

Sequential Recommendation via Stochastic Self-Attention

Ziwei Fan*, Zhiwei Liu,
 Yu Wang
 Department of Computer Science,
 University of Illinois at Chicago
 USA
 {zfan20,zliu213,ywang617}@uic.edu

Alice Wang, Zahra Nazari
 Spotify
 USA
 {alicew,zahran}@spotify.com

Lei Zheng
 Pinterest Inc
 USA
 lzheng@pinterest.com

Hao Peng
 School of Cyber Science and
 Technology, Beihang University
 China
 penghao@act.buaa.edu.cn

Philip S. Yu
 Department of Computer Science,
 University of Illinois at Chicago
 USA
 psyu@uic.edu

ABSTRACT

Sequential recommendation models the dynamics of a user’s previous behaviors in order to forecast the next item, and has drawn a lot of attention. Transformer-based approaches, which embed items as vectors and use dot-product self-attention to measure the relationship between items, demonstrate superior capabilities among existing sequential methods. However, users’ real-world sequential behaviors are *uncertain* rather than deterministic, posing a significant challenge to present techniques. We further suggest that dot-product-based approaches cannot fully capture *collaborative transitivity*, which can be derived in item-item transitions inside sequences and is beneficial for cold start items. We further argue that BPR loss has no constraint on positive and sampled negative items, which misleads the optimization.

We propose a novel STOchastic Self-Attention (STOSA) to overcome these issues. STOSA, in particular, embeds each item as a stochastic Gaussian distribution, the covariance of which encodes the uncertainty. We devise a novel Wasserstein Self-Attention module to characterize item-item position-wise relationships in sequences, which effectively incorporates uncertainty into model training. Wasserstein attentions also enlighten the collaborative transitivity learning as it satisfies triangle inequality. Moreover, we introduce a novel regularization term to the ranking loss, which assures the dissimilarity between positive and the negative items. Extensive experiments on five real-world benchmark datasets demonstrate the superiority of the proposed model over state-of-the-art baselines, especially on cold start items. The code is available in <https://github.com/zfan20/STOSA>.

*Part of this work is done in Spotify Research.

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.

WWW '22, April 25–29, 2022, Lyon, France.

© 2022 Association for Computing Machinery.
 ACM ISBN 978-1-4503-9096-5/22/04...\$15.00
<https://doi.org/10.1145/XXXXXX.XXXXXX>

CCS CONCEPTS

- Information systems → Collaborative filtering; Recommender systems; Personalization.

KEYWORDS

Sequential Recommendation, Transformer, Self-Attention, Uncertainty

ACM Reference Format:

Ziwei Fan*, Zhiwei Liu, Yu Wang, Alice Wang, Zahra Nazari, Lei Zheng, Hao Peng, and Philip S. Yu. 2022. Sequential Recommendation via Stochastic Self-Attention. In *Proceedings of the ACM Web Conference 2022 (WWW '22)*, April 25–29, 2022, Lyon, France. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/XXXXXX.XXXXXX>

1 INTRODUCTION

Recommender systems [6, 9, 25, 26, 44] become crucial components in web applications [23], which provide personalized item lists by modeling interactions between users and items. Sequential recommendation (SR) attracts a lot of attention from both the academic community and industry due to its success and scalability. SR methods format each user’s historical interactions as a sequence by sorting interactions chronologically. The goal of SR is to characterize users’ evolving interests and predict the next preferred item.

SR encodes users’ dynamic interests by modeling item-item transition relationships in sequences. Recent advancements in Transformer [24, 42] introduce the self-attention mechanism to reveal the position-wise item-item relationships, which leads to the state-of-the-art performance in SR. SASRec is the pioneering work in proposing Transformer for sequential recommendation, which applies scaled dot-product self-attention to learn item-item correlation weights. BERT4Rec [38] adopts bi-directional modeling in sequences. TiSASRec [21] and SSE-PT [45] extend SASRec with additional time interval information and user regularization, respectively.

Despite the success of self-attention in sequential recommendation, we argue that methods based on dot-product self-attention fail to incorporate: 1) *dynamic uncertainty* and 2) *collaborative transitivity*.

Firstly, existing SR methods assume that dynamic user interests are deterministic. As such, the inferred user embeddings are fixed vectors in the latent space, which are insufficient to represent multifarious user interests, especially in the real-world dynamic environment. Item transitions, which reflect the evolving process of user sequential behaviors, are sometimes hard to understand, and the two items in one item transition may not even lie in the same product category. As such, if a user has a significant portion of unexpected item transitions, modeling this user with a deterministic process achieves sub-optimal recommendations. For example, in books recommendation, a user interested in science-fiction, romance, and biography is more uncertain than another user interested in thriller, horror, and fantasy. Moreover, even two users share the same interest topics, the user with more fluctuated interests (e.g., items in item transitions are in different topics) is more uncertain. Intuitively, users with greater interest dynamic variability are more uncertain. Therefore, dynamic uncertainty is a crucial component when we model user interests in a sequential environment.

Another limitation of the existing self-attention mechanism is that it fails to incorporate *collaborative transitivity* in sequences. Collaborative transitivity can realize the latent similarity between items appearing in the same item-item transition pair but also inductively introduces additional collaborative similarities beyond limited item-item transitions in datasets. Thus, the collaborative transitivity can further alleviate cold-start item issues with the help of extra inducted collaborative similar items. For example, given item transition pairs $(i_x \rightarrow i_y)$ and $(i_y \rightarrow i_z)$, we can conclude that i_x and i_y are close to each other and so as for i_y and i_z . According to collaborative transitivity, i_x and i_z should also be close. However, existing dot-product self-attention is unable to realize this collaborative closeness. For example, given embeddings of $i_x = [0, 2]$, $i_y = [1, 1]$, $i_z = [2, 0]$, the dot-products of (i_x, i_y) and (i_y, i_z) are both 2, however, the dot-product between i_x and i_z is 0 because i_x and i_z transition pair is not observed. This issue becomes worse for unpopular items (*item cold start problem*) as the insufficient data for cold start items limits the set of collaborative neighbors.

However, it is rather challenging to resolve the dynamic uncertainty and incorporate the collaborative transitivity into a SR model. Firstly, characterizing dynamic uncertainty among item transition relationships is still under-explored. Most existing self-attention SR models, e.g., SASRec [17] and BERT4Rec [38], represent items as fixed vector embeddings, ignoring the uncertainty in sequential correlations. A recent work DT4SR [7] represents items as distributions, which proposes the mean and covariance embedding to model uncertainty in items. However, DT4SR is incapable of modeling dynamic uncertainty as it models item transition relationships via dot-product attention, which cannot incorporate such dynamic uncertainty.

On top of modeling dynamic uncertainty, inducing collaborative transitivity remains a challenge. Existing works [13, 17, 34, 38, 39] based on the dot-product fail to satisfy the triangle inequality and consequently cannot accomplish the collaborative transitivity. Different from dot product, distances typically satisfy triangle inequality¹, which transits additional collaborative closeness and benefits

¹ $d(i_x, i_z) \leq d(i_x, i_y) + d(i_y, i_z)$

a lot in item cold start issue. This assumption is theoretically supported in TransRec [10] and will also be empirically demonstrated in the following experimental results section. Although some metric learning frameworks [10, 22, 49] propose distance functions to guarantee the triangle inequality, none of them can model dynamic uncertainty as well as collaborative transitivity in the sequential setting. The choice of distance function is pivotal to collaborative transitivity modeling.

Moreover, we argue that collaborative transitivity optimized by standard Bayesian Personalized Ranking (BPR)² [34] loss fails to guarantee the dissimilarity between positive and negative items in BPR. BPR measures the difference between a user's preference scores on the positive item and a randomly sampled negative item. However, there is no guarantee that the positive item is further away from the negative item in the latent space. The incorporation of the distance between positive items and negatively sampled items is reasonable and necessary. Typically, negative items are sampled from items that the user never shows interest with, even if randomly sampled negative items are not necessarily hard negative items [47] due to data bias [2]. A proper way to sample negative items is an important topic in the recommendation, but it is beyond this paper's scope.

To this end, we propose a new framework **STOchastic Self-Attention (STOSA)**, which comprises of three modules: (1) stochastic embeddings, (2) Wasserstein self-attention, and (3) regularization term in BPR loss. Specifically, we model items as Gaussian distributions with stochastic embeddings, consisting of mean (for base interests) and covariance (for the variability of interests) embeddings. On top of stochastic embeddings, we propose to use distances to measure item transitions, which is originated from metric learning [16, 31]. We propose a novel Wasserstein Self-attention layer, which measures attentions as scaled Wasserstein distances between items. We also introduce a novel regularization term in BPR loss to consider the distance between positive items and negatively sampled items. The contributions of this work are as follows:

- To the best of our knowledge, STOSA is the first work proposing a Wasserstein Self-Attention to consider collaborative transitivity in SR.
- We introduce stochastic embeddings to measure both base interests and the variability of interests inherent in user behaviors and improve BPR loss for SR with an additional regularization for constraining the distance between positive items and negatively sampled items.
- STOSA outperforms the state-of-the-art recommendation methods. The experimental results also demonstrate the effectiveness of STOSA on cold start items.
- Several visualizations verify the effectiveness of Wasserstein self-attention over the traditional scaled dot-product self-attention and justify the improvements for cold start items by collaborative transitivity.

²We only discuss BPR loss in this paper because it is the most widely used ranking loss for the top-N recommendation, but the same issue also happens to other losses, such as Hinge Loss.

2 RELATED WORK

Several topics are closely related to our research problem. We first introduce some relevant works in the sequential recommendation, which is the primary problem setting in this paper. As we use distance rather than dot-product as the metric, recommendation methods with metric learning will also be discussed. Finally, we will introduce some relevant works about using distributions as representations.

2.1 Sequential Recommendation

Sequential Recommendation (SR) recommends the next item based on the chronological sequence of the user’s historical interactions. The fundamental idea of SR is learning sequential patterns within consecutive interacted items. One representative work of SR is FPMC [35]. With the inspiration of Markov Chain’s capability of learning item-item transitions, FPMC fuses the idea of Markov Chains with matrix factorization. FPMC learns the first-order item transition matrix, assuming that the next-item prediction is only relevant to the previous one item. Fossil [11] extends this idea and considers higher order of item transitions. Another line of SR uses convolutional neural networks for sequence modeling, such as Caser [39]. Caser regards the embedding matrix of items in the sequence as an image and applies convolution operators with the motivation of capturing local item-item transitions.

The advancements of sequence modeling developed in deep neural networks inspire the adoption of Recurrent Neural Network (RNN) [15, 27, 30, 33, 46, 48] and Self-Attention mechanisms into SR [17, 21, 38]. For example, GRU4Rec [14] proposes to use Gated Recurrent Units in the session-based recommendation. The success of self-attention-based Transformer [42] and BERT [4] inspires the community to investigate the possibility of self-attention in the sequential recommendation. Unlike Markov chain and RNN methods, self-attention utilizes attention scores from all item-item pairs in the sequence. SASRec [17] and BERT4Rec [38] both demonstrate the effectiveness of self-attention with the state-of-the-art performance in next-item recommendation.

2.2 Metric Learning for Recommendation

Metric learning explores a proper distance function to measure the dissimilarity between objects, such as Euclidean distance, Mahalanobis distance [29] and Graph distance [8]. A crucial property that differentiates distance and dot-product as metrics is that distances usually satisfy the triangle inequality. Triangle inequality, as an inductive bias [31] for distances, is useful when data sparsity issue exists [10]. One early work on metric learning for recommendation is CML [16]. CML proposes a hinge loss on minimizing the L2 distance between embeddings of the user and interacted items. LRML [40] then demonstrates the geometric restriction of CML and introduces a latent relation as a translation vector in the distance calculation. TransRec [10] borrows the idea of knowledge embedding and also develops a translation vector for the sequential recommendation. SML [22] is the state-of-the-art metric learning recommendation method. It introduces an additional item-centric metric and the adaptive margin on top of CML.

2.3 Distribution Representations

Representing objects (e.g., words, nodes, and items) as distributions has been attracting interest from the research community [1, 12, 37, 43]. Distribution representations introduce uncertainties and provide more flexibility compared with one single fixed embedding. DVNE [50] utilizes Gaussian distribution as the node embedding in graphs and proposes a deep variational model for higher-order proximity information propagation. TIGER [32] and [43] represent words as Gaussians, and TIGER also introduces Gaussian attention for better learning the entailment relationship among words. PMLAM [28] and DDN [49] both propose to use Gaussian distributions to represent users and items. DDN learns the mean and covariance embeddings with two neural networks. DT4SR [7] is the most relevant work, which represents items as distributions and learns mean and covariance with separate Transformers.

3 PRELIMINARIES AND DISCUSSIONS

In this section, we first formulate the SR problem and then introduce the self-attention mechanism for solving this problem.

3.1 Problem Definition

Given a set of users \mathcal{U} and items \mathcal{V} , and their associated interactions, we can sort the interacted items of each user $u \in \mathcal{U}$ chronologically in a sequence as $S^u = [v_1^u, v_2^u, \dots, v_{|S^u|}^u]$, where $v_i^u \in \mathcal{V}$ denotes the i -th interacted item in the sequence. The goal of SR is to recommend a top-N ranking list of items as the potential next items in a sequence. Formally, we should predict $p(v_{|S^u|+1}^{(u)} = v | S^u)$.

3.2 Self-Attention for Recommendation

Since we adopt the self-attention mechanism as the backbone of sequence encoder, we introduce it before proposing our model. The intuition of the self-attention mechanism is that items in sequences are correlated but of distinct importance to the items at different positions in a sequence. Specifically, given a user’s action sequence S^u and the maximum sequence length n , the sequence is first truncated by removing earliest items if $|S^u| > n$ or padded with 0s to get a fixed length sequence $s = (s_1, s_2, \dots, s_n)$. An item embedding matrix $\mathbf{M} \in \mathbb{R}^{|\mathcal{V}| \times d}$ is defined, where d is the number of latent dimensions. A trainable positional embedding $\mathbf{P} \in \mathbb{R}^{n \times d}$ is added to sequence embedding matrix as:

$$\hat{\mathbf{E}}_{S^u} = [\mathbf{m}_{s_1} + \mathbf{p}_{s_1}, \mathbf{m}_{s_2} + \mathbf{p}_{s_2}, \dots, \mathbf{m}_{s_n} + \mathbf{p}_{s_n}]. \quad (1)$$

Specifically, self-attention uses dot-products between items in the sequence to infer their correlations, which are as follows:

$$\text{SA}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}}\right)\mathbf{V}, \quad (2)$$

where $\mathbf{Q} = \hat{\mathbf{E}}_{S^u} \mathbf{W}^Q$, $\mathbf{K} = \hat{\mathbf{E}}_{S^u} \mathbf{W}^K$, and $\mathbf{V} = \hat{\mathbf{E}}_{S^u} \mathbf{W}^V$. As both \mathbf{Q} and \mathbf{K} use the same input sequence, the scaled dot-product component can learn the latent correlation between items. Additionally, other components in Transformer are utilized in SASRec, including the point-wise feed-forward network, residual connection, and layer normalization.

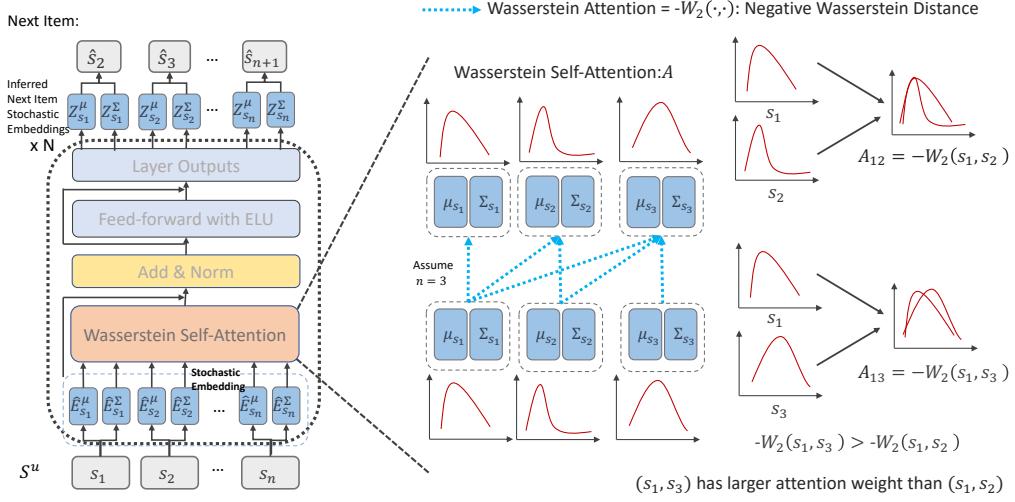


Figure 1: Model Architecture of the proposed STOSA. s_i denotes the item in the position i and \hat{s}_{i+1} indicates the output inferred next item in $(i+1)$ -th position. We propose stochastic embeddings to consider dynamic uncertainty information and introduce a novel Wasserstein Self-Attention layer for capturing collaborative transitivity signals. We introduce Feed-forward networks with ELU activation and guarantee the positive definite property of covariances.

4 PROPOSED MODEL

In this section, we introduce stochastic self-attention (STOSA) to overcome limitations of existing dot-product self-attention, as shown in Figure 1. We first represent items as stochastic embeddings with Elliptical Gaussian distributions, comprised of the mean embedding and covariance embedding. Then we develop a novel Wasserstein self-attention module based on the Wasserstein distance to infer the stochastic sequence embeddings. A Wasserstein distance is adopted to measure the dissimilarity between items in the sequence with uncertainty signals. Finally, we incorporate a novel regularization term measuring the distance between positive and negative items into the standard BPR loss.

4.1 Stochastic Embedding Layers

We introduce uncertainty into item embeddings by representing items as distributions. Differing deterministic vector representation, modeling items as stochastic distributions covers larger space for including more collaborative neighbors. Specifically, we use multi-dimensional elliptical Gaussian distributions to represent items. An elliptical Gaussian distribution is governed by a mean vector and a covariance vector³, where covariance introduces the potential uncertainty of the item. For all items, we define a mean embedding table $M^\mu \in \mathbb{R}^{|\mathcal{V}| \times d}$ and the covariance embedding table $M^\Sigma \in \mathbb{R}^{|\mathcal{V}| \times d}$. As mean and covariance identify different signals, we thus introduce separate positional embeddings for mean and covariance $P^\mu \in \mathbb{R}^{n \times d}$ and $P^\Sigma \in \mathbb{R}^{n \times d}$, respectively. In analogy to Eq. (1), we

can obtain mean and covariance sequence embeddings of user u as:

$$\begin{aligned}\hat{E}_{S^u}^\mu &= [\hat{E}_{s_1}^\mu, \hat{E}_{s_2}^\mu, \dots, \hat{E}_{s_n}^\mu] = [m_{s_1}^\mu + p_{s_1}^\mu, m_{s_2}^\mu + p_{s_2}^\mu, \dots, m_{s_n}^\mu + p_{s_n}^\mu], \\ \hat{E}_{S^u}^\Sigma &= [\hat{E}_{s_1}^\Sigma, \hat{E}_{s_2}^\Sigma, \dots, \hat{E}_{s_n}^\Sigma] = [m_{s_1}^\Sigma + p_{s_1}^\Sigma, m_{s_2}^\Sigma + p_{s_2}^\Sigma, \dots, m_{s_n}^\Sigma + p_{s_n}^\Sigma].\end{aligned}\quad (3)$$

For example, for the first item s_1 in the sequence, its stochastic embedding is represented as a d -dimensional elliptical Gaussian distribution $\mathcal{N}(\mu_{s_1}, \Sigma_{s_1})$, where $\mu_{s_1} = \hat{E}_{s_1}^\mu$ and $\Sigma_{s_1} = \text{diag}(\hat{E}_{s_1}^\Sigma) \in \mathbb{R}^{d \times d}$.

4.2 Wasserstein Self-Attention Layer

There remain challenges in modeling sequential dynamics with stochastic embeddings. First, it remains problematic to model dynamics of item transitions with distributions while still satisfying the triangle inequality. Secondly, the aggregation of these sequential signals to obtain the sequence's representation (i.e., the user's representation) is still not resolved. To tackle both challenges, we introduce Wasserstein distances as attention weights to measure the pair-wise relationships between items in the sequence, and we also adopt the linear combination property of Gaussian distributions [5] to aggregate historical items and obtain the sequence representation.

4.2.1 Wasserstein Attention. We propose a novel variant of self-attention adaptive to stochastic embeddings. We first denote $A \in \mathbb{R}^{n \times n}$ as the self-attention values. A_{kt} denotes the attention value between item s_k and item s_t in k -th and t -th positions in the sequence, where $k \leq t$ with the consideration of causality, respectively. According to Eq. (2), the attention weight of traditional self-attention is calculated as:

$$A_{kt} = Q_k K_t^\top / \sqrt{d}. \quad (4)$$

³The covariance of elliptical Gaussian distribution is a diagonal matrix, therefore the diagonal values can be viewed as a vector.

However, dot-product is not designed for measuring the discrepancy between distributions (i.e., stochastic embeddings) and fails to satisfy triangle inequality. Instead, we adopt Wasserstein distance⁴ [36] to measure the distance between stochastic embeddings of two items. Formally, given two items s_k and s_t , the corresponding stochastic embeddings are $\mathcal{N}(\mu_{s_k}, \Sigma_{s_k})$ and $\mathcal{N}(\mu_{s_t}, \Sigma_{s_t})$, where $\mu_{s_k} = \hat{\mathbf{E}}_{s_k}^\mu W_K^\mu$, $\Sigma_{s_k} = \text{ELU}\left(\text{diag}(\hat{\mathbf{E}}_{s_k}^\Sigma W_K^\Sigma)\right) + 1$, $\mu_{s_t} = \hat{\mathbf{E}}_{s_t}^\mu W_Q^\mu$, $\Sigma_{s_t} = \text{ELU}\left(\text{diag}(\hat{\mathbf{E}}_{s_t}^\Sigma W_Q^\Sigma)\right) + 1$. Exponential Linear Unit (ELU) maps inputs into $[-1, +\infty)$. It is used to guarantee the positive definite property of covariance. We define the attention weight as the negative 2-Wasserstein distance $W_2(\cdot, \cdot)$ is measured as follows,

$$\begin{aligned} \mathbf{A}_{kt} &= -(W_2(s_k, s_t)) \\ &= -\left(||\mu_{s_k} - \mu_{s_t}||_2^2 + \text{trace}\left(\Sigma_{s_k} + \Sigma_{s_t} - 2(\Sigma_{s_k}^{1/2} \Sigma_{s_k} \Sigma_{s_t}^{-1/2})^{1/2}\right)\right), \end{aligned} \quad (5)$$

Why Wasserstein distance? There are several advantages of using Wasserstein distance. First, Wasserstein distance measures the distance between distributions, with the capability of measuring the dissimilarity of items with uncertainty information. Secondly, Wasserstein distance satisfies triangle inequality [3] and can capture collaborative transitivity inductively in sequence modeling. Finally, Wasserstein distance also enjoys the advantage of a more stable training process as it provides a smoother measurement when two distributions are non-overlapping [19], which in SR means two items are far away from each other. However, KL divergence will produce an infinity distance, causing numerical instability.

Note that Eq. (5) can be computed with batch matrix multiplications without sacrificing computation and space efficiency compared with traditional self-attention, which will be discussed in the complexity analysis section of Appendices.

4.2.2 Wasserstein Attentive Aggregation. The output embedding of the item in each position of the sequence is the weighted sum of embeddings from previous steps, where weights are normalized attention values $\tilde{\mathbf{A}}$ as:

$$\tilde{\mathbf{A}}_{kt} = \frac{\mathbf{A}_{kt}}{\sum_{j=1}^t \mathbf{A}_{jt}}. \quad (6)$$

As each item is represented as a stochastic embedding with both mean and covariance, the aggregations of mean and covariance are different. We adopt the linear combination property of Gaussian distribution [5], which is as follows,

$$\mathbf{z}_{s_t}^\mu = \sum_{k=1}^t \tilde{\mathbf{A}}_{kt} \mathbf{v}_k^\mu, \text{ and } \mathbf{z}_{s_t}^\Sigma = \sum_{k=0}^t \tilde{\mathbf{A}}_{kt}^2 \mathbf{v}_k^\Sigma, \quad (7)$$

where $\mathbf{V}_{s_k}^\mu = \hat{\mathbf{E}}_{s_k}^\mu W_V^\mu$, $\mathbf{V}_{s_k}^\Sigma = \text{diag}(\hat{\mathbf{E}}_{s_k}^\Sigma) W_V^\Sigma$, and $k \leq t$ for causality. The outputs $\mathbf{Z}^\mu = (\mathbf{z}_{s_1}^\mu, \mathbf{z}_{s_2}^\mu, \dots, \mathbf{z}_{s_n}^\mu)$ and $\mathbf{Z}^\Sigma = (\mathbf{z}_{s_1}^\Sigma, \mathbf{z}_{s_2}^\Sigma, \dots, \mathbf{z}_{s_n}^\Sigma)$ together form the newly generated sequence's stochastic embeddings, which aggregates historical sequential signals with awareness of uncertainty.

⁴We also tried Kullback–Leibler (KL) divergence, but it achieves worse performance and inferior inference efficiency as well as violates triangle inequality.

4.3 Feed-Forward Network and Layer Outputs

The self-attention and the aggregation learn relationships in linear transformation. However, non-linearity can capture more complex relationships. We apply two point-wise fully connected layers with an ELU activation to introduce non-linearity in learning stochastic embeddings:

$$\begin{aligned} \text{FFN}^\mu(\mathbf{z}_{s_t}^\mu) &= \text{ELU}(\mathbf{z}_{s_t}^\mu W_1^\mu + b_1^\mu) W_2^\mu + b_2^\mu, \\ \text{FFN}^\Sigma(\mathbf{z}_{s_t}^\Sigma) &= \text{ELU}(\mathbf{z}_{s_t}^\Sigma W_1^\Sigma + b_1^\Sigma) W_2^\Sigma + b_2^\Sigma, \end{aligned} \quad (8)$$

where $W_1^* \in \mathbb{R}^{d \times d}$, $W_2^* \in \mathbb{R}^{d \times d}$, $b_1^* \in \mathbb{R}^d$, and $b_2^* \in \mathbb{R}^d$ are learnable parameters and $*$ can be μ or Σ . We adopt ELU instead of ReLU because of the numerical stability of ELU. We also adopt other components like [17, 38, 42], such as residual connection, layer normalization, and dropout layers, the layer outputs are,

$$\begin{aligned} \mathbf{Z}_{s_t}^\mu &= \mathbf{z}_{s_t}^\mu + \text{Dropout}(\text{FFN}^\mu(\text{LayerNorm}(\mathbf{z}_{s_t}^\mu))), \\ \mathbf{Z}_{s_t}^\Sigma &= \text{ELU}\left(\mathbf{z}_{s_t}^\Sigma + \text{Dropout}(\text{FFN}^\Sigma(\text{LayerNorm}(\mathbf{z}_{s_t}^\Sigma)))\right) + 1. \end{aligned} \quad (9)$$

We adopt ELU activation and ones addition to covariance embeddings to guarantee the positive definite property of covariance. Note that if we stack more layers, \mathbf{Z}^μ and \mathbf{Z}^Σ can be inputs of the next Wasserstein self-attention layer. We ignore the layer superscript for avoiding over-complex symbolization.

4.4 Prediction Layer

We predict the next item based on output embeddings \mathbf{Z}^μ and \mathbf{Z}^Σ from last layer if we stack several layers. We adopt the similar shared item embedding strategy in [17, 21, 38] for reducing model size and the risk of overfitting. Formally, for the item s_t in the t -th position of the sequence, the prediction score of next item j at $(t+1)$ -th position is formulated as 2-Wasserstein distance of two distributions $\mathcal{N}(\mu_{s_t}, \Sigma_{s_t})$ and $\mathcal{N}(\mu_j, \Sigma_j)$,

$$d_{s_t, j} = W_2(s_t, j), \quad (10)$$

where $\mu_{s_t} = \mathbf{Z}_{s_t}^\mu$ and $\Sigma_{s_t} = \mathbf{Z}_{s_t}^\Sigma$ are inferred representations given (s_1, s_2, \dots, s_t) , $1 \leq t \leq n$; $\mu_j = \mathbf{M}_j^\mu$ and $\Sigma_j = \mathbf{M}_j^\Sigma$ are embeddings indexed from input stochastic embedding tables \mathbf{M}^μ and \mathbf{M}^Σ .

For evaluation, different from dot-product methods, a smaller distance score indicates a higher probability of the next item. We thus generate the top-N recommendation list by sorting scores in ascending order.

4.5 BPR Loss with Positive v.s. Negative

We adopt the standard BPR loss [34] as base loss for measuring the ranking prediction error. However, BPR loss fails to consider the distance between the positive item and the negative sampled item. Therefore, we introduce a regularization term to enhance such distances as follows:

$$\ell_{pon}(s_t, j^+, j^-) = [d_{s_t, j^+} - d_{j^+, j^-}]_+, \quad (11)$$

where $[x]_+ = \max(x, 0)$ is the standard hinge loss, j^+ is the ground truth next item, and j^- is a randomly sampled negative item from items that the user never interacts with. The intuition behind $\ell_{pon}(t, j^+, j^-)$ is that the distance between positive item and negative item d_{j^+, j^-} has to be larger than the prediction distance d_{s_t, j^+} . Otherwise, when $d_{j^+, j^-} < d_{s_t, j^+}$, it becomes counter-intuitive. The

inequality $d_{j^+, j^-} < d_{s_t, j^+}$ indicates that the positive item j^+ is closer with negative item j^- than with s_t while s_t , as the previous item of j^+ , should have a smaller distance with j^+ instead. We incorporate this hinge loss as a regularization term with BPR loss into the final loss as follows,

$$L = \sum_{S^u \in \mathcal{S}} \sum_{t=1}^{|S^u|} -\log(\sigma(d_{s_t, j^+} - d_{s_t, j^-})) + \lambda \ell_{pvn}(s_t, j^+, j^-) + \beta \|\Theta\|_2^2. \quad (12)$$

We minimize L and optimize all learnable parameters Θ with Adam optimizer [18]. In the ideal case, the second term $\lambda \ell_{pvn}(s_t, j^+, j^-)$ becomes 0, which means s_t is close to j^+ but both s_t and j^+ are far away from j^- .

5 EXPERIMENTS

In this section, we validate the effectiveness of the proposed STOSA in several aspects by presenting experimental results and comparisons. The experiments answer the following research questions (RQs):

- **RQ1:** Does STOSA provide better recommendations than baselines?
- **RQ2:** What is the influence of Wasserstein self-attention and the regularization term ℓ_{pvn} ?
- **RQ3:** Does STOSA help the item cold start issue?
- **RQ4:** What is the distinction between dot-product and Wasserstein attention?
- **RQ5:** Why STOSA can alleviate the item cold-start issue?

5.1 Datasets

We evaluate the proposed STOSA on five public benchmark datasets from Amazon review datasets across various domains, with more than 1.2 million users and 63k items in total. Amazon datasets are known for high sparsity and have several categories of rating reviews. We use timestamps of each rating to sort interactions of each user to form the sequence. The latest interaction is used for testing, and the last second one is used for validation. Following [7, 10, 17, 21, 38, 45], we also adopt the 5-core settings by filtering out users with less than 5 interactions. We treat the presence of ratings as positive interactions. Details of datasets statistics and preprocessing steps are in the Appendices.

5.2 Evaluation

For each user, we sort the prediction scores calculated by Eq. (10) in ascending order to generate the top-N recommendation list. We **rank all items** instead of the biased sampling evaluation [20]. We adopt the standard top-N ranking evaluation metrics, Recall@N, NDCG@N, and Mean Reciprocal Rank (MRR). Recall@N measures the average number of positive items being retrieved in the generated top-N recommendation list for each user. NDCG@N extends Recall@N by also considering the positions of retrieved positive items in the top-N list. MRR measures the ranking performance in the entire ranking list instead of top-N. We report the averaged metrics over all users. We report the performances when $N = 1$ and $N = 5$.

5.3 Baselines

We compare the proposed STOSA with the following baselines in three groups. The first group includes static recommendation methods, which ignore the sequential order, including BPR [34] and LightGCN [13]. The second group of baselines consist of recommendation methods based on metric learning, including CML [16], SML [22]. The third group includes sequential recommendation methods: TransRec [10], Caser [39], SASRec [17], DT4SR [7], and BERT4Rec [38].

For all baselines, we search the embedding dimension in {64, 128}. As the proposed model has both mean and covariance embeddings, we only search for {32, 64} for STOSA for the fair comparison. More details of hyper-parameters grid search are in Appendices.

5.4 Overall Comparison (RQ1 and RQ2)

We compare the performance of all models in Table 1 and demonstrate the effectiveness of STOSA. We interpret the results with the following observations:

- STOSA obtains the best performance against all baselines in all metrics, especially in top-1 recommendation (Recall@1). The relative improvements range from 3.16% to 35.11% in all metrics, demonstrating the superiority of STOSA. We can also observe that the improvements are consistent in MRR for measuring the entire recommendation list, ranging from 5.68% to 11.54%. We attribute improvements to several factors of STOSA: (1). the distribution representations help expand the latent interaction space of items to better understand uncertainty and flexibility; (2). the consideration of collaborative transitivity enhances the discovery and induction of collaborative signals inherent in item-item transitions; (4). the newly introduced ℓ_{pvn} loss poses an additional constraint, which restrains the distances between positive items and negative items to be no larger than the ones of positive item transitions.
- Static methods, including BPRMF, LightGCN, CML, and SML, perform worse than sequential methods. This phenomenon verifies the necessity of temporal order information for the recommendation. Among all static methods, BPRMF, LightGCN achieve the best performances in different datasets. BPRMF and LightGCN perform better in Tools and Office datasets while achieving comparative results in other datasets. The reason for inferior performances of metric learning methods (CML and SML) might be the norm constraint of embeddings as both models will normalize all embeddings to have the norm of one.
- Among all sequential baselines, SASRec and DT4SR perform the best. DT4SR outperforms SASRec in three of five datasets, indicating the necessity of distribution representations and modeling uncertainty information in sequential recommendation. BERT4Rec fails to achieve satisfactory performance potentially due to the loss inconsistency between the adopted Cloze objective and the recommendation ranking task. The comparison between Caser and Transformer-based methods demonstrates the effectiveness of self-attention in sequential modeling for the recommendation.

Table 1: Overall Performance Comparison Table. The best and second-best results are bold and underlined, respectively. ‘OOM’ means the out-of-memory error. ‘Improve.’ is the relative improvement against the second-best baseline performance.

Dataset	Metric	BPRMF	LightGCN	CML	SML	TransRec	Caser	SASRec	DT4SR	BERT4Rec	STOSA	Improv.
Home	Recall@1	0.0029	0.0026	0.0025	0.0026	0.0018	OOD	<u>0.0046</u>	0.0029	0.0029	0.0053	+13.63%
	Recall@5	0.0096	0.0095	0.0076	0.0084	0.0063	OOD	<u>0.0127</u>	0.0129	0.0105	0.0133	+3.16%
	NDCG@5	0.0062	0.0060	0.0059	0.0056	0.0040	OOD	<u>0.0087</u>	0.0082	0.0067	0.0093	+6.76%
	MRR	0.0073	0.0071	0.0062	0.0061	0.0052	OOD	<u>0.0094</u>	0.0093	0.0092	0.0100	+6.17%
Beauty	Recall@1	0.0082	0.0064	0.0072	0.0069	0.0085	0.0112	0.0129	<u>0.0143</u>	0.0119	0.0193	+35.11%
	Recall@5	0.0300	0.0287	0.0249	0.0279	0.0321	0.0309	0.0416	<u>0.0449</u>	0.0396	0.0504	+12.15%
	NDCG@5	0.0189	0.0174	0.0184	0.0173	0.0204	0.0214	0.0274	<u>0.0296</u>	0.0257	0.0351	+18.45%
	MRR	0.0216	0.0203	0.0198	0.0191	0.0236	0.0231	0.0291	<u>0.0323</u>	0.0294	0.0360	+11.54%
Tools	Recall@1	0.0062	0.0071	0.0048	0.0055	0.0059	0.0056	<u>0.0103</u>	0.0103	0.0059	0.0120	+15.81%
	Recall@5	0.0216	0.0231	0.0129	0.0156	0.0210	0.0129	0.0284	<u>0.0289</u>	0.0189	0.0312	+7.85%
	NDCG@5	0.0139	0.0152	0.0096	0.0107	0.0134	0.0091	0.0194	<u>0.0196</u>	0.0123	0.0217	+11.04%
	MRR	0.0154	0.0170	0.0107	0.0118	0.0152	0.0106	0.0207	0.0206	0.0160	0.0226	+9.81%
Toys	Recall@1	0.0084	0.0077	0.0072	0.0102	0.0062	0.0089	0.0193	<u>0.0202</u>	0.0110	0.0240	+18.88%
	Recall@5	0.0301	0.0266	0.0249	0.0283	0.0222	0.0240	<u>0.0551</u>	0.0550	0.0300	0.0577	+4.86%
	NDCG@5	0.0194	0.0173	0.0154	0.0195	0.0143	0.0210	<u>0.0377</u>	0.0360	0.0206	0.0412	+9.45%
	MRR	0.0216	0.0200	0.0178	0.0210	0.0166	0.0221	0.0385	<u>0.0387</u>	0.0244	0.0415	+7.35%
Office	Recall@1	0.0073	0.0088	0.0096	0.0090	0.0100	0.0069	0.0198	<u>0.0206</u>	0.0137	0.0234	+13.59%
	Recall@5	0.0214	0.0226	0.0249	0.0190	0.0343	0.0302	<u>0.0656</u>	0.0630	0.0485	0.0677	+3.20%
	NDCG@5	0.0144	0.0157	0.0172	0.0140	0.0219	0.0186	<u>0.0428</u>	0.0421	0.0309	0.0461	+7.71%
	MRR	0.0162	0.0181	0.0191	0.0164	0.0263	0.0268	0.0457	0.0475	0.0408	0.0502	+5.68%

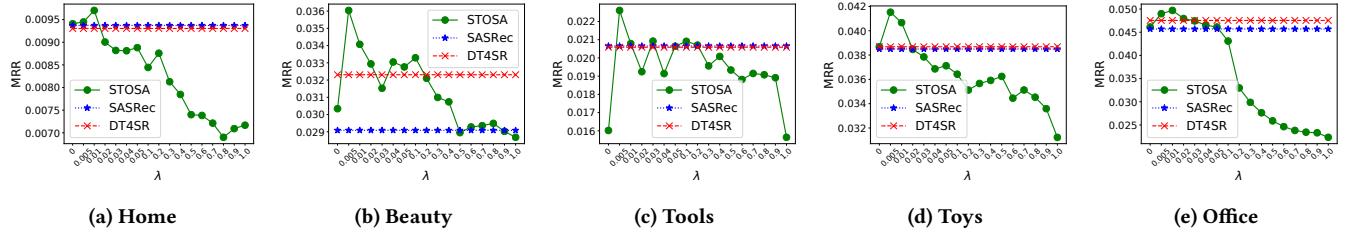


Figure 2: MRR performances over various λ on all datasets.

5.5 Parameters Sensitivity (RQ2)

In this section, we investigate the performance sensitivity of the weight λ on the additional regularization ℓ_{pan} across all datasets. Recall that the λ in Eq. (12) constraints the distance between the positive items and sampled negative items to be no less than the ground truth prediction distance. The trends are shown in Figure 2.

With a proper selection of λ , STOSA can perform better than SASRec and DT4SR. We can observe that as the values of λ become larger, the MRR performance first improves then drops. Another observation is that the values of λ can significantly affect the performance. A properly selected λ can dramatically improve the performance, indicating the necessity of the consideration of distances between positive items and sample negative items. However, λ should not be too large as negative items are not strictly negative rather than sampled negative items, which is the potential reason why the performance drops when λ increases.

One special case is setting $\lambda = 0$, which is also an ablation study of verifying the effectiveness of STOSA. When $\lambda = 0$, STOSA still outperforms SASRec in most datasets, except the Tools dataset. This indicates the superiority of Wasserstein Self-Attention and the

necessity of modeling uncertainty information and collaborative transitivity for the sequential recommendation.

5.6 Improvements Analysis (RQ3)

In this section, we analyze the sources of performance gains by comparing with SASRec on different groups of users and items. The analysis verifies the effectiveness of uncertainty information in user modeling and cold start items issue alleviation.

5.6.1 Performances w.r.t Sequence Lengths. We separate users into groups based on their number of interactions in the training portion, which is also the training sequence lengths of users. We report the average NDCG@5 on each group of users. Figure 3 shows the sizes of each group of users and the corresponding NDCG@5 performances. The group with the shortest sequence length has the most users, and the sizes decrease as the sequence lengths become longer.

From Figure 3, STOSA achieves the most significant improvements in users within the largest sequence length interval, compared with short sequences. The relative improvements on the longest sequence interval range from 9.70% to 54.45% across all

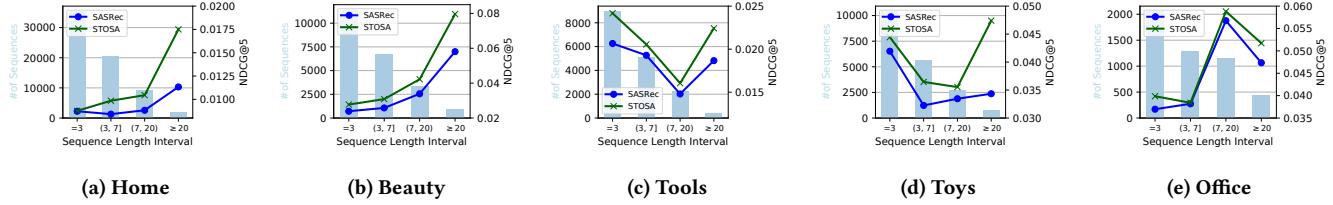


Figure 3: NDCG@5 performances on different sequence lengths (*i.e.*, number of training interactions of users) on all datasets.

datasets. The intuition behind these improvements is that users with more interactions are more likely to have diverse interests, indicating more uncertain behaviors. It demonstrates the effectiveness of stochastic representations in capturing uncertainty in user behaviors. We can also observe that STOSA can achieve comparative and better performances in most sequence length intervals, demonstrating the superiority of STOSA for the sequential recommendation.

5.6.2 Performances w.r.t Item Popularity. We investigate the performances on different groups of items based on the popularity to demonstrate that the collaborative transitivity and the proposed regularization ℓ_{pon} both help alleviate the cold-start items issue. We report the sizes and average NDCG@5 on each group of items, and each group is separated based on popularity. The distribution of sizes is similar to the one of users, where most items are unpopular items.

The performances comparison on all datasets is shown in Figure 4. For all datasets, the best improvements are from items with interactions no more than 3 (*i.e.*, cold start items). This observation supports the effectiveness of Wasserstein Self-Attention in capturing collaborative transitivity and the additional regularization ℓ_{pon} in generalizing the latent item transitions discovery. However, the performance becomes worse for popular items in Beauty and Toys datasets. We believe the reason might be the noisy neighbors of popular items in stochastic representations, which introduces potentially larger space for collaborative neighborhoods discovery.

5.7 Qualitative Analysis

In this section, we qualitatively visualize the attention weights and some examples of similar items retrieval. Specifically, we conduct a case study on a specific user with a long sequence length while her test item is a cold start item, which helps identify the significant difference between STOSA and SASRec. We also analyze the proposed STOSA by comparing the prediction lists with SASRec and STOSA, which can be found in Appendices.

5.7.1 Wasserstein Self-Attentions Visualization (RQ4). Figure 5 illustrates the heat maps of self-attention weights on the last 20 positions, learned by SASRec and STOSA, respectively. Recall that one of the critical differences between SASRec and STOSA is the calculation of attention weights, where SASRec adopts dot-product, and STOSA uses negative 2-Wasserstein distance.

We can observe some commonalities and differences between two attention weights heat maps. The attention weights of STOSA and SASRec are shown in Figure 5a and Figure 5b, respectively.

Both attention weights give larger weights to more recent behaviors, where the values in the bottom right corner are significant. However, STOSA has a more uniform attention weights distribution than SASRec as SASRec only highlights a small set of items in the sequence. The reason behind this difference is potentially the consideration of collaborative transitivity, which connects co-occurred items more tightly and introduces more collaborative neighbors in item-item transitions modeling.

5.7.2 Item Embeddings Visualization Comparison (RQ5). We use T-SNE [41] to visualize the latent spaces of items learned by SASRec and STOSA, respectively, as shown in Figure 6. We color items based on the items' popularity. We have the following observations: (1). the distributions of popular items (items with popularity more than 7) are significantly different between SASRec and STOSA, and (2). the distributions of cold items (popularity <= 3) are mostly uniform in both SASRec and STOSA. In SASRec, popular items are mostly located far away from the center and not closely connected to each other. STOSA learns a significantly different distribution of popular items, forcing them to locate close to the center and form a denser connected group. Moreover, due to the limited data for cold items, most cold items are uniformly distributed.

This difference is attributed to the collaborative transitivity learned in STOSA, which SASRec ignores and fails to generalize collaborative signals to cold items. The collaborative transitivity helps cold items to retrieve more collaborative neighbors as popular items can be related to triangle inequality. It again demonstrates the necessity and superiority of collaborative transitivity signals for the sequential recommendation.

6 CONCLUSION

This work proposes a novel stochastic self-attention sequential model STOSA for modeling dynamic uncertainty and capturing collaborative transitivity. We also introduce a novel regularization to BPR loss, guaranteeing a large distance between the positive item and negative sampled items. Extensive results and qualitative analysis on five real-world datasets demonstrate the effectiveness of STOSA and also well support the superiority of STOSA in alleviating cold start item recommendation issues.

7 ACKNOWLEDGMENTS

This work was supported in part by NSF under grants III-1763325, III-1909323, III-2106758, and SaTC-1930941. Hao Peng was supported by S&T Program of Hebei through grant 21340301D. For any correspondence, please refer to Ziwei Fan and Hao Peng.

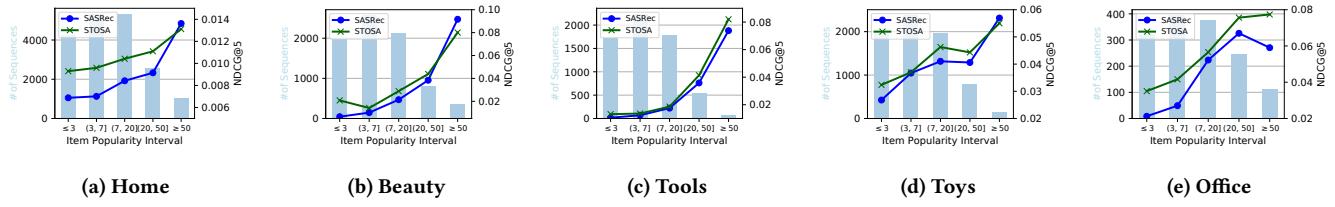


Figure 4: NDCG@5 performances on different item popularity (i.e., number of training interactions of items) on all datasets.

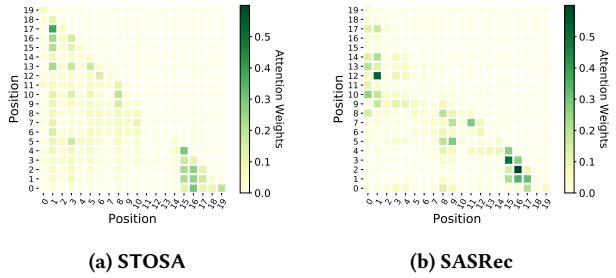


Figure 5: Attention Weights Visualizations of STOSA and SASRec on the user A278LEQK1TEPVB in Office dataset.

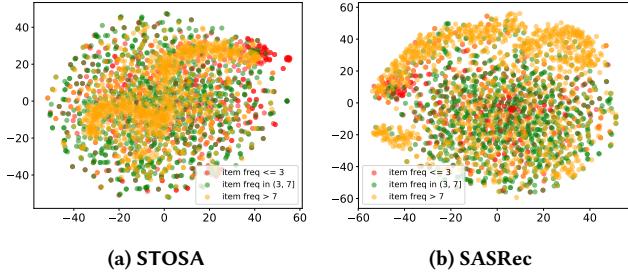


Figure 6: T-SNE Visualizations of item embeddings in Office dataset learned from STOSA and SASRec. The figures are best viewed in color.

REFERENCES

- [1] Aleksandar Bojchevski and Stephan Günnemann. 2018. Deep Gaussian Embedding of Graphs: Unsupervised Inductive Learning via Ranking. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=r1ZdkJ-0W>
- [2] Jiawei Chen, Hande Dong, Xiang Wang, Fuli Feng, Meng Wang, and Xiangnan He. 2020. Bias and debias in recommender system: A survey and future directions. *arXiv preprint arXiv:2010.03240* (2020).
- [3] Philippe Clement and Wolfgang Desch. 2008. An elementary proof of the triangle inequality for the Wasserstein metric. *Proc. Amer. Math. Soc.* 136, 1 (2008), 333–339.
- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [5] Paul S Dwyer. 1958. Generalizations of a Gaussian theorem. *The Annals of Mathematical Statistics* 29, 1 (1958), 106–117.
- [6] Ziwei Fan, Zhiwei Liu, Jiawei Zhang, Yun Xiong, Lei Zheng, and Philip S. Yu. 2021. *Continuous-Time Sequential Recommendation with Temporal Graph Collaborative Transformer*. Association for Computing Machinery, New York, NY, USA, 433–442. <https://doi.org/10.1145/3459637.3482242>
- [7] Ziwei Fan, Zhiwei Liu, Lei Zheng, Shen Wang, and Philip S Yu. 2021. Modeling Sequences as Distributions with Uncertainty for Sequential Recommendation. *arXiv preprint arXiv:2106.06165* (2021).
- [8] Xinbo Gao, Bing Xiao, Dacheng Tao, and Xuelong Li. 2010. A survey of graph edit distance. *Pattern Analysis and applications* 13, 1 (2010), 113–129.
- [9] Jibing Gong, Shen Wang, Jinlong Wang, Wenzheng Feng, Hao Peng, Jie Tang, and Philip S Yu. 2020. Attentional graph convolutional networks for knowledge concept recommendation in moocs in a heterogeneous view. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 79–88.
- [10] Ruining He, Wang-Cheng Kang, and Julian McAuley. 2017. Translation-based recommendation. In *Proceedings of the eleventh ACM conference on recommender systems*. 161–169.
- [11] Ruining He and Julian McAuley. 2016. Fusing similarity models with markov chains for sparse sequential recommendation. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*. IEEE, 191–200.
- [12] Shizhu He, Kang Liu, Guoliang Ji, and Jun Zhao. 2015. Learning to represent knowledge graphs with gaussian embedding. In *Proceedings of the 24th ACM international conference on information and knowledge management*. 623–632.
- [13] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgc: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*. 639–648.
- [14] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939* (2015).
- [15] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [16] Cheng-Kang Hsieh, Longqi Yang, Yin Cui, Tsung-Yi Lin, Serge Belongie, and Deborah Estrin. 2017. Collaborative metric learning. In *Proceedings of the 26th international conference on world wide web*. 193–201.
- [17] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 197–206.
- [18] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1412.6980>
- [19] Soheil Kolouri, Phillip E. Pope, Charles E. Martin, and Gustavo K. Rohde. 2019. Sliced Wasserstein Auto-Encoders. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=H1xaJn05FQ>
- [20] Walid Krichene and Steffen Rendle. 2020. On sampled metrics for item recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1748–1757.
- [21] Jiacheng Li, Yujie Wang, and Julian McAuley. 2020. Time interval aware self-attention for sequential recommendation. In *Proceedings of the 13th international conference on web search and data mining*. 322–330.
- [22] Mingming Li, Shuai Zhang, Fuqing Zhu, Wanhai Qian, Liangjun Zang, Jizhong Han, and Songlin Hu. 2020. Symmetric metric learning with adaptive margin for recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 4634–4641.
- [23] Xu Lin, Panagiotis Ilia, and Jason Polakis. 2020. Fill in the blanks: Empirical analysis of the privacy threats of browser form autofill. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*. 507–519.
- [24] Zhiwei Liu, Ziwei Fan, Yu Wang, and Philip S. Yu. 2021. *Augmenting Sequential Recommendation with Pseudo-Prior Items via Reversely Pre-Training Transformer*. Association for Computing Machinery, New York, NY, USA, 1608–1612. <https://doi.org/10.1145/3404835.3463036>
- [25] Zhiwei Liu, Xiaohan Li, Ziwei Fan, Stephen Guo, Kannan Acham, and S Yu Philip. 2020. Basket recommendation with multi-intent translation graph neural network. In *2020 IEEE International Conference on Big Data (Big Data)*. IEEE, 728–737.
- [26] Zhiwei Liu, Liangwei Yang, Ziwei Fan, Hao Peng, and Philip S Yu. 2021. Federated Social Recommendation with Graph Neural Network. *arXiv preprint arXiv:2111.10778* (2021).
- [27] Chen Ma, Peng Kang, and Xue Liu. 2019. Hierarchical gating networks for sequential recommendation. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. 825–833.

- [28] Chen Ma, Liheng Ma, Yingxue Zhang, Ruiming Tang, Xue Liu, and Mark Coates. 2020. Probabilistic metric learning with adaptive margin for top-K Recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1036–1044.
- [29] Geoffrey J McLachlan. 1999. Mahalanobis distance. *Resonance* 4, 6 (1999), 20–26.
- [30] Bo Peng, Zhiyun Ren, Srinivasan Parthasarathy, and Xia Ning. 2021. HAM: hybrid associations models for sequential recommendation. *IEEE Transactions on Knowledge and Data Engineering* (2021).
- [31] Silviu Pitis, Harris Chan, Kiarash Jamali, and Jimmy Ba. 2020. An Inductive Bias for Distances: Neural Nets that Respect the Triangle Inequality. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=HJeIDpVFP>
- [32] Chen Qian, Fuli Feng, Lijie Wen, and Tat-Seng Chua. 2021. Conceptualized and Contextualized Gaussian Embedding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 13683–13691.
- [33] Massimo Quadrana, Alexandros Karatzoglou, Balázs Hidasi, and Paolo Cremonesi. 2017. Personalizing session-based recommendations with hierarchical recurrent neural networks. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*. 130–137.
- [34] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2012. BPR: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618* (2012).
- [35] Steffen Rendle, Christoph Freudenthaler, and Lars Schmidt-Thieme. 2010. Factorizing personalized markov chains for next-basket recommendation. In *Proceedings of the 19th international conference on World wide web*. 811–820.
- [36] Ludger Rüschendorf. 1985. The Wasserstein distance and approximation theorems. *Probability Theory and Related Fields* 70, 1 (1985), 117–129.
- [37] Chi Sun, Hang Yan, Xipeng Qiu, and Xuanjing Huang. 2018. Gaussian word embedding with a wasserstein distance loss. *arXiv preprint arXiv:1808.07016* (2018).
- [38] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. BERT4Rec: Sequential recommendation with bidirectional encoder representations from transformer. In *Proceedings of the 28th ACM international conference on information and knowledge management*. 1441–1450.
- [39] Jiaxi Tang and Ke Wang. 2018. Personalized top-n sequential recommendation via convolutional sequence embedding. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*. 565–573.
- [40] Yi Tay, Luu Anh Tuan, and Siu Cheung Hui. 2018. Latent relational metric learning via memory-based attention for collaborative ranking. In *Proceedings of the 2018 World Wide Web Conference*. 729–739.
- [41] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).
- [42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [43] Luke Vilnis and Andrew McCallum. 2014. Word representations via gaussian embedding. *arXiv preprint arXiv:1412.6623* (2014).
- [44] Yu Wang, Zhiwei Liu, Ziwei Fan, Lichao Sun, and Philip S. Yu. 2021. *DSKReG: Differentiable Sampling on Knowledge Graph for Recommendation with Relational GNN*. Association for Computing Machinery, New York, NY, USA, 3513–3517. <https://doi.org/10.1145/3459637.3482092>
- [45] Liwei Wu, Shuqing Li, Cho-Jui Hsieh, and James Sharpnack. 2020. SSE-PT: Sequential recommendation via personalized transformer. In *Fourteenth ACM Conference on Recommender Systems*. 328–337.
- [46] An Yan, Shuo Cheng, Wang-Cheng Kang, Mengting Wan, and Julian McAuley. 2019. CosRec: 2D convolutional neural networks for sequential recommendation. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 2173–2176.
- [47] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. 2018. Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 974–983.
- [48] Lei Zheng, Ziwei Fan, Chun-Ta Lu, Jiawei Zhang, and Philip S Yu. 2019. Gated Spectral Units: Modeling Co-evolving Patterns for Sequential Recommendation. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1077–1080.
- [49] Lei Zheng, Chaozhuo Li, Chun-Ta Lu, Jiawei Zhang, and Philip S Yu. 2019. Deep Distribution Network: Addressing the Data Sparsity Issue for Top-N Recommendation. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1081–1084.
- [50] Dingyuan Zhu, Peng Cui, Daixin Wang, and Wenwu Zhu. 2018. Deep variational network embedding in wasserstein space. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2827–2836.

Appendices

A COMPLEXITY ANALYSIS

We analyze the space and time complexity of STOSA and demonstrate that STOSA has similar asymptotic space and time complexity as SASRec [17]. Note that even if stochastic embeddings are comprised of mean and covariance embeddings, we still use the same latent size as SASRec by equally separating dimensions to mean and covariance. For example, if we use $d = 128$ in SASRec, we use $d_\mu = d/2 = 64$ and $d_\Sigma = d/2 = 64$ for fair comparisons.

A.1 Space Complexity:

The learnable parameters in STOSA are from the stochastic embeddings and parameters in the Wasserstein self-attention layers, feed-forward networks and layer normalization. The overall number of parameters is $O(2|\mathcal{V}|\frac{d}{2} + 2n\frac{d}{2} + 2(\frac{d}{2})^2) = O(|\mathcal{V}|d + nd + d^2/2)$, which is slightly smaller than the complexity of SASRec, which is $O(|\mathcal{V}|d + nd + d^2)$ [17].

A.2 Time Complexity:

The computational complexity of STOSA is dominated by the Wasserstein self-attention layer and the feed-forward networks. The Wasserstein self-attention defined in Eq. (5) can be converted to using batch matrix multiplications. The second term in Eq. (5) can be transformed as a calculation of Euclidean norm as follows

$$\text{trace} \left(\Sigma_{s_t} + \Sigma_{s_k} - 2(\Sigma_{s_k}^{1/2} \Sigma_{s_t} \Sigma_{s_k}^{1/2})^{1/2} \right) = \|\Sigma_{s_t}^{1/2} - \Sigma_{s_k}^{1/2}\|_F^2, \quad (13)$$

where $\|\cdot\|_F^2$ is Frobenius norm that can be calculated by matrix multiplications. Also, as Σ_{s_t} and Σ_{s_k} are both diagonal matrices, we can further reduce the computational complexity to $O(\frac{nd}{2} + \frac{n^2d}{2} + 2n^2)$. And the Euclidean norm of mean embeddings part in Eq. (5) can also be calculated by matrix multiplications with the same time complexity. Therefore, the overall Wasserstein self-attention time complexity is $O(nd + n^2d + 4n^2)$. By also considering the feed-forward networks, we obtain the final asymptotic computational complexity as $O(nd + n^2d + 4n^2 + \frac{nd^2}{2})$. The computation complexity of traditional self-attention [17] is $O(n^2d + nd^2)$. Note that both complexities are typically dominated by the $O(n^2d)$ term as d is typically much larger than 4. It indicates that STOSA has asymptotic similarly time complexity with SASRec.

Table 2: Datasets Statistics

Dataset	#users	#items	#interactions	density	interactions per user	avg.
Home	66,519	28,237	551,682	0.03%	8.3	
Beauty	22,363	12,101	198,502	0.05%	8.3	
Toys	19,412	11,924	167,597	0.07%	8.6	
Tools	16,638	10,217	134,476	0.08%	8.1	
Office	4,905	2,420	53,258	0.44%	10.8	

B DATASETS AND PREPROCESSING

Details of datasets statistics⁵ are presented in Table 2.

C IMPLEMENTATION DETAILS AND BASELINES GRID SEARCH

We implement STOSA with Pytorch in a Nvidia 3090 GPU with 64GB system memory. We grid search all parameters and report the test performance based on the best validation results. For all baselines, we search the embedding dimension in {64, 128}. As the proposed model has both mean and covariance embeddings, we only search for {32, 64} for STOSA for the fair comparison. We also search max sequence length from {50, 100}. We tune the learning rate in $\{10^{-3}, 10^{-4}\}$, search the L2 regularization weight from $\{10^{-1}, 10^{-2}, 10^{-3}\}$, dropout rate from $\{0.3, 0.5, 0.7\}$. For sequential methods, we search number of layers from {1, 2, 3}, and number of heads in {1, 2, 4}. We adopt the early stopping strategy that model optimization stops when the validation MRR does not increase for 50 epochs. The followings are the model specific hyper-parameters search ranges of baselines: The third group consists of sequential recommendation methods:

- **BPR**⁶: BPR is the most classical collaborative filtering method for personalized ranking with implicit feedbacks. We search the learning rate in $\{10^{-3}, 10^{-4}\}$ and L2 regularization weight from $\{10^{-1}, 10^{-2}, 10^{-3}\}$.
- **LightGCN**⁷: LightGCN is the state-of-the-art static recommendation method, which considers high-order collaborative signals in user-item graph. We search number of layers from {1, 2, 3}, and node dropout from $\{0.1, 0.3, 0.5, 0.7\}$.
- **CML**⁸: One of the earliest works that adopt distance metrics to measure the affinity between users and items for recommendation. We search covariance loss weight from {0.1, 0.3}, and margin from {3.0, 4.0, 5.0}.
- **SML**⁹: This method is the state-of-the-art metric learning recommendation method. It extends CML to additionally consider item-centric distances. We search the γ from {5, 10, 20}, λ from {0.01, 0.1, 1.0}, and margin {0.1, 0.3, 0.5}.
- **TransRec**¹⁰: A metric learning-based sequential recommendation that proposes translation vectors to encode the item transition relationships. We search the λ from {0.001, 0.01, 0.1}.
- **Caser**¹¹: A CNN-based sequential recommendation method that views the sequence embedding matrix as an image and applies convolution operators to it. We search the length L from {5, 10}, and T from {1, 3, 5}.
- **SASRec**¹²: The state-of-the-art sequential method that depends on the Transformer architecture. We search the dropout rate from {0.3, 0.5, 0.7}.
- **DT4SR**¹³: A metric learning-base sequential method that models items as distributions and proposes mean and covariance Transformers. We search the dropout rate from {0.3, 0.5, 0.7}.

⁵<https://jmcauley.ucsd.edu/data/amazon/>

⁶https://github.com/xiangwang1223/neural_graph_collaborative_filtering

⁷<https://github.com/kuandeng/LightGCN>

⁸<https://github.com/changun/CollMetric>

⁹<https://github.com/MingmingLie/SML>

¹⁰<https://github.com/YifanZhou95/Translation-based-Recommendation>

¹¹https://github.com/graytowne/caser_pytorch

¹²<https://github.com/RUCAIBox/CIKM2020-S3Rec>

¹³<https://github.com/DyGRec/DT4SR>

- **BERT4Rec**¹⁴: This method extends SASRec to model bidirectional item transitions with standard Cloze objective. We search the mask probability from the range of {0.1, 0.2, 0.3, 0.5, 0.7}.

D QUALITATIVE ANALYSIS

D.1 Predictions Comparison

Table 3: Prediction lists comparison of user A278LEQK1TEPVB. The ground truth item is in red and each item is associated with its ID and the popularity. The names of items are described in the following second paragraph.

Model	SASRec	STOSA
Rank-1	(Wraparound Labels, 12)	(WF-7620 Printer, 3)
Rank-2	(WF-3640 Printer, 3)	(Wraparound Labels, 12)
Rank-3	(Locker, 3)	(WF-3640 Printer, 3)
Rank-4	(Binders, 1)	(Speaker Phone, 8)
Rank-5	(WF-7620 Printer, 3)	(Ring Binder, 10)

We show the differences of top-5 predicted ranking lists for the user A278LEQK1TEPVB between SASRec and STOSA in Table.3. The user’s last five interacted items are $\mathcal{S}[-5:] = \text{[Ink Refillable', 'Write 'N Wipe', 'Magnetic Whiteboard', 'Pencil Cup Holder', 'WF-4640 Inkjet Printer']}$. Moreover, STOSA can prioritize more relevant items even when items are cold start. The last interacted item of the user A278LEQK1TEPVB is a Inkjet Printer. The rank-1 item of STOSA WF-7620 Printer is a better version of Inkjet Printer while the rank-1 item of SASRec is a Wraparound Labels, which does not match with the user’s real interests. The reason might be that SASRec prefers more popular items and Wraparound Labels has more interactions than WF-7620 Printer. We can conclude that stochastic representations in STOSA can introduce different neighbors by utilizing collaborative transitivity information and also upvote relevant but cold start items in recommendation list.

¹⁴<https://github.com/FeiSun/BERT4Rec>