

Knowledge-aware Complementary Product Representation Learning

Da Xu*
Chuanwei Ruan*
Walmart Labs

Sunnyvale, California, USA
{Da.Xu,Chuanwei.Ruan}@walmartlabs.com

Jason Cho, Evren Korpeoglu, Sushant Kumar,
Kannan Achan
Walmart Labs
Sunnyvale, California, USA
{HCho,EKorpeoglu,SKumar4,KAchan}@walmartlabs.com

ABSTRACT

Learning product representations that reflect complementary relationship plays a central role in e-commerce recommender system. In the absence of the product relationships graph, which existing methods rely on, there is a need to detect the complementary relationships directly from noisy and sparse customer purchase activities. Furthermore, unlike simple relationships such as similarity, complementariness is asymmetric and non-transitive. Standard usage of representation learning emphasizes on only one set of embedding, which is problematic for modelling such properties of complementariness.

We propose using knowledge-aware learning with dual product embedding to solve the above challenges. We encode contextual knowledge into product representation by multi-task learning, to alleviate the sparsity issue. By explicitly modelling with user bias terms, we separate the noise of customer-specific preferences from the complementariness. Furthermore, we adopt the dual embedding framework to capture the intrinsic properties of complementariness and provide geometric interpretation motivated by the classic separating hyperplane theory. Finally, we propose a Bayesian network structure that unifies all the components, which also concludes several popular models as special cases.

The proposed method compares favourably to state-of-art methods, in downstream classification and recommendation tasks. We also develop an implementation that scales efficiently to a dataset with millions of items and customers.

CCS CONCEPTS

• **Information systems** → **Collaborative filtering**; **Personalization**; *Content ranking*; *Probabilistic retrieval models*.

KEYWORDS

Representation Learning; Dual Embedding; Complementary Product; Recommender System; User Modelling

*Both authors contributed equally to this 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.

WSDM '20, February 3–7, 2020, Houston, TX, USA

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-6822-3/20/02...\$15.00

<https://doi.org/10.1145/3336191.3371854>

ACM Reference Format:

Da Xu, Chuanwei Ruan, and Jason Cho, Evren Korpeoglu, Sushant Kumar, Kannan Achan. 2020. Knowledge-aware Complementary Product Representation Learning. In *The Thirteenth ACM International Conference on Web Search and Data Mining (WSDM '20)*, February 3–7, 2020, Houston, TX, USA. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3336191.3371854>

1 INTRODUCTION

Modern recommender systems aim at providing customers with personalized contents efficiently through the massive volume of information. Many of them are based on customers' explicit and implicit preferences, where content-based [30] and collaborative filtering systems [14] have been widely applied to e-commerce [17], online media [28] and social network [13]. In recent years representation learning methods have quickly gained popularity in online recommendation literature. Alibaba [34] and Pinterest [36] have deployed large-scale recommender system based on their trained product embeddings. Youtube also use their trained video embeddings as part of the input to a deep neural network [6]. Different from many other use cases, e-commerce platforms often offer a vast variety of products, and nowadays customers shop online for all-around demands from electronics to daily grocery, rather than specific preferences on narrow categories of products. Therefore understanding the intrinsic relationships among products [38] while taking care of individual customer preferences motivates our work. A topic modelling approach was recently proposed to infer complements and substitutes of products as link prediction task using the extracted product relation graphs of Amazon.com [18].

For e-commerce, product complementary relationship is characterized by co-purchase patterns according to customer activity data. For example, many customers who purchase a new *TV* often purchase *HDMI cables* next and then purchase *cable adaptors*. Here *HDMI cables* are complementary to *TV*, and *cable adaptors* are complementary to *HDMI cable*. We use \rightarrow as a shorthand notation for 'complement to' relationship. By this example, we motivate several properties of the complementary relationship:

- **Asymmetric.** *HDMI cable* \rightarrow *TV*, but *TV* \nrightarrow *HDMI cable*.
- **Non-transitive.** Though *HDMI cables* \rightarrow *TV*, and *cable adaptors* \rightarrow *HDMI cable*, but *cable adaptors* \nrightarrow *TV*.
- **Transductive.** *HDMI cables* are also likely to complement other *TVs* with similar model and brand.
- **Higher-order.** Complementary products for meaningful product combos, such as (*TV*, *HDMI cable*), are often different from their individual complements.

Although the above properties for complementary relationship highlight several components for designing machine learning algorithms, other manipulations are also needed to deal with the **noise** and **sparsity** in customer purchase activity data. Firstly, the low signal-to-noise ratio in customer purchase sequences causes difficulty in directly extracting complementary product signals from them. Most often, we only observe a few complementary patterns in customer purchase sequences or baskets. We list a customer’s single day purchase as a sequence for illustration:

$\{Xbox, games, T-shirt, toothbrush, pencil, notepad\}$.

Notice that among the fifteen pairs of possible item combinations, only two pairs can be recognized as complementary, i.e. $games \rightarrow Xbox, notepad \rightarrow pencil$. The rest purchases are out of the customer’s personal interest. As a solution to the noise issue, we introduce a customer-product interaction term to directly take account of the noises, while simultaneously modelling personalized preferences. We use \mathcal{U} to denote the set of users (customers) and \mathcal{I} to denote the set of items (products). Let $U \in \mathcal{U}$ be a user categorical random variable and $\{I_{t-1}, \dots, I_{t-k}\}$ be a sequence of item categorical random variables that represents k consecutive purchases before time t . If we estimate the conditional probability $p(I_{t+1}|U, I_t, \dots, I_{t-k})$ with *softmax* classifier under score function $S(\cdot)$, then previous arguments suggest that $S(I_{t+1}, (U, I_t, \dots, I_{t-k}))$ should consist of a user-item preference term and a item complementary pattern term:

$$S(I_{t+1}, (U, I_t, \dots, I_{t-k})) = f_{UI}(I_{t+1}, U) + f_I(I_{t+1}, I_t, \dots, I_{t-k}). \quad (1)$$

Now we have $f_{UI}(\cdot)$ accounting for user-item preference (bias) and $f_I(\cdot)$ characterizing the strength of item complementary pattern. When the complementary pattern of an purchase sequence is weak, i.e. $f_I(I_{t+1}, I_t, \dots, I_{t-k})$ is small, the model will enlarge the user-item preference term, and vice versa.

On the other hand, the sparsity issue is more or less standard for recommender systems and various techniques have been developed for both content-based and collaborative filtering systems [12, 22]. Specifically, it has been shown that modelling with contextual information boosts performances in many cases [1, 11, 19], which also motivates us to develop our context-aware solution.

Representation learning with shallow embedding gives rises to several influential works in natural language processing (NLP) and geometric deep learning. The *skip-gram* (SG) and *continuous bag of words* (CBOW) models [20] as well as their variants including *GLOVE* [23], *fastText* [5], *Doc2vec* (paragraph vector) [16] have been widely applied to learn word-level and sentence-level representations. While classical node embedding methods such as *Laplacian eigenmaps* [3] and *HOPE* [21] arise from deterministic matrix factorization, recent work like *node2vec* [10] explore from the stochastic perspective using random walks. However, we point out that the word and sentence embedding models target at semantic and syntactic similarities while the node embedding models aim at topological closeness. These relationships are all symmetric and mostly transitive, as opposed to the complementary relationship. Furthermore, the **transductive** property of complementarity requires that similar products should have similar complementary products, suggesting that product similarity should also be modelled

explicitly or implicitly. We approach this problem by encoding contextual knowledge to product representation such that contextual similarity is preserved.

We propose the novel idea of using context-aware dual embeddings for learning complementary products. While both sets of embedding are used to model complementarity, product similarities are implicitly represented on one of the embeddings by encoding contextual knowledge. Case studies and empirical testing results on the public *Instacart* and a proprietary e-commerce dataset show that our dual product embeddings are capable of capturing the desired properties of complementarity, and achieve cutting edge performance in classification and recommendation tasks.

2 CONTRIBUTIONS AND RELATED WORKS

Compared to the previously published work on learning product representations and modelling product relationships for e-commerce, our contributions are summarized below.

Learning higher-order product complementary relationship with customer preferences - Since product complementary patterns are mostly entangled with customer preference, we consider both factors in our work. And by directly modelling whole purchase sequences $f_I(I_{t+1}, I_t, \dots, I_{t-k})$, we are able to capture **higher-order** complementary relationship. The previous work of inferring complements and substitutes with topic modelling [18] and deep neural networks [37] on Amazon data relies on the extracted graph of product relationships. Alibaba proposes an approach by first constructing the weighted product co-purchase count graph from purchase records and then implementing a node embedding algorithm [34]. However, individual customer information is lost after the aggregation. The same issue occurs in item-based collaborative filtering [27] and product embedding [31]. A recent work on learning grocery complementary relationship for next basket predict models (item, item, user) triplets extracted from purchases sequences, and thus do not account for higher-order complementarity [33].

Use context-aware dual product embedding to model complementarity - Single embedding space may not capture **asymmetric** property of complementarity, especially when they are treated as projections to lower-dimensional vector spaces where inner products are symmetric. Although Youtube’s approach takes both user preference and video co-view patterns into consideration, they do not explore complementary relationship [6]. PinSage, the graph convolutional neural network recommender system of Pinterest, focus on symmetric relations between their pins [36]. The triplet model [33] uses two sets of embedding, but contextual knowledge is not considered. To our best knowledge, we are the first to use context-aware dual embedding in modelling product representations.

Propose a Bayesian network structure to unify all components, and conclude several classic models as special cases - As we show in Table 1, the classic collaborative filtering models for product recommendation [17], sequential item models such as *item2vec* [2] and *metapath2vec* [7] for learning product representation, the recent *triplet2vec* model [33] as well as the *prod2vec* [9] model which considers both purchase sequence and user, can all be

viewed as special cases within the Bayesian network representation of our model.

Besides the above major contributions, we also propose a fast inference algorithm to deal with the **cold start** challenge [15, 29], which is crucial for modern e-commerce platforms. Different from the cold-start solutions for collaborative filtering using matrix factorization methods [4, 40], we infer from product contextual features. The strategy is consistent with recent work, which finds contextual features playing an essential role in mitigating the cold-start problem [8, 28]. We also show that cold-start product representations inferred from our algorithm empirically perform better than simple context similarity methods in recommendation tasks.

3 METHOD

In this section, we introduce the technical details of our method. We first describe the Bayesian network representation of our approach, specify the factorized components and clarify their relationships. We then introduce the loss function in Section 3.1. In Section 3.2, we define the various types of embeddings and how we use them to parameterize the score functions. We then discuss two ranking criteria for recommendation tasks using embeddings obtained from the proposed approach, in Section 3.3. The geometric interpretation of using dual product embeddings and user bias term is discussed in Section 3.4. Finally, we present our fast inference algorithm for cold-start products in Section 3.5.

3.1 Probability model for purchase sequences

	Special cases	Proposed model
Item-item CF models [17]		
Sequential models (<i>item2vec</i> [2])		
<i>triplet2vec</i> [33]		
<i>prod2vec</i> [9]		

Table 1: The Bayesian network representation of joint probability distribution of observation, for our proposed approach (right panel) and classic recommendation algorithms (left panels).

Let \mathcal{U} be the set of N users and \mathcal{I} be the set of M items. The contextual knowledge features for items such as brand, title and description are treated as tokenized words. Product categories and discretized continuous features are also treated as tokens. Without loss of generality, for each item we concatenate all the tokens into a vector denoted by \mathbf{w}_i for item $i \in \mathcal{I}$. We use \mathcal{W}

to denote the whole set of tokens which has n_w instances in total. Similarly, each user u has a knowledge feature vector of tokens \mathbf{x}_u and \mathcal{X} denotes the whole set of n_x user features. A complete observation for user u with the first purchase after time t and previous k (k may vary) consecutive purchases is given by $\{u, i_t, i_{t-1}, \dots, i_{t-k}, \mathbf{x}_u, \mathbf{w}_{i_t}, \mathbf{w}_{i_{t-1}}, \dots, \mathbf{w}_{i_{t-k}}\}$.

Bayesian factorization. The model is optimized for predicting the next purchase according to the user and the most recent purchases. To encode contextual knowledge into item/user representations, we consider the contextual knowledge prediction task:

- $p(I_t | \mathbf{U}, I_{t-1}, \dots, I_{t-k})$ - Predict next-to-buy item given user and recent k purchases.
- $p(\mathbf{W}_I | \mathbf{I})$ - Predict items' contextual knowledge features. By assuming the conditional independence of contextual features, we adopt the factorization such that:

$$p(\mathbf{W}_I | \mathbf{I}) = \prod_{\mathbf{W} \in \mathbf{W}_I} p(\mathbf{W} | \mathbf{I}).$$

- $p(\mathbf{X}_U | \mathbf{U})$ - Predict users' features. We also factorize this term into $\prod_{\mathbf{X} \in \mathbf{X}_U} p(\mathbf{X} | \mathbf{U})$.

To put together all the components, we first point out that the joint probability function of the full observation has an equivalent Bayesian network representation [32] (right panel of Table 1). The log probability function now has the factorization (decomposition) shown in (2). Since we do not model the marginals of $p(\mathbf{I})$ and $p(\mathbf{U})$ with embeddings, we treat them as constant terms denoted by C in (2).

$$\begin{aligned} & \log p(\mathbf{U}, I_t, I_{t-1}, \dots, I_{t-k}, \mathbf{X}_U, \mathbf{W}_I, \mathbf{W}_{I_{t-1}}, \dots, \mathbf{W}_{I_{t-k}}) \\ &= \log \{ p(I_t | \mathbf{U}, I_{t-1}, \dots, I_{t-k}) p(\mathbf{X}_U | \mathbf{U}) \prod_{j=0}^k \log p(\mathbf{W}_{I_{t-j}} | I_{t-j}) \} + C \\ &= \log p(I_t | \mathbf{U}, I_{t-1}, \dots, I_{t-k}) \\ &+ \underbrace{\sum_{\mathbf{X} \in \mathbf{X}_U} \log p(\mathbf{X} | \mathbf{U}) + \sum_{j=0}^k \sum_{\mathbf{W} \in \mathbf{W}_{I_{t-j}}} \log p(\mathbf{W} | I_{t-j})}_{\text{Contextual knowledge}} + C. \end{aligned} \quad (2)$$

Log-odds and loss function. As we have pointed out before, several related models can be viewed as special cases of our approach under the Bayesian network structure. Visual presentations of these models are shown in Table 1. Classic item-item collaborative filtering for recommendation works on conditional probability of item pairs $p(I_t | I_{t-1})$. Sequential item models such as *item2vec* work on $p(I_t | I_{t-1}, \dots, I_{t-k})$. The recently proposed triplet model focus on (item, item, user) triplet distribution: $p(I_t | I_{t-1}, U)$, and *prod2vec* models user purchase sequence as $p(I_t | I_{t-1}, \dots, I_{t-k}, U)$.

In analogy to *SG* and *CBOV* model, each term in (2) can be treated as as multi-class classification problem with *softmax* classifier. However, careful scrutiny reveals that such perspective is not suitable under e-commerce setting. Specifically, using softmax classification function for $p(I_t | \mathbf{U}, I_{t-1}, \dots, I_{t-k})$ implies that given the user and recent purchases we only predict one next-to-buy complementary item. However, it is very likely that there are several complementary items for the purchase sequence. The same

argument holds when modelling item words and user features with multi-class classification. Instead, we treat each term in (2) as the log of probability for the Bernoulli distribution. Now the semantic for $p(\mathbf{I}_t|\mathbf{U}, \mathbf{I}_{t-1}, \dots, \mathbf{I}_{t-k})$ becomes that given the user and recent purchases, whether I_t is purchased next. Similarly, $p(\mathbf{W}|\mathbf{I})$ now implies whether or not the item has the target context words.

Let $s(i, u)$ models personalized preference of user u on item i and $s(i_t, t_{-1:t-k})$ measures the high-order item complementariness given k previously purchased items. A larger value of the two score functions implies stronger user preferences and complementary patterns, respectively. Following the decomposition (1), the log-odds of next-to-buy complementary item is given by (3),

$$\log \left(\frac{P(\mathbf{I}_t = i_t | \mathbf{U}, \mathbf{U}_{t-1} \dots \mathbf{I}_{t-k})}{P(\mathbf{I}_t \neq i_t | \mathbf{U}, \mathbf{U}_{t-1} \dots \mathbf{I}_{t-k})} \right) = s(i_t, u) + s(i_t, t_{-1:t-k}). \quad (3)$$

Similarly, we define $s(w, i)$ as the relatedness of item feature w and item i , $s(x, u)$ as the relatedness of user feature x and user u . Therefore, by treating the observed instance as the positive label and all other instances as negative labels, the loss function for each complete observation under binary logistic loss can be formulated according to the log-likelihood function decomposition in (2):

$$\begin{aligned} \ell = & \log(1 + e^{-s(i_t, u) - s(i_t, i_{t-1:t-k})}) + \sum_{\tilde{i} \in \mathcal{I} \setminus i_t} \log(1 + e^{s(\tilde{i}, u) + s(\tilde{i}, i_{t-1:t-k})}) \\ & + \sum_{j=0}^k \sum_{w \in \mathcal{W}_{i_{t-j}}} \left\{ \log(1 + e^{-s(w, i_{t-j})}) + \sum_{\tilde{w} \in \mathcal{W} \setminus \mathcal{W}_{i_{t-j}}} \log(1 + e^{s(\tilde{w}, i_{t-j})}) \right\} \\ & + \sum_{x \in \mathcal{X}_u} \left\{ \log(1 + e^{-s(x, u)}) + \sum_{\tilde{x} \in \mathcal{X} \setminus x} \log(1 + e^{s(\tilde{x}, u)}) \right\}. \end{aligned} \quad (4)$$

We notice that binary classification loss requires summing over all possible negative instances, which is computationally impractical. We use negative sampling as an approximation for all binary logistic loss terms, with frequency-based negative sampling schema proposed in [20]. For instance, the first line in (4) is now approximated by (5) where $\text{Neg}(\mathcal{I})$ denotes a set of negative item samples.

$$\log(1 + e^{-s(i_t, u) - s(i_t, i_{t-1:t-k})}) + \sum_{\tilde{i} \in \text{Neg}(\mathcal{I})} \log(1 + e^{s(\tilde{i}, u) + s(\tilde{i}, i_{t-1:t-k})}) \quad (5)$$

3.2 Parameterization with context-aware dual embeddings

In last section we give the definitions of the log-odds function: $s(i, u)$, $s(i_t, (i_{t-1}, \dots, i_{t-k}))$, $s(w, i)$ and $s(x, u)$. However, exact estimation of all these functions is impractical as the their amount grows quadratically with the number of items and users. Motivated by previous work in distributed representation learning, we embed the item i , user u , item feature w and user feature x into low dimensional vector space and apply the dot product to approximate the log-odds functions.

Item embeddings in dual space. Let $Z^I \in \mathbb{R}^{M \times P}$ and $\tilde{Z}^I \in \mathbb{R}^{M \times P}$ be the dual item embeddings, such that $s(i_1, i_2) = \langle \tilde{z}_{i_1}^I, z_{i_2}^I \rangle$ and $s(i_2, i_1) = \langle \tilde{z}_{i_2}^I, z_{i_1}^I \rangle$. We refer to Z^I as item-in embedding and

\tilde{Z}^I as item-out embedding. By employing the dual embeddings, the model is able to capture the **asymmetric** and **non-transitive** property. Notice that inner products with only one set of embeddings is inherently symmetric and transitive according to their definition in the Euclidean space.

Knowledge-aware item embedding and user embedding. Let $Z^W \in \mathbb{R}^{n_w \times P}$ be the set of word embeddings, such that $s(w, i) = \langle z_w^W, z_i^I \rangle$. Similarly, let $Z^X \in \mathbb{R}^{n_x \times P}$ be the set of discretize user feature embeddings such that $s(x, u) = \langle z_x^X, z_u^U \rangle$. By imposing the contextual constraints in the loss function (4), contextually similar items and users are enforced to have embeddings close to each other. Due to the **transductive** property of the complementariness relationship, imposing the contextual knowledge constraints helps alleviate the sparsity issue since unpopular items could leverage the information from their similar but more popular counterparts.

User-item interaction term. Let $Z^U \in \mathbb{R}^{N \times P}$ be the set of embeddings for users and $Z^{IU} \in \mathbb{R}^{M \times P}$ be the set of embeddings for item-user context, such that $s(i, u) = \langle \tilde{z}_i^I, z_u^U \rangle$. By capturing the user specific bias in the observed data, learning the user-item interactions helps discover a more general and intrinsic complementary pattern between items.

Higher-order complementariness. We then define the score function when the input is an item sequence, i.e $s(i_t, (i_{t-1}, \dots, i_{t-k}))$. Making prediction based on item sequence has also been spotted in other embedding-based recommender systems. PinSage has experimented on mean pooling and max pooling as well as using LSTM [36]. Youtube uses mean pooling [6], and another work from Alibaba for learning user interests proposes using the attention mechanism [39]. We choose using simple pooling methods over others for interpretability and speed. Our preliminary experiments show that mean pooling constantly outperforms max pooling, so we settle down to the mean pooling shown in (6).

$$s(i_t, i_{t-1:t-k}) \equiv s(i_t, (i_{t-1}, \dots, i_{t-k})) = \frac{\langle \tilde{z}_{i_t}^I, \sum_{j=1}^k z_{i_{t-j}}^I \rangle}{k}. \quad (6)$$

It is now straightforward to see that $s(i, u)$ and $s(i_t, (i_{t-1}, \dots, i_{t-k}))$ meet the demand of a realization of $f_{UI}(\cdot)$ and $f_I(\cdot)$ in (1).

3.3 Ranking criteria for recommendation

In this section, we briefly discuss how to conduct online recommendation with the optimized embeddings, where top-ranked items are provided for customers based on their recent purchases and/or personal preferences.

The first criteria only relies on the score function $s(i_t, i_{t-1:t-k})$ given by (6), which indicates the strength of complementary signal between the previous purchase sequence and target item.

The second criterion accounts for both user preference and complementary signal strength, where target items are ranked according to $s(i_t, u) + s(i_t, i_{t-1:t-k})$. According to the decomposition in (1), the second ranking criterion can be given by:

$$S(\mathbf{I}_{t+1}, (\mathbf{U}, \mathbf{I}_t, \dots, \mathbf{I}_{t-k})) = \langle \tilde{z}_{i_t}^I, z_u^U \rangle + \frac{\langle \tilde{z}_{i_t}^I, \sum_{j=1}^k z_{i_{t-j}}^I \rangle}{k}. \quad (7)$$

The advantage for the second criterion is that user preferences are explicitly considered. However, it might drive the recommendations away from the most relevant complementary item to compensate for user preference. Also, the computation cost is doubled compared with the first criterion. To combine the strength of both methods, we use the first criterion to recall a small pool of candidate items and then use the second criterion to re-rank the candidates.

3.4 Geometric Interpretation

Here we provide our geometric interpretation for the additional item-out embedding \tilde{Z}^I and the user bias term. The use of dual item embedding is not often seen in embedding-based recommender system literature, but it is essential for modelling the intrinsic properties of complementarity. Suppose all embeddings are fixed except \tilde{Z}^I . According to classical separating hyperplane theory, the vector \tilde{z}_j is actually the normal of the separating hyperplane for item j with respect to the embeddings of the positive and negative purchase sequences. In other words, for item j as a next-to-buy complementary item, the hyperplane tries to identify which previous purchase sequences are positive and which are not.

Consider the total loss function $L = \sum_{q=1}^n \ell_q$, where n is the number of observed purchase sequences and the subscript q gives the index of the observation. In the loss function (4), item-out embeddings \tilde{Z}^I only appears in the first two terms. Since L is separable, we collect all terms that involve the item-out embedding of item j as shown in (8). For clarity purpose we let $k = 1$ in (8), i.e only the most recent purchase is included.

$$\mathcal{L}_j = \sum_{(i,u) \in \mathcal{P}} \log(1 + e^{-b_{j,u} - \langle \tilde{z}_j^I, z_i^I \rangle}) + \sum_{(\tilde{i},u) \in \mathcal{N}} \log(1 + e^{b_{j,u} + \langle \tilde{z}_j^I, z_{\tilde{i}}^I \rangle}). \quad (8)$$

In (8), \mathcal{P} represents the whole set of item-user (i, u) pairs in the observed two-item purchase sequences $\{j, i\}$ for user u . \mathcal{N} denotes other (\tilde{i}, u) pairs where item j is used as one of the negative samples in (5). The scalar $b_{j,u} \equiv s(j, u)$ is the preference of user u on item j . It is then obvious that optimizing \tilde{z}_j^I in (8) is equivalent to solving a logistic regression with \tilde{z}_j^I as parameters. The design matrix is constructed from the fixed item-in embeddings $\{z_i^I\}_{(i,u) \in \mathcal{P} \cup \mathcal{N}}$. One difference from ordinary logistic regression is that we have a fixed intercept terms $b_{j,u}$ for each user. Analytically this means that we use the user's preferences on item j , i.e. $s(j, u)$, as the intercept, when using his/her purchase sequences to model the complementary relationship between item j and other items. In logistic regressions, the regression parameter vector \tilde{z}_j^I gives the normal of the optimized separating hyperplane. We provide a sketch of the concept in Figure 1.

When $k \geq 2$, we replace z_i^I in (8) with mean pooling $\frac{1}{k} \sum_{q=1}^k z_{i_q}^I$. Geometrically speaking, we now optimize the separating hyperplane with respect to the centroids of the positive and negative purchase sequences in item-in embedding space. The sketch in also provided in Figure 1. Representing sequences by their centroids helps us capturing the **higher-order** complementarity beyond the pairwise setting.

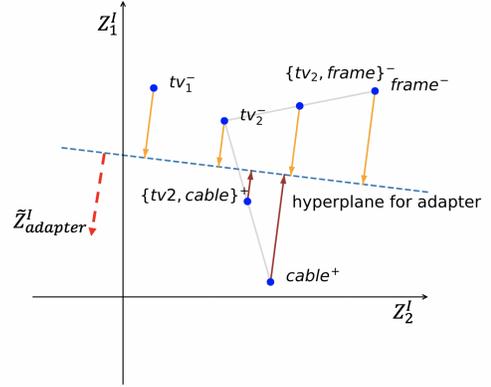


Figure 1: Geometric illustration for the separating hyperplane interpretation with $k = 1, 2$ and $P_I = 2$. For adaptor as complementary item, TV^- , $frame^-$ are negative samples and $cable^+$ is positive sample. Also, $(TV_1, cable)^+$ combo is positive sample and $(TV_2, frame)^-$ is negative sample, where combo is represented by the centroid of the two items. $\tilde{z}_{adaptor}^I$ is optimized such that positive/negative samples have positive/negative distances to the separating hyperplane characterized by its normal vector $\tilde{z}_{adaptor}^I$ and intercept $b_{adaptor}$ (not specified here).

3.5 Cold Start Item Inference

Another advantage of encoding contextual knowledge to item embedding is that we are able to infer item-in embedding for cold-start items using only the trained context embedding Z^W . Therefore, the model gives better coverage and saves us the effort of retraining the whole model, especially for large e-commerce platforms, where customers frequently explore only a tiny fraction of products.

Recall that $\langle \tilde{Z}^I, Z^I \rangle$ measures complementarity and $\langle \tilde{Z}^I, Z^W \rangle$ accounts for contextual relationship. The **transductive** property is preserved on item-in embedding space because items with similar complementary products and contexts are now embedded closer to each other. Although there may not be enough data to unveil the complementary patterns for cold-start products, we can still embed them by leveraging the contextual relationship. For item i_0 with contextual features w_{i_0} , we infer the item-in embedding $z_{i_0}^I$ according to (9) so that it stays closely to contextually similar items in item-in embedding space. As a result, they are more likely to share similar complementary items in recommendation tasks.

$$\hat{z} = \arg \max_{z \in \mathbb{R}^P} \prod_{w \in w_{i_0}} p(w|z), \text{ where } p(w|z) \approx \frac{e^{s(w, i_0)}}{\sum_{\tilde{w} \in \text{Neg}(\mathcal{W})} e^{s(\tilde{w}, i_0)}}. \quad (9)$$

Inferring item-out embedding for cold-start item, however, is not of pressing demand since the platforms rarely recommend cold-start products to customers. We examine the quality of inferred representations in recommendations tasks in Section 4.

4 EXPERIMENT AND RESULT

We thoroughly evaluate the proposed approach on two tasks: product classification and recommendation using only the optimized

product representations. While the classification tasks aim at examining the **meaningfulness** of product representation, recommendation tasks reveal their **usefulness** [33]. We then examine the representations inferred for cold-start products on the same tasks. Ablation studies are conducted to show the importance of incorporating contextual knowledge and user bias term, and the sensitivity analysis are provided for the embedding dimensions. We present several case studies to show that our approach captures the desired properties of complementariness from Section 1.

4.1 Datasets

We work on two real-world datasets: the public *Instacart* dataset and a proprietary e-commerce dataset from *Walmart.com*.

- **Instacart.** *Instacart.com* provides an online service for same-day grocery delivery in the US. The dataset contains around 50 thousand grocery products with shopping records of ~ 200 thousand users, with a total of ~ 3 million orders. All products have contextual information, including name, category, department. User contexts are not available. Time information is not presented, but the order of user activities are available.
- **Walmart.com.** We obtain a proprietary dataset of the *Walmart.com* whose online shopping catalogue covers more than 100 million items. The dataset consists of user add-to-cart and transaction records (with temporal information) over a specific time span, with ~ 8 million products, ~ 20 million customers and a total of ~ 0.1 billion orders. Product contextual information, including name, category, department and brand are also available. Also, user-segmentation (persona) labels are available for a small fraction of users.

Data preprocessing. We follow the same preprocessing procedure for *Instacart* dataset described in [33]. Items with less than ten transactions are removed. Since temporal information is not available, we use each user’s last order as testing data, the second last orders as validation data and the other orders as training data. For the *Walmart.com* dataset we also filter out items with less than ten transactions. We split the data into training, validation and testing datasets using cutoff times in the chronological order. To construct purchase sequence, we use purchases made in past d_1 days as purchase sequence, where d_1 can be treated as a tuning parameter that varies for different use cases. For the *Instacart* dataset we have to make the compromise of using past k purchases since time is not provided. Here k can be treated as sliding window size.

4.2 Implementation

After observing that the loss function of the proposed approach is separable and sparse, we adopt the Hogwild! [24] as optimization algorithm and implement it in C++ for best performance.

Our implementation only takes 15 minutes to train on *Instacart* dataset and 11 hours to train on the proprietary dataset, which contains ~ 2 billion observations from the ~ 5 million users and ~ 2 million items after filtering, in a Linux server with 16 CPU threads and 40 GB memory.

4.3 Results on *Instacart* dataset

We compare the performance of the proposed approach on product classification and with-in basket recommendation against several

state-of-the-art product representation learning methods: *item2vec*, *prod2vec*, *metapath2vec* and *triplet2vec*. We also conduct ablation studies and sensitivity analysis. In ablation studies, we report the performances where item context or user bias term is removed. In the sensitivity analysis, we focus on the dimensions of the item and user embedding. Finally, we evaluate the proposed approach for item cold-start inference. We randomly remove a fraction of items from training dataset as cold-start items and infer their representations after training with the remaining data. The inferred representations are then evaluated using the same two tasks.

Tasks and metric. In product classification tasks we use items’ department and categories as labels to evaluate the product representations under coarse-grained and fine-grained scenarios, respectively. We apply the one-vs-all linear logistics regression with the learned item embeddings as input features. The label fraction is 0.5 and l_2 regularization term is selected from 5-fold cross-validation. For this imbalanced multi-class classification problem, we report the average **micro-F1** and **macro-F1** scores. To prevent information leaks, we exclude department and category information and only use product name and brand for the proposed method.

For the within-basket recommendation task, we rank candidate products according to their complementariness scores to the current basket, by taking the average item-in embeddings as basket representation. The Area Under the ROC Curve (**AUC**) [25] and Normalized Discounted Cumulative Gain (**NDCG**) designed for evaluating recommendation outcome are used as metrics. While **AUC** examines the average performance without taking account of exact ranking status, **NDCG** is top-biased, i.e. higher ranked recommendations are given more credit. The classification and recommendations performances are reported in Table 2.

First we take a look at how k (sequence length) influence the performances in both tasks (Figure 2). For fair comparisons with the baseline models reported in [33], we also set product embedding dimension to 32. In Figure 2 we see that smaller window size leads to better performances on classification tasks, and for the recommendation task the performances are quite stable under all k . The former result suggests that for grocery shopping on *Instacart.com*, more recent purchases play major roles in fulfilling the **meaningfulness** of product representation. However, the latter result indicates that for the recommendation task, our approach is more robust against the length of input sequences, even though longer sequences may introduce more noise. Without further notice, we use $k = 2$ in the following experiments and analysis.

From Table 2 we observe that the proposed approach (**full model**) outperforms all baselines, in both classification and recommendation tasks, which suggests that our model can better capture the **usefulness** and **meaningfulness** of product representations. By taking out the user bias term, we observe a drop in performances (**no user** in Table 2) that indicates the importance of explicitly modelling with user preference. Although removing product context leads to worse performance on product classification (**no context** in Table 2), it gives similar or slightly better results in the recommendation task. We conclude the reasons as follow. Firstly, it is apparent that including context information can improve the **meaningfulness** of product representation since more useful information is encoded. Secondly, as we discussed before, the impact of context knowledge on recommendation can depend on the

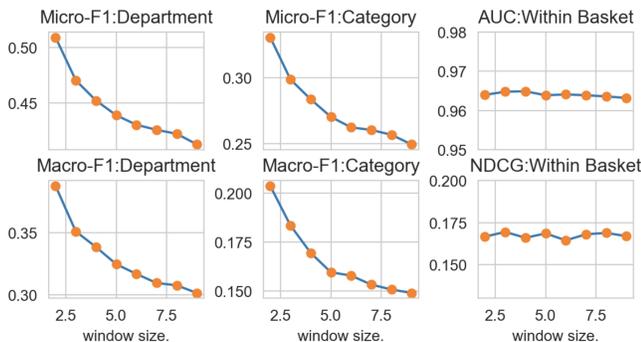


Figure 2: The performances of classification and recommendation tasks under different window size k on *Instacart* dataset without item context. The results are averaged over five runs.

sparsity of user-item interactions. Product-context interactions are introduced to deal with the sparsity issue by taking advantage of the **transductive** property, i.e. similar products may share complementary products. As a consequence, context information may not provide much help in the denser *Instacart* dataset which covers only 50 thousand items. We show later that context information is vital for recommendation performances in the sparse proprietary dataset. Finally, the inference algorithm (9) of cold-start products gives reasonable results, as they achieve comparable performances with baseline methods in both tasks. It also outperforms the heuristic method which assigns representations to cold start products according to their most similar product (according to the Jaccard similarity of the contexts).

Method	Classification		Recom.		In-basket	
	Department	Category	micro	macro	AUC	NDCG
<i>item2vec</i>	0.377	0.283	0.187	0.075	0.941	0.116
<i>prod2vec</i>	0.330	0.218	0.106	0.030	0.941	0.125
<i>m.2vec</i> *	0.331	0.221	0.155	0.036	0.944	0.125
<i>triple2vec</i>	0.382	0.294	0.189	0.082	0.960	0.127
no context	0.509	0.470	0.452	0.438	0.964	<u>0.166</u>
no user	0.661	0.545	0.615	0.529	0.954	0.160
full model	<u>0.666</u>	<u>0.553</u>	<u>0.619</u>	<u>0.535</u>	<u>0.965</u>	0.151
cold start	0.301	0.207	0.304	0.214	0.846	0.114
-infer	-	-	-	-	0.801	0.084

Table 2: Performances on multiple tasks on Instacart data. Methods in bold fonts represents our method under several configurations. The left 4 columns provide micro and macro f-1 scores for the product classification tasks under coarse-grained (department) and fine-grained (category) levels. The right two columns gives the AUC and NDCG for within-basket recommendations. The '-infer' row gives cold start recommendation results based only on context Jaccard similarity without using inference algorithm. **m.2vec* is the shorthand for *metapath2vec*. The results are averaged over five runs.

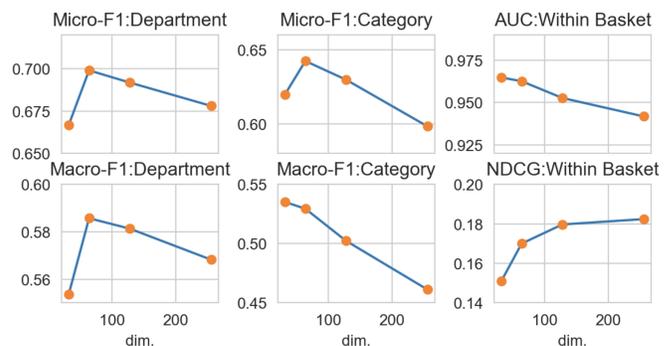


Figure 3: Sensitivity analysis on embedding dimension for product classification and within-basket recommendation tasks on *Instacart* dataset. The results are averaged over five runs.

4.4 Results on the *Walmart.com* dataset

We focus on examining the **usefulness** of our product representation learning method with the *Walmart.com* dataset. As we mentioned before, the availability of timestamps and the lack of **confounding** factors such as loyalty make the *Walmart.com* dataset ideal for evaluating recommendation performances. This also encourages us to compare with baseline models such as collaborative filtering (*CF*) for recommendation. We also include state-of-the-art supervised implicit-feedback recommendation baselines such as factorizing personalized Markov chain (*FPMC*) [26] and Bayesian personalized ranking (*BPR*) [25]. Last but not least, we further compare with the graph-based recommendation method adapted from *node2vec* [34]. On the other hand, without major modifications, the aforementioned product embedding methods do not scale to the size of the proprietary dataset. Still, we manage to implement a similar version of *prod2vec* based on the *word2vec* implementation for a complete comparison.

Tasks and metric. The general next-purchase recommendations (includes within-basket and next-basket) are crucial for real-world recommender systems and are often evaluated by the top-K hitting rate (*Hit@K*) and top-K normalized discounted cumulative gain (*NDCG@K*) according to customers' purchase in next d_2 days. With timestamps available for the proprietary dataset we can carry out these tasks and metrics.

The performance of the proposed approach under different d_1 and d_2 are first provided in Figure 4, where we observe that larger d_2 and smaller d_1 leads to better outcomes. We notice that different d_1 gives very similar results, which is in accordance with the window size in *Instacart* dataset. They both suggest that our approach is robust against the input sequence length. The fact that larger d_2 leads to better results is self-explanatory. Without further notice we use $d_1 = 3$ to construct input sequences and $d_2 = 7$ to extract target products for all models. For fair comparisons with baselines which cannot take advantage of user context information, we also exclude user contextual features in our model in all experiments.

The recommendation performances are provided in Table 3. The proposed approach outperforms all baseline methods, both standard recommendation algorithms and product embedding methods, by significant margins on all metrics. Due to the high sparsity and

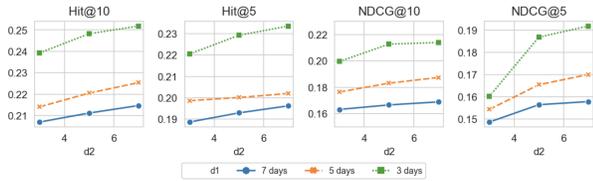


Figure 4: Sensitivity analysis for the d_1 and d_2 , where purchases up to d_1 days before the current time are used as input sequences and use the user’s purchases in next d_2 days are used as labels, with $P = 100$. The results are averaged over five runs.

Model	Hit@10	Hit@5	NDCG@10	NDCG@5
<i>CF</i>	0.0962	0.0716	0.0582	0.0506
<i>FPMC</i>	0.1473	0.1395	0.1386	0.1241
<i>BPR</i>	0.1578	0.1403	0.1392	0.1239
<i>node2vec</i>	0.1103	0.0959	0.0874	0.0790
<i>prod2vec</i>	0.1464	0.1321	0.1235	0.1162
no context	0.1778	0.1529	0.1346	0.1198
no user	0.1917	0.1726	0.1485	0.1301
full model	<u>0.2518</u>	<u>0.2336</u>	<u>0.2139</u>	<u>0.1919</u>
cold start	0.1066	0.0873	0.0759	0.0562
-infer	0.0734	0.0502	0.0531	0.0428

Table 3: Comparison with baseline models under $P = 100$, $d_1 = 3$ and $d_2 = 7$. The results are averaged over five runs.

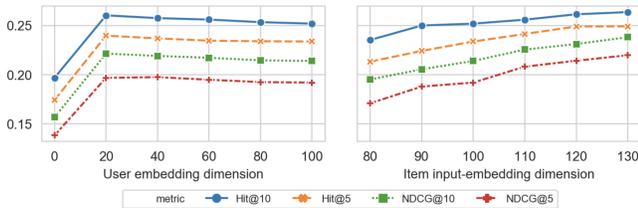


Figure 5: Sensitivity analysis for user and item embeddings dimensions on the proprietary dataset averaged over five runs.

low signal-to-noise ratio of the Walmart.com dataset, the baseline methods may not be able to capture co-purchase signals without accounting for contextual information and user bias.

As expected, after removing product context terms the performances drop significantly, which again confirms the usefulness of contexts for large sparse datasets. Taking out user bias term also deteriorates the overall performance, which matches what we observe from the *Instacart* dataset. Besides, the representations inferred for cold-start products give reasonable performances that are comparable to the outcome of some baselines on standard products.

4.5 Sensitivity analysis

We provide a thorough sensitivity analysis on embedding dimensions on both datasets in Figure 3 and 5. Specifically, for *Instacart* dataset, after increasing product embedding dimension over 80, the proposed model starts to lose classification accuracies. On the other hand, the NDCG for recommendation task keeps improving

with larger dimensions (within the range considered), though the AUC slightly decreases. For the Walmart.com dataset, we focus on both user and product embedding dimensions. Increasing user embedding dimension over 20 starts to cause worse performances. Moreover, similar to the results on *Instacart* dataset, larger product embedding dimensions give better outcomes.

In fact, recent work reveals the connection between embedding size and the bias-variance trade-off [35]. Embedding with large dimensions can overfit training data, and with small dimensions it might underfit. It may help explain our results where a larger product embedding dimension causes overfitting issues in classification tasks on *Instacart* dataset. For the Walmart.com dataset, increasing user embedding dimension over the threshold may also cause overfitted the user preferences in training data.

Anchor item	Top recommended complementary items		
Bath curtain	Bath curtain hooks	Bath curtain rod	Bath mat
Bath mat	Towels	Mops	Detergent
XBox Player	Wireless controller	Video Game	Controller stack
TV	Protection plan	HDMI cable	TV mount frame
TV mount frame	HDMI cable	DVD hold with cable	Cable cover
Round cake pan	Straight spatula	Cupcake pan	Lever tool
Straight spatula	Spatula set	Wood spatula	Steel spatula
Round cake pan +	Cake carrier	Cake decorating set	Icing smoother
Straight spatula			
Sofa	Sofa pillows	Coffee table	Sofa cover
Sofa +	Floor lamp	Area rug	Sofa pillows
Coffee table			

Table 4: Case studies for product recommendation.

4.6 Case study

We provide several case studies in Table 4 to show that the proposed approach captures the desired properties of complementarity, i.e. **asymmetric**, **non-transitive** and **higher-order**. The **transductive** property has been explicitly modelled by including contextual information, so we do not provide concrete examples here.

By virtue of the examples, we observe the **Asymmetric** property via $TV \rightarrow TV\ mount\ frame$ but $TV\ mount\ frame \not\rightarrow TV$. The **Non-transitive** is also reflected by $TV \rightarrow TV\ mount\ frame$, $TV\ mount\ frame \rightarrow cable\ cover$ but $TV \not\rightarrow cable\ cover$. As for the **Higher-order** property, we first note that when *round cake pan* is combined with *straight spatula*, the recommendations are not a simple mixture of individual recommendations. *Cake carrier* and *cake decorating set*, which complement the combo as a whole, are recommended. Similar higher-order pattern is also observed for *sofa* and *coffee table* combination. They provide evidence that our model captures the **higher-order** complementary patterns.

5 CONCLUSION

In this work, we thoroughly investigate the properties of product complementary relationship: **asymmetric**, **non-transitive**, **transductive** and **higher-order**. We propose a novel representation learning model for complementary products with dual embeddings

that are compatible with the above properties. Contextual information is leveraged to deal with the sparsity in customer-product interactions, and user preferences are explicitly modelled to take account of noise in user purchase sequences. A fast inference algorithm is developed for cold-start products. We demonstrate the effectiveness of the product representations obtained from the proposed approach through classification and recommendation tasks on *Instacart* and the proprietary e-commerce dataset.

Future work includes exploring populational contexts such as product popularity as well as dynamic factors including trend and seasonality. Furthermore, the idea of context-aware representation learning for complex relationships can be extended beyond e-commerce setting for learning general event representations from user-event interaction sequences.

REFERENCES

- [1] Marko Balabanović and Yoav Shoham. 1997. Fab: content-based, collaborative recommendation. *Commun. ACM* 40, 3 (1997), 66–72.
- [2] Oren Barkan and Noam Koenigstein. 2016. Item2vec: neural item embedding for collaborative filtering. In *2016 IEEE 26th International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 1–6.
- [3] Mikhail Belkin and Partha Niyogi. 2003. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation* 15, 6 (2003), 1373–1396.
- [4] Jesús Bobadilla, Fernando Ortega, Antonio Hernandez, and Jesús Bernal. 2012. A collaborative filtering approach to mitigate the new user cold start problem. *Knowledge-Based Systems* 26 (2012), 225–238.
- [5] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics* 5 (2017), 135–146.
- [6] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems*. ACM, 191–198.
- [7] Yuxiao Dong, Nitesh V Chawla, and Ananthram Swami. 2017. metapath2vec: Scalable representation learning for heterogeneous networks. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*. ACM, 135–144.
- [8] Zeno Gantner, Lucas Drumond, Christoph Freudenthaler, Steffen Rendle, and Lars Schmidt-Thieme. 2010. Learning attribute-to-feature mappings for cold-start recommendations. In *Data Mining (ICDM), 2010 IEEE 10th International Conference on*. IEEE, 176–185.
- [9] Mihajlo Grbovic, Vladan Radosavljevic, Nemanja Djuric, Narayan Bhamidipati, Jaikrit Savla, Varun Bhagwan, and Doug Sharp. 2015. E-commerce in your inbox: Product recommendations at scale. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1809–1818.
- [10] Aditya Grover and Jure Leskovec. 2016. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 855–864.
- [11] John Hannon, Mike Bennett, and Barry Smyth. 2010. Recommending twitter users to follow using content and collaborative filtering approaches. In *Proceedings of the fourth ACM conference on Recommender systems*. ACM, 199–206.
- [12] Zan Huang, Hsinchun Chen, and Daniel Zeng. 2004. Applying associative retrieval techniques to alleviate the sparsity problem in collaborative filtering. *ACM Transactions on Information Systems (TOIS)* 22, 1 (2004), 116–142.
- [13] Ioannis Konstas, Vassilios Stathopoulos, and Joemon M Jose. 2009. On social networks and collaborative recommendation. In *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*. ACM, 195–202.
- [14] Yehuda Koren and Robert Bell. 2015. *Advances in collaborative filtering*. In *Recommender systems handbook*. Springer, 77–118.
- [15] Xuan Nhat Lam, Thuc Vu, Trong Duc Le, and Anh Duc Duong. 2008. Addressing cold-start problem in recommendation systems. In *Proceedings of the 2nd international conference on Ubiquitous information management and communication*. ACM, 208–211.
- [16] Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *International Conference on Machine Learning*. 1188–1196.
- [17] Greg Linden, Brent Smith, and Jeremy York. 2003. Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet computing* 1 (2003), 76–80.
- [18] Julian McAuley, Rahul Pandey, and Jure Leskovec. 2015. Inferring networks of substitutable and complementary products. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 785–794.
- [19] Prem Melville, Raymond J Mooney, and Ramadass Nagarajan. 2002. Content-booster collaborative filtering for improved recommendations. *Aaai/iaai* 23 (2002), 187–192.
- [20] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*. 3111–3119.
- [21] Mingdong Ou, Peng Cui, Jian Pei, Ziwei Zhang, and Wenwu Zhu. 2016. Asymmetric transitivity preserving graph embedding. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 1105–1114.
- [22] Manos Papagelis, Dimitris Plexousakis, and Themistoklis Kutsuras. 2005. Alleviating the sparsity problem of collaborative filtering using trust inferences. In *Trust management*. Springer, 224–239.
- [23] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 1532–1543.
- [24] Benjamin Recht, Christopher Re, Stephen Wright, and Feng Niu. 2011. Hogwild: A Lock-Free Approach to Parallelizing Stochastic Gradient Descent. In *Advances in Neural Information Processing Systems* 24, J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger (Eds.). Curran Associates, Inc., 693–701. <http://papers.nips.cc/paper/4390-hogwild-a-lock-free-approach-to-parallelizing-stochastic-gradient-descent.pdf>
- [25] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian personalized ranking from implicit feedback. In *Proceedings of the twenty-fifth conference on uncertainty in artificial intelligence*. AUAI Press, 452–461.
- [26] 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*. ACM, 811–820.
- [27] Badrul Sarwar, George Karypis, Joseph Konstan, and John Riedl. 2001. Item-based collaborative filtering recommendation algorithms. In *Proceedings of the 10th international conference on World Wide Web*. ACM, 285–295.
- [28] Martin Saveski and Amin Mantrach. 2014. Item cold-start recommendations: learning local collective embeddings. In *Proceedings of the 8th ACM Conference on Recommender systems*. ACM, 89–96.
- [29] Andrew I Schein, Alexandrin Popescul, Lyle H Ungar, and David M Pennock. 2002. Methods and metrics for cold-start recommendations. In *Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 253–260.
- [30] Ian Soboroff and Charles Nicholas. 1999. Combining content and collaboration in text filtering. In *Proceedings of the IJCAI*, Vol. 99. sn, 86–91.
- [31] Flavian Vasile, Elena Smirnova, and Alexis Conneau. 2016. Meta-prod2vec: Product embeddings using side-information for recommendation. In *Proceedings of the 10th ACM Conference on Recommender Systems*. ACM, 225–232.
- [32] Martin J Wainwright, Michael I Jordan, et al. 2008. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning* 1, 1–2 (2008), 1–305.
- [33] Mengting Wan, Di Wang, Jie Liu, Paul Bennett, and Julian McAuley. 2018. Representing and Recommending Shopping Baskets with Complementarity, Compatibility and Loyalty. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. ACM, 1133–1142.
- [34] Jizhe Wang, Pipei Huang, Huan Zhao, Zhibo Zhang, Binqiang Zhao, and Dik Lun Lee. 2018. Billion-scale Commodity Embedding for E-commerce Recommendation in Alibaba. *arXiv preprint arXiv:1803.02349* (2018).
- [35] Zi Yin and Yuanyuan Shen. 2018. On the dimensionality of word embedding. In *Advances in Neural Information Processing Systems*. 895–906.
- [36] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. 2018. Graph Convolutional Neural Networks for Web-Scale Recommender Systems. *arXiv preprint arXiv:1806.01973* (2018).
- [37] Yin Zhang, Haokai Lu, Wei Niu, and James Caverlee. 2018. Quality-aware Neural Complementary Item Recommendation. In *Proceedings of the 12th ACM Conference on Recommender Systems (RecSys '18)*. ACM, New York, NY, USA, 77–85. <https://doi.org/10.1145/3240323.3240368>
- [38] Jiaqian Zheng, Xiaoyuan Wu, Junyu Niu, and Alvaro Bolivar. 2009. Substitutes or complements: another step forward in recommendations. In *Proceedings of the 10th ACM conference on Electronic commerce*. ACM, 139–146.
- [39] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, 1059–1068.
- [40] Ke Zhou, Shuang-Hong Yang, and Hongyuan Zha. 2011. Functional matrix factorizations for cold-start recommendation. In *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*. ACM, 315–324.

A RESEARCH METHODS

Since we optimize the embeddings with stochastic gradient descent, we first present the gradients for the embeddings with respect to the loss function on single complete observation

$\{u, i_t, i_{t-1}, \dots, i_{t-k}, \mathbf{x}_u, \mathbf{w}_{i_t}, \mathbf{w}_{i_{t-1}}, \dots, \mathbf{w}_{i_{t-k}}\}$ presented in (4) under negative sampling.

Let $Q_{i,u,i_{t-1:t-k}} = s(i, u) + s(i, i_{t-1:t-k})$, $Q_{w,i} = s(w, i)$ and $Q_{x,u} = s(x, u)$. The gradient for \mathbf{z}_u^I can be computed by

$$\begin{aligned} \nabla_{\mathbf{z}_u^I} \ell = & -\mathbf{z}_{i_t}^I \frac{e^{-Q_{i,u,i_{t-1:t-k}}}}{1 + e^{-Q_{i,u,i_{t-1:t-k}}}} + \sum_{\tilde{i} \in \text{Neg}(I)} \mathbf{z}_{\tilde{i}}^I \frac{e^{Q_{i,u,i_{t-1:t-k}}}}{1 + e^{Q_{i,u,i_{t-1:t-k}}}} \\ & - \sum_{x \in \mathcal{X}_u} \mathbf{z}_x^X \frac{e^{-Q_{x,u}}}{1 + e^{-Q_{x,u}}} + \sum_{\tilde{x} \in \mathcal{X} \setminus \mathcal{X}_u} \mathbf{z}_{\tilde{x}}^X \frac{e^{Q_{\tilde{x},u}}}{1 + e^{Q_{\tilde{x},u}}}. \end{aligned} \quad (10)$$

The gradient for $\tilde{\mathbf{z}}_{i_t}^I$ and $\mathbf{z}_{i_t}^I$ are given by

$$\nabla_{\tilde{\mathbf{z}}_{i_t}^I} \ell = -\left(\mathbf{z}_u^I - \frac{1}{k} \sum_{j=1}^k \mathbf{z}_{i_{t-j}}^I\right) \frac{e^{-Q_{i,u,i_{t-1:t-k}}}}{1 + e^{-Q_{i,u,i_{t-1:t-k}}}}, \quad (11)$$

and

$$\nabla_{\mathbf{z}_{i_t}^I} \ell = -\sum_{w \in \mathcal{W}_{i_t}} \mathbf{z}_w^W \frac{e^{-Q_{w,i_t}}}{1 + e^{-Q_{w,i_t}}} + \sum_{\tilde{w} \in \text{Neg}(\mathcal{W})} \mathbf{z}_{\tilde{w}}^W \frac{e^{-Q_{\tilde{w},i_t}}}{1 + e^{-Q_{\tilde{w},i_t}}}. \quad (12)$$

For $\mathbf{z}_{i_{t-j}}^I, j = 1, \dots, k$, we compute their gradients as:

$$\begin{aligned} \nabla_{\mathbf{z}_{i_{t-j}}^I} \ell = & -\mathbf{z}_{i_t}^I \frac{e^{-Q_{i,u,i_{t-1:t-k}}}}{k(1 + e^{-Q_{i,u,i_{t-1:t-k}}})} + \sum_{\tilde{i} \in \text{Neg}(I)} \mathbf{z}_{\tilde{i}}^I \frac{e^{Q_{i,u,i_{t-1:t-k}}}}{k(1 + e^{Q_{i,u,i_{t-1:t-k}}})} \\ & - \sum_{w \in \mathcal{W}_{i_{t-j}}} \mathbf{z}_w^W \frac{e^{-Q_{w,i_{t-j}}}}{1 + e^{-Q_{w,i_{t-j}}}} + \sum_{\tilde{w} \in \text{Neg}(\mathcal{W})} \mathbf{z}_{\tilde{w}}^W \frac{e^{-Q_{\tilde{w},i_{t-j}}}}{1 + e^{-Q_{\tilde{w},i_{t-j}}}}. \end{aligned} \quad (13)$$

As for the embeddings for item contextual feature (word) token w and users feature x that show up in the observation, their gradients can be computed by:

$$\nabla_{\mathbf{z}_w^W} \ell = \sum_{\substack{i:w \in i, \\ i \in \{i_t, \dots, i_{t-k}\}}} -\mathbf{z}_i^I \frac{e^{-Q_{w,i}}}{1 + e^{-Q_{w,i}}} \quad (14)$$

and

$$\nabla_{\mathbf{z}_x^X} \ell = -\mathbf{z}_u^I \frac{e^{-Q_{x,u}}}{1 + e^{-Q_{x,u}}}. \quad (15)$$

For the \tilde{i} , \tilde{w} and \tilde{x} that are being sampled in (10) to (15), their gradients can be calculated in similar fashion. We do not present them here to avoid unnecessary repetition.

During training, we linearly decay the learning rate after each batch update from the starting value α_0 to 0. This practice is in accordance with other widely-used shallow embedding models such as *word2vec*, *doc2vec*, *fastText*, etc. In our synchronous stochastic gradient descent implementation, each observation constitutes its own batch. Therefore with e epochs and n observations, the adaptive learning rate at step l is given by

$$\alpha_l = \frac{l}{n \times e} \alpha_0. \quad (16)$$

Algorithm 1 Framework for learning embeddings

Input: Training dataset \mathcal{D} with n observations, training epochs $e = 30$;

Output: The embeddings $Z^I, \tilde{Z}^I, Z^U, Z^W$ and Z^X ;

- 1: Shuffle \mathcal{D} ;
 - 2: **for** epoch = 1, ..., e **do**
 - 3: **for** $q = 1, \dots, n$ **do**
 - 4: Construct negative sample sets $\text{Neg}(I)$, $\text{Neg}(W)$ and $\text{Neg}(X)$ by (17);
 - 5: Update item/user/word/user feature embedding of the q^{th} observation according to (10)-(15) with learning rate $\alpha_{\text{epoch} \times n + q}$ in (16);
 - 6: Update embeddings of the instances in the above negative sample sets with same learning rate.
 - 7: **end for**
 - 8: **end for**
-

When implementing frequency-based negative sampling, the items are sampled according to the probability computed by

$$p(i) \propto 1 - \sqrt{\frac{10^{-5}}{f(i)}}, \quad (17)$$

where $f(i)$ is the frequency of item i in training dataset [20]. Negative sampling for words and user features are conducted in the same way. We set the number of negative samples to five under all cases in training. The pseudo-code of our algorithm is given in Algorithm (1).