavior sequence as an individual concept since each item has its unique fine-grained characteristics. In this way, we obtain the concept sequence $C = X \in \mathbb{R}^{t \times d}$, where $X = g_\theta(x_{u,t+1})$ denotes the vectorial representations of the behavior sequence. In essence, concept scores indicate to what extent these concepts are important to represent the user's interest. Since there is no groundtruth for

one user's real interest, we use the target item y_{t+1} as the indicator:

$$p_i^{item} = \phi_\theta(c_i, y), \quad (5)$$

where c_i indicates the representation of i th concept in C , and y indicates the representation of the target item. ϕ_θ is the similarity function, and we empirically use dot product for its effectiveness in the experiment. p_i^{item} is thus the score for the i th concept.

Interest-level Concepts. However, such a solution may incur redundancy in concepts since some items may share similar semantics, and might deteriorate the capability of modeling higher-order relationships between items. To this end, besides the item-level concepts, we introduce interest-level concepts by leveraging the attention mechanism [56] to extract interest-level concepts. Formally, $X \in \mathbb{R}^{t \times d}$, we obtain the following attention matrix:

$$A = \text{softmax}(\mathbf{W}_2 \tanh(\mathbf{W}_1 X^\top))^\top, \quad (6)$$

where $\mathbf{W}_1 \in \mathbb{R}^{d_a \times d}$ and $\mathbf{W}_2 \in \mathbb{R}^{K \times d_a}$ are trainable transformation matrices. K is thus the number of concepts that is pre-defined. A is of shape $\mathbb{R}^{t \times K}$. We obtain the concept sequence as the following:

$$C = A^\top X, \quad (7)$$

Since interest-level concepts and the target item are not naturally embedded in the same space, we compute the concept score by the weighted sum of item-level scores:

$$p^{interest} = A^\top \phi_\theta(X, y). \quad (8)$$

For both item-level concepts and interest-level concepts, we treat the top half concepts with the highest scores as indispensable concepts and the remaining half concepts as dispensable concepts. This strategy is mainly designed to prevent the number of indispensable or dispensable concepts from being too small. We leave finding more effective solutions as future works as illustrated in Section 5.

3.3.3 Counterfactual Transformation. The proposed counterfactual transformation aims to construct out-of-distribution user sequences by transforming the original user sequence for one user. We here

use ***user sequence*** to generally denote the concept sequence, which can be either the commonly known item-level concept sequence, *i.e.*, the original user behavior sequence, or the interest-level concept sequence. Based on the inductive bias described in Section 3.3.1, we propose to replace the identified indispensable/dispensable concepts at the rate of r_{rep} to construct counterfactually negative/positive user sequences, respectively. We note that directly dropping indispensable/dispensable concepts also seems feasible, but replacement has the advantage of not affecting overall sequence length and the relative positions of remaining concepts. We maintain a first-in-first-out queue as a concept memory for each level and use dequeued concepts as substitutes. We enqueue the concepts extracted from the current mini-batch. We denote the user sequence with indispensable/dispensable concepts being replaced as counterfactually negative/positive user sequence.

3.3.4 User Encoders. We note that item-level concept representations can be trained along with the item embedding matrix in most recommendation frameworks. Therefore, the above item-level concept identification and counterfactual transformation processes can be performed without any modification on the user encoder in the original baseline model, *i.e.*, a model-agnostic and non-intrusive design. We denote the architecture solely considering item-level concepts as CauseRec-Item. CauseRec-Item obtains counterfactually positive/negative user representations $\{\mathbf{x}^{+,m}\}_{m=1,\dots,M}/\{\mathbf{x}^{-,n}\}_{n=1,\dots,N}$ from counterfactual item-level concept sequences using the original user encoder f_θ .

We denote the architecture solely considering interest-level concepts as CauseRec-Interest. Different from CauseRec-Item, interest-level concepts are constructed with learnable parameters, *i.e.*, \mathbf{W}_1 and \mathbf{W}_2 in Equation 6. Therefore, CauseRec-Interest is an intrusive design, and the inputs to the user encoder should be the interest-level concept sequence rather than the behavior sequence at the item-level. We note that there are no further modifications, and the architecture of the user encoder can stay the same as in the original baseline model. CauseRec-Interest obtains counterfactually positive/negative user representations from counterfactual interest-level concept sequence using the original user encoder f_θ .

We denote the architecture that considers counterfactual transformation on both the item-level concept sequence and the interest-level concept sequence as CauseRec-H(ierarchical). CauseRec-H is also an intrusive design with interest-level concepts as the inputs of the user encoder. Different from CauseRec-Interest, CauseRec-H further considers counterfactual transformations performed on item-level concepts. The counterfactually transformed item-level sequence will be forwarded to construct interest-level concept sequence using Equation 6-7. We note that counterfactual transformations will not be performed on these two levels simultaneously, which might introduce unnecessary noises. In other words, each counterfactual user representation is constructed with transformation on sequence solely from one level.

3.3.5 Learning Objectives. Besides the original recommendation loss $\mathcal{L}_{matching}$ described near Equation 2, we propose several contrastive learning objectives that are especially designed for learning accurate and robust user representations.

Contrast between Counterfactual and Observation. As discussed in Section 3.3.1, a ***robust*** user representation should be less sensitive to (possible attack on) dispensable concepts. Therefore, the user representations learned from counterfactual sequences with indispensable concepts transformed should be intuitively pushed away from the original user representation. Similar in spirit, an ***accurate*** representation should trust more on indispensable concepts. Therefore, user representations learned from counterfactual sequences with dispensable concepts transformed should be intuitively pulled closer to the original user representation. Under these intuitions, we derive inspiration from the recent success of contrastive learning in CV [17, 68] and NLP [7], we use triplet margin loss to measure the relative similarity between samples:

$$\mathcal{L}_{co} = \sum_{m=1}^M \sum_{n=1}^N \max \{d(\mathbf{x}^q, \mathbf{x}^{+,m}) - d(\mathbf{x}^q, \mathbf{x}^{-,n}) + \Delta_{co}, 0\}, \quad (9)$$

where \mathbf{x}^q denotes the original user representation. We set the distance function d as the L2 distance since user representations generated by the same user encoder are in the same embedding space. We empirically set the margin $\Delta_{co} = 1$.

Contrast between Interest and Items. The above objective considers the user representation side solely, and we further capitalize on the target item y_t , which also enhances the user representation learning. Formally, given the L2-normalized representation of the target item $\tilde{\mathbf{y}}$ and user representation $\tilde{\mathbf{x}}$, we have:

$$\mathcal{L}_{ii} = \sum_{m=1}^M 1 - \tilde{\mathbf{x}}^{+,m} \cdot \tilde{\mathbf{y}} + \sum_{n=1}^N \max (0, \tilde{\mathbf{x}}^{-,n} \cdot \tilde{\mathbf{y}} - \Delta_{ii}), \quad (10)$$

This objective also has the advantage of preventing the user encoder from learning trivial representations for counterfactual user sequences. We set the margin $\Delta_{ii} = 0.5$ in the experiment. Finally, the loss for training the whole framework can be written as:

$$\mathcal{L}_{cause} = \mathcal{L}_{matching} + \lambda_1 \mathcal{L}_{co} + \lambda_2 \mathcal{L}_{ii}. \quad (11)$$

During testing/serving, only the backbone model that generates the user representation is needed. The identification of indispensable/dispensable concepts and the counterfactual transformation processes are disregarded. Noteworthy, the computation of proposal scores which depends on the target item does not belong to the backbone model and is not required during testing.

4 EXPERIMENTS

We conduct experiments on real-world public datasets and mainly aim to answer the following three research questions:

- **RQ1:** How does CauseRec perform compared to the base model and various SOTA sequential recommenders?
- **RQ2:** How do the proposed building blocks and different hyper-parameter settings affect CauseRec?
- **RQ3:** How do user representations benefit from modeling the counterfactual world and contrastive representation learning?

4.1 Experimental Setup

To demonstrate the generalization capability on learning users' representations of the proposed CauseRec architecture, we employ an

Table 1: Statistics of the Datasets.

Dataset	#Users	#Items	#Interactions	#Density
Amazon Books	459, 133	313, 966	8, 898, 041	0.00063
Yelp	31, 668	38, 048	1, 561, 406	0.00130
Gowalla	52, 643	91, 599	2, 984, 108	0.00084

evaluation framework [5, 35, 41] where models should confront unseen user behavior sequences. Specifically, the users of each dataset are split into training/validation/test subset by the proportion of 8 : 1 : 1. For training sequential recommenders, we incorporate a commonly used setting by viewing each item in the behavior sequence as a potential target item and using behaviors that happen before the target item to generate the user’s representation, as defined in Section 3.1. For evaluation, only users in the validation/test set are considered, and we choose to generate users’ representations on the first 80% behaviors, which are unseen during training. Such a framework can help justify whether models can learn accurate and robust user representations that can generalize well. We mainly focus on the *matching* phase of recommendation and accordingly choose the datasets, comparison methods, and evaluation metrics.

Datasets We consider three challenging recommendation datasets, of which the statistics are shown in Table 1.

- **Amazon Books.** We take Books category from the product review datasets provided by [44], for evaluation. For each user, we keep at most 20 behaviors that are chronologically ordered.
- **Yelp2018.** Yelp challenge (2018 edition) releases the review data for small businesses (*e.g.*, restaurants). We view these businesses as items and use a 10-core setting [20, 65] where each item/user has at least ten interactions.
- **Gowalla.** A widely used check-in dataset [34] from the Gowalla platform. Similarly, we use the 10-core setting [19].

Comparison Methods We mainly consider sequential recommenders for comparison since models are required to confront unseen behaviors for each user. Therefore, factorization-based and graph-based methods are not considered. The compared state-of-the-art models are listed as the following:

- **POP.** A naive baseline that always recommends items with the most number of interactions.
- **YouTube DNN [8].** A successful industrial recommender that generates one user’s representation by pooling the embeddings of historically interacted items.
- **GRU4Rec [21].** An early attempt to introduce recurrent neural networks into recommendation.
- **MIND [32].** The first framework that extracts multiple interest vectors for one user based on the capsule network.
- **ComiRec-DR [5].** A recently proposed SOTA framework following MIND to extract diverse interests using dynamic routing and incorporate a controllable aggregation module to balance recommendation diversity and accuracy.
- **ComiRec-SA [5].** ComiRec-SA differs from ComiRec-DR by using self-attention to model interests.

Evaluation Metrics We employ three broadly used numerical criteria for the matching phase, *i.e.*, *Recall*, *Normalized Discounted*

Cumulative Gain (NDCG), and *Hit Rate*. We report metrics computed on the top 20/50 recommended candidates. Higher values indicate better performance for all metrics.

Implementation Details We use Adam [30] for optimization with learning rate of 0.003/0.005 for Books/Yelp and Gowalla, $\beta_1 = 0.9$, $\beta_2 = 0.99$, $\epsilon = 1 \times 10^{-8}$, weight decay of 1×10^{-5} . We train CauseRec-Item for (maximum) 10 epochs and CauseRec-Interest/CauseRec-H for (maximum) 30 epochs with mini-batch size 1024. All models are with embedding size 64. We set hyper-parameters $\lambda_1 = \lambda_2 = 1$ and do not tune them with bells and whistles. As illustrated in Section 3.2, the item encoder is a plain embedding lookup matrix, and the user encoder is a three-layer perceptron with hidden size 256. We set $N = 8$, $M = 1$, $r_{rep} = 0.5$ for CauseRec-Item/-Interest and $N = 16$, $M = 2$ for CauseRec-H to accommodate transformation on two levels, as illustrated in Section 3.3.4. We set $K = 20$ for CauseRec-Interest/-H.

4.2 Performance Analysis (RQ1)

The comparison results of CauseRec with SOTA sequential recommenders are listed in Table 2. We report three architectures of CauseRec including CauseRec-Item (CauseItem), CauseRec-Interest (CauseIn), and CauseRec-Hierarchical (CauseH), as described in Section 3.3.4. In a nutshell, we observe a clear improvement of these architectures over various comparison methods and across three different metrics. Notably, CauseRec-H improves the previous SOTA ComiRec-SA/DR by +.0299 (relatively 22.1%) concerning NDCG@50 on the Amazon Books dataset and +.0179 (relatively 8.64%) concerning Recall@20 on the Gowalla dataset. Among the comparison methods, ComiRec mostly yields the best performance by modeling multiple interests for a given user. However, only modeling the noisy historical behaviors might result in diverse but noisy interests that may not accurately represent users, finally leading to inferior results. GRU4Rec achieves comparably good results with ComiRec on the Gowalla dataset. GRU4Rec can effectively model the sequential dependency between items in the behavior sequence. However, it might be more likely to suffer from the noises due to the strict step-by-step encoding process. In contrast, CauseRec architectures confront the noises within users’ behaviors by pushing the user representation away from counterfactually negative user representations and pulling it closer to counterfactually positive user representations. Besides, these results demonstrate the generalization capability of CauseRec on confronting out-of-distribution user sequences by modeling the counterfactual world.

Among three CauseRec architectures, CauseRec-Item is a model-agnostic and non-intrusive design, which means it can be applied to any other sequential recommender without any modification on the original user encoder, and solely functions in the training stage without sacrificing inference efficiency. CauseRec-Interest constructs interest-level concepts by grouping items that may belong to a certain interest (*e.g.*, chocolate and cake belong to sweets) into one holistic concept. Compared to CauseRec-Item, CauseRec-Interest has the advantage of reducing concept redundancy and modeling higher-order relationships between items, and thus improving CauseRec-Item. To combine the merits of CauseRec-Interest and CauseRec-Item, CauseRec-Hierarchical considers both

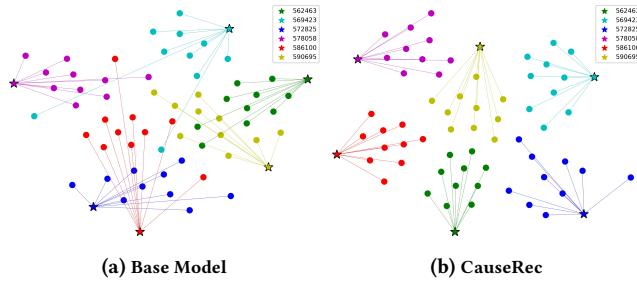


Figure 3: Visualization of randomly sampled users (shown as stars) with their interacted items (shown as points of the same color) from the Amazon Books dataset. We perform the t-SNE transformation on the representations learned by the base model (left) and CauseRec (right).

t-SNE transformation on the user/item representations learned by the base model (as shown in Figure 3a) and CauseRec-Item (as shown in Figure 3b). The connectivities of users and test items in the embedding space can help reflect whether the model learns accurate and robust user representations. From Figure 3b, we observe that users with their corresponding test items easily form clusters and show small intra-cluster distances and large inter-cluster distances. By jointly comparing the same users (e.g., 590695, and 586100) in Figures 3a and 3b, we can see that CauseRec-Item helps the user encoder learn representations that are closer to their corresponding test items. These results qualitatively demonstrate the effectiveness of CauseRec on learning accurate and robust user representations.

We also present a recommendation result from the Amazon Books test datasets in Figure 4. We list the historical behaviors, the top five books recommended by the base model and CauseRec-Item, and books interacted by the corresponding user in the test set. We mainly visualize the books' covers and categories for better clarity. We note that the side information is generally not considered in training matching models (both the base model and CauseRec). As shown in Figure 4, we observe that CauseRec yields more consistent recommendation results to the books in the test set. Supposing historical behaviors consist of noisy ones, and behaviors in the test accurately reflect users' interest for the current state, CauseRec successfully captures users' interests, i.e., Children's Books, and Literature&Fictions. In contrast, the base model is more likely to be affected by noisy behaviors that appear only a few times, such as the Biographies&Memories, and Education&Reference. These results further demonstrate that CauseRec can learn accurate and robust user representations that are less distracted by noisy behaviors.

5 CONCLUSION AND FUTURE WORK

In this work, we propose to model the counterfactual data distribution to confront the sparsity and noise nature of observed user interactions in recommender systems. The proposed CauseRec conditionally sample counterfactually positive and negative user sequences with transformations on the dispensable/indispensable concepts. We propose multiple structures (-item, -interest, -hierarchical) to confront both fine-grained item-level concepts and abstract interest-level concepts. Several contrastive objectives are devised to

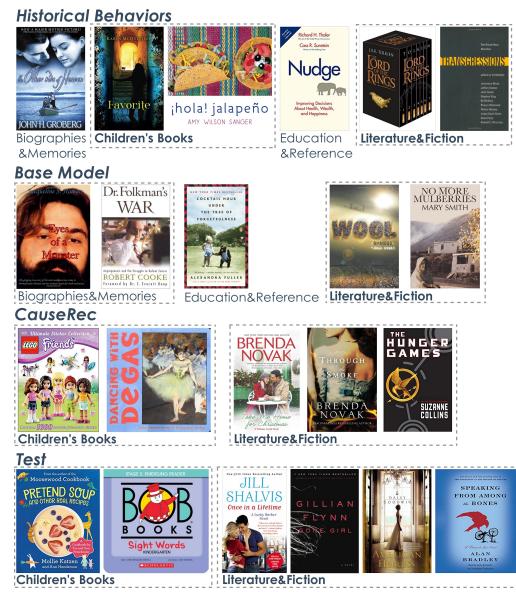


Figure 4: Case study by visualizing a real-world sample from the Amazon Books testing set. We mainly show the books' covers and categories for clarity.

contrast the counterfactual with the observational to learn accurate and robust user representations. Among several proposed architectures, CauseRec-Item has the advantage of being non-intrusive, i.e., solely functioning at training while not affecting serving efficiency. With a naive matching baseline, CauseRec achieves a considerable improvement over it and SOTA sequential matching recommenders. Extensive experiments help to justify the strengths of CauseRec as being both simple in design and effective in performance.

This work can be viewed as an initiative to exploit the joint power of constative learning and counterfactual thinking for recommendation. We believe that such a simple and effective idea can be inspirational to future developments, especially in model-agnostic and non-intrusive designs. CauseRec-Item is compatible with various user encoders within most existing sequential recommenders. We choose a naive baseline to better demonstrate the effectiveness of this work, and we plan to explore its strengths in more models. Another future direction is to whether more effective solutions of identifying indispensable/dispensable concepts exist, including both the computation of concept scores and the determination of indispensable or dispensable for each concept based on the scores. Lastly, we will explore the strengths of CauseRec for the ranking phase of recommendation. Counterfactual transformations designed with various auxiliary features and complex model architectures will open up new research possibilities.

6 ACKNOWLEDGMENTS

The work is supported by the National Key R&D Program of China (No. 2020YFC0832500), NSFC (61625107, 61836002, 62072397), Zhejiang Natural Science Foundation (LR19F020006), and Fundamental Research Funds for the Central Universities (2020QNA5024).

