

Neural Representations in Hybrid Recommender Systems: Prediction versus Regularization

Ramin Raziperchikolaei, Tianyu Li, and Young-joo Chung

Rakuten Institute of Technology, Rakuten, Inc.

{*ramin.raziperchikola, tianyu.li, youngjoo.chung*}@rakuten.com

Abstract

Autoencoder-based hybrid recommender systems have become popular recently because of their ability to learn user and item representations by reconstructing various information sources, including users' feedback on items (e.g., ratings) and side information of users and items (e.g., users' occupation and items' title). However, existing systems still use representations learned by matrix factorization (MF) to predict the rating, while using representations learned by neural networks as the regularizer. In this paper, we define the neural representation for prediction (NRP) framework and apply it to the autoencoder-based recommendation systems. We theoretically analyze how our objective function is related to the previous MF and autoencoder-based methods and explain what it means to use neural representations as the regularizer. We also apply the NRP framework to a direct neural network structure which predicts the ratings without reconstructing the user and item information. We conduct extensive experiments on two MovieLens datasets and two real-world e-commerce datasets. The results confirm that neural representations are better for prediction than regularization and show that the NRP framework, combined with the direct neural network structure, outperforms the state-of-the-art methods in the prediction task, with less training time and memory.

1 Introduction

The goal of recommender systems is to help users identify the items that best fit their personal tastes from a large set of items [20]. To achieve this goal, recommender systems use different kinds of user and item information. One important source of information is the feedback of users on items, which could be implicit (e.g., click on a link, purchase, etc.) [4, 19, 20] or explicit (e.g., a rating between 1 and 5) [13, 20]. On e-commerce platforms, predicting users' explicit feedback (e.g., ratings on reviews) is more desirable because it provides better insight about users' preferences. Therefore, we focus on predicting explicit feedback (i.e., ratings).

The focus of this paper is on hybrid recommender systems, which use both feedback of the users on items and *side information* to make prediction [1]. Side information of the users and items include the content of items (e.g., category, title, description, etc.) and profile of users (e.g., age, location, gender, etc.), respectively. Using the feedback and side information jointly helps the hybrid methods to overcome the limitation of other approaches that use either feedback or side information, and achieve state-of-the-art results [15, 16].

More specifically, assume we have a sparse rating matrix $\mathbf{R} \in \mathbb{R}^{m \times n}$, where m and n are the number of users and items, respectively, $R_{jk} > 0$ is the rating of the user j on the item k , and $R_{jk} = 0$ means the rating is unknown. Assume the side information of all the users and items are represented by \mathbf{X} and \mathbf{Y} , respectively. *The goal of hybrid methods is to predict the unknown ratings using the known ratings and the user and item side information.*

Deep neural networks have become popular in designing hybrid recommendation methods, mainly because of their ability in learning good representations. The most widely-used neural network structure in recommender systems has been (denoising) autoencoders [3, 14, 15, 21, 23, 25, 27]. These methods define $\mathbf{g}^u()$, $\mathbf{f}^u()$, $\mathbf{g}^i()$, $\mathbf{f}^i()$ as the user’s encoder, user’s decoder, item’s encoder, and item’s decoder, respectively. The outputs of the encoders $\mathbf{g}^u()$ and $\mathbf{g}^i()$ are the learned neural representations of the users and items. They also consider $\mathbf{U} \in \mathbb{R}^{m \times d}$ and $\mathbf{V} \in \mathbb{R}^{n \times d}$ as the d -dimensional representations of the users and items, respectively. Their objective function can then be written as:

$$\min_{\mathbf{U}, \mathbf{V}, \theta} L(\mathbf{f}^u(\mathbf{g}^u(\mathbf{R}, \mathbf{X}))) + L(\mathbf{f}^i(\mathbf{g}^i(\mathbf{R}, \mathbf{Y}))) + \lambda_1 \sum_{j,k} \mathbb{1}(R_{jk} > 0) \|R_{jk} - \mathbf{U}_{j,:} \mathbf{V}_{k,:}^T\|^2 + \lambda_2 \|\mathbf{U} - \mathbf{g}^u(\mathbf{R}, \mathbf{X})\|^2 + \lambda_3 \|\mathbf{V} - \mathbf{g}^i(\mathbf{R}, \mathbf{Y})\|^2 + \text{reg. terms}, \quad (1)$$

where $\theta = [\theta_{fu}, \theta_{gu}, \theta_{fi}, \theta_{gi}]$ contains all the parameters of the two autoencoders and $\mathbf{U}_{j,:}$ denotes the i th row of the matrix \mathbf{U} . The indicator function $\mathbb{1}(arg)$ returns 1 when arg is true, and 0 otherwise. The rating of the user j on item k is approximated by the dot product of the $\mathbf{U}_{j,:}$ and $\mathbf{V}_{k,:}$.

We divide the objective function of Eq. (1) into three parts:

1. The reconstruction losses in the first two terms try to reconstruct the ratings and the side information of the users and items.
2. The MF in the third term decomposes the rating matrix into user and item representations, which will be used for the prediction later.
3. The fourth and fifth terms try to keep the representations learned by MF in some distance from the neural representations. As we argue in Section. 3, these terms play the role of regularizer, which keep the representations from converging to the solution of MF. The hyper-parameters λ_2 and λ_3 determine how close the two representations should be from each other.

Note that the reconstruction loss in Eq. (1) is slightly different from one work to the others. In [14], each autoencoder reconstructs the repeated versions of the side information. In [3], each autoencoder reconstructs both the ratings and the side information, where the side information has been added to each layer of the network. Li et al. [15] proposed to reconstruct multiple sources of side information of the users and items, in addition to the ratings.

Autoencoder-based methods optimize the objective of Eq. (1) by alternating over the following two steps: 1) fix the network parameters and optimize over \mathbf{U} and \mathbf{V} and 2) fix \mathbf{U} and \mathbf{V} and train the parameters of the two autoencoders.

This approach has three main issues. First, the motivation behind using neural representation for the regularization purpose is unclear. Also, it is difficult to decide how far/close the neural and MF representations should be from each other, i.e., it is difficult to set the hyper-parameters λ_2 and λ_3 . Second, optimization is difficult and time-consuming because 1) the autoencoders

have many parameters, 2) the matrices \mathbf{U} and \mathbf{V} are huge, and 3) the matrices and parameters need to be optimized several times in alternation. Third, the dot product to predict ratings from representations \mathbf{U} and \mathbf{V} might not be sufficient to combine the two representations.

To solve the above issues, we introduce the **Neural Representation for Prediction (NRP)** framework that learns one set of user and item representations from the neural networks and uses them for the prediction directly, instead of using them as the regularizer. Here are the contributions of our paper:

- We propose the NRP framework and apply it to the autoencoder structure, analyze its objective function and optimal solution, and compare it with the previous approaches based on autoencoder and MF.
- We introduce a direct neural network structure integrated with our framework to obtain more expressive power and training efficiency.
- We conduct experiments on two MovieLens datasets and two real-world e-commerce datasets and demonstrate that 1) neural representations perform better in prediction than regularization and 2) the direct neural network structure combined with the NRP framework outperforms previous methods, while having faster training and less memory usage.

2 Related work

Factorization models. Matrix factorization (MF) decomposes the rating matrix into two low-rank user and item matrices (representations), such that the dot product of the representations approximates the rating matrix [12, 13, 24]. Since MF only uses the ratings, its performance degrades significantly when the rating matrix is highly sparse or when we have new users and items in the system (cold-start problem) [20].

Autoencoders for hybrid recommender systems. These methods use the autoencoder structure to extract features from side information and combine such features with the feedback pattern to predict ratings. [3, 14, 15, 23, 25, 27]. Recent works [3, 14, 15] utilize all the side information and ratings of users/items by training two autoencoders. As explained in the introduction, these methods use MF’s representations for the prediction. There are also cases where the MF and neural representations are mixed for prediction. In CDL [25] and AutoSVD [27], the item’s representation comes from the autoencoder and the user’s representation comes from the MF. These methods use the side information of the items as the only source of information, and use dot product to combine the representations. In our method, both users’ and items’ representations are generated from the neural network, where the inputs are the user/item rating vectors and side information.

Deep learning to model implicit feedback in collaborative filtering. DMF [26], NeuMF [10], and DeepCF [4] are examples of the recent works which combine the deep networks and collaborative filtering (CF) to predict the implicit feedback. These methods do not use the autoencoder structure, as the goal of autoencoders is to extract representations from the side information. Our method differs from these methods in several ways. First, our approach is a general framework that can be applied to different network structures, such as autoencoders. Second, our method uses both side information and ratings to predict the ratings, while the input to the CF-based

methods is ratings and/or ids. Third, our approach learns a single representation per user and item, while NeuMF and DeepCF learn two representations. Lastly, our method uses an MLP to map the representations to prediction, while DMF uses dot product.

Deep learning for content-based recommendation. The main idea is to extract representations from the side information such as users’ profiles and items’ descriptions. DSSM [11] maximizes/minimizes the similarity between the representations of the query and the clicked/not-clicked documents. MV-DNN [5] is an extension of the DSSM, where the document has multiple views. The problem definition here is different from ours, as these methods do not consider the previous ratings of the users on items.

Deep Learning for high-order interactions. In web/app recommender systems, the input usually comes in the form of high-dimensional and sparse categorical features. Factorization Machines [18] proposed to model high-order interactions by learning an embedding vector per feature. This idea has been combined with the embedding layer of the neural networks in several recent works, such as Wide&Deep[2], NFM [9], and DeepFM[6]. Similar to the content-based methods, the problem definition is different from us, as these methods do not consider the past ratings of the users on items.

3 Our proposed method

Our main idea is to remove the MF terms and use the neural representations of the users and items for the prediction task. For concision, we call our method **NRP**, which stands for **N**eural **R**epresentation for **P**rediction. We first apply our framework to the autoencoder structure and show how it is related to the previous autoencoder-based methods. Then, we integrate this framework with a direct neural network structure.

3.1 NRP with autoencoders

Similar to the previous works, our model contains two autoencoders, one for the users and one for the items. The difference is that the encoders’ outputs are the only user/item representations in our model. Here is our objective function:

$$\min_{\theta} L(\mathbf{f}^u(\mathbf{g}^u(\mathbf{R}, \mathbf{X}))) + L(\mathbf{f}^i(\mathbf{g}^i(\mathbf{R}, \mathbf{Y}))) + \lambda_1 \sum_{j,k} \mathbb{1}(R_{jk} > 0) ||R_{jk} - \mathbf{g}^u(\mathbf{R}_{j,:}, \mathbf{X}_{j,:})^T \mathbf{g}^i(\mathbf{R}_{:,k}, \mathbf{Y}_{k,:})||^2. \tag{2}$$

To have an apple-to-apple comparison between our approach and the previous ones, we use the same loss functions and the decoder and encoder structures as aSDAE [3] and DHA [15].

Our objective in Eq. (2) gives three advantages over the previous works: 1) the hyper-parameters λ_2 and λ_3 , from Eq. (1), are removed, 2) the number of parameters decreased as we removed \mathbf{U} and \mathbf{V} , which helps in faster training and saving memory, and 3) the network can be trained end-to-end, as there is no need to optimize over \mathbf{U} and \mathbf{V} . In Section 4. we will see that these advantages lead to better prediction performance.

We now analyze the objective function of Eq. (2), compare it with the one in Eq. (1), and explain why the neural representations act as a regularizer in previous works. First, we rewrite our objective in (2) as follows:

$$\begin{aligned} \min_{\boldsymbol{\theta}, \mathbf{U}, \mathbf{V}} Q(\boldsymbol{\theta}, \mathbf{U}, \mathbf{V}) &= L(\mathbf{f}^u(\mathbf{g}^u(\mathbf{R}, \mathbf{X}))) + L(\mathbf{f}^i(\mathbf{g}^i(\mathbf{R}, \mathbf{Y}))) + \lambda_1 \sum_{j,k} \mathbb{1}(R_{jk} > 0) \|R_{jk} - \mathbf{U}_{j,:}^T \mathbf{V}_{k,:}\|^2 \\ \text{s.t. } \mathbf{U} &= \mathbf{g}^u(\mathbf{R}, \mathbf{X}) \quad \text{and} \quad \mathbf{V} = \mathbf{g}^i(\mathbf{R}, \mathbf{Y}). \end{aligned} \quad (3)$$

The objective functions in Equations (2) and (3) are equivalent, so we focus on comparing (3) with (1).

We consider two special cases of the objective function in Eq. (1). First, consider the case where $\lambda_2 = \lambda_3 = 0$, which makes the last two terms 0. The first two terms can also be removed since they do not contain the user/item representations \mathbf{U} and \mathbf{V} . So only the MF term remains.

The second case is when $\lambda_2 = \lambda_3 \rightarrow \infty$. The following theorem shows that in this case the two objective functions in (1) and (3) will be equivalent (i.e. they have the same optimal solution).

theorem 1. *The objective function of the Eq. (1), with $\lambda_2 = \lambda_3 \rightarrow \infty$, has the same optimal solution as the objective function of the Eq. (3).*

Proof. Let us define a new vector $\mathbf{p} = [\boldsymbol{\theta}, \mathbf{U}, \mathbf{V}]$, containing all the parameters of the problem. Note that $Q(\mathbf{p})$, defined in Eq. (3), contains the first three terms of Eq. (1) and contains all the terms of Eq. (3). We assume that $\bar{\mathbf{p}} = [\bar{\boldsymbol{\theta}}, \bar{\mathbf{U}}, \bar{\mathbf{V}}]$ is the optimal solution of (3) and $\mathbf{p}^* = [\boldsymbol{\theta}^*, \mathbf{U}^*, \mathbf{V}^*]$ is the optimal solution of the Eq. (1) when $\lambda_2 = \lambda_3 \rightarrow \infty$. We replace \mathbf{p}^* and $\bar{\mathbf{p}}$ in Eq. (1) to get the following inequality:

$$\begin{aligned} \lim_{\lambda_2 = \lambda_3 \rightarrow \infty} Q(\mathbf{p}^*) + \lambda_2 \|\mathbf{U}^* - \mathbf{g}^u(\mathbf{R}, \mathbf{X}; \boldsymbol{\theta}_{g^u}^*)\|^2 + \lambda_3 \|\mathbf{V}^* - \mathbf{g}^i(\mathbf{R}, \mathbf{Y}; \boldsymbol{\theta}_{g^i}^*)\|^2 \leq \\ Q(\bar{\mathbf{p}}) + \lambda_2 \|\bar{\mathbf{U}} - \mathbf{g}^u(\mathbf{R}, \mathbf{X}; \bar{\boldsymbol{\theta}}_{g^u})\|^2 + \lambda_3 \|\bar{\mathbf{V}} - \mathbf{g}^i(\mathbf{R}, \mathbf{Y}; \bar{\boldsymbol{\theta}}_{g^i})\|^2 = Q(\bar{\mathbf{p}}). \end{aligned} \quad (4)$$

The right side of the above inequality is $Q(\bar{\mathbf{p}})$ since $\bar{\mathbf{p}}$ is the optimal solution of Eq. (3) and satisfies the constraints. By rearranging the last equation in (4), we obtain

$$\|\mathbf{U}^* - \mathbf{g}^u(\mathbf{R}, \mathbf{X}; \boldsymbol{\theta}_{g^u}^*)\|^2 + \|\mathbf{V}^* - \mathbf{g}^i(\mathbf{R}, \mathbf{Y}; \boldsymbol{\theta}_{g^i}^*)\|^2 \leq \lim_{\lambda_2 \rightarrow \infty} \frac{Q(\bar{\mathbf{p}}) - Q(\mathbf{p}^*)}{\lambda_2} = 0. \quad (5)$$

To satisfy the inequality, i.e., making the sum of the norms 0, both norms have to be 0. This means that \mathbf{p}^* is a feasible vector, i.e. $\mathbf{U}^* = \mathbf{g}^u(\mathbf{R}, \mathbf{X}; \boldsymbol{\theta}_{g^u}^*)$ and $\mathbf{V}^* = \mathbf{g}^i(\mathbf{R}, \mathbf{Y}; \boldsymbol{\theta}_{g^i}^*)$. By considering feasibility of \mathbf{p}^* in Eq. (4) (i.e., replacing \mathbf{U}^* by $\mathbf{g}^u(\mathbf{R}, \mathbf{X}; \boldsymbol{\theta}_{g^u}^*)$ and \mathbf{V}^* by $\mathbf{g}^i(\mathbf{R}, \mathbf{Y}; \boldsymbol{\theta}_{g^i}^*)$ in Eq. (4)), we get $Q(\mathbf{p}^*) \leq Q(\bar{\mathbf{p}})$. Note that we assumed that $\bar{\mathbf{p}}$ is the optimal solution of Eq. (3), so for two feasible points $\bar{\mathbf{p}}$ and \mathbf{p}^* we have $Q(\mathbf{p}^*) \geq Q(\bar{\mathbf{p}})$. So the conclusion is that $Q(\mathbf{p}^*) = Q(\bar{\mathbf{p}})$. In other words, for $\lambda_2 = \lambda_3 \rightarrow \infty$, the objective functions of the Eq. (1) and Eq. (3) become equivalent. \square

Fig. 1 shows a simple visualization of the objective functions and their optimal solutions. In this figure, the green contours correspond to the MF (by setting $\lambda_2 = \lambda_3 = 0$ in Eq. (1)) and the magenta contours correspond to $Q(\cdot)$, the main term of our objective function in Eq. (3). The feasible set, which satisfies our constraints in Eq. (3), has been shown by a blue rectangle. This feasible set contains the low-dimensional representations that can be created by the user and item encoders. The optimal solution of our objective function in Eq. (3) lies where the contour line of $Q(\cdot)$ with the smallest value intersects the feasible region.

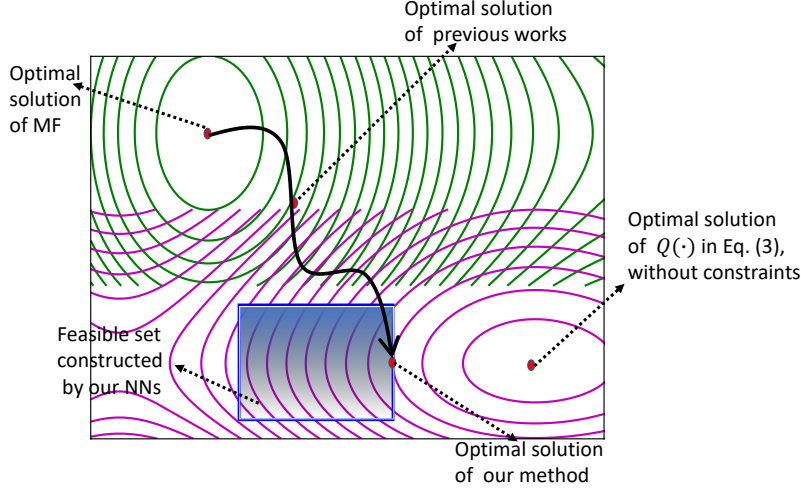


Figure 1: Visualization of the optimal solutions of different methods. The figure shows the contours over the users and items representations (\mathbf{U} and \mathbf{V}). *Green contours*: contours of the MF, which is achieved by setting $\lambda_2 = \lambda_3 = 0$ in Eq. (1). *Magenta contours*: contours of $Q(\cdot)$, the main term of our objective function in Eq. (3). There is a path between the optimal solution of MF and ours. Previous approaches find a solution somewhere in the middle of the path. Our solution is always inside the feasible set created by the encoders.

By setting $\lambda_2 = \lambda_3 = 0$ and increasing it to $\lambda_2 = \lambda_3 \rightarrow \infty$, a path of solutions will be created, between the solution of the MF and our NRP autoencoder. The previous autoencoder methods use a fixed $\lambda_2 > 0$ and $\lambda_3 > 0$, so their optimal solution lies somewhere on the path. The smaller (larger) these hyper-parameters, the closer (farther away) the solution of the autoencoder-based methods will be to the MF’s solution. *We believe the neural representations act as the regularizer in previous works since they are only used to keep \mathbf{U} and \mathbf{V} away from the MF’s optimal solution.*

A question arises here: can we optimize the objective of Eq. (1) with $\lambda_2 = \lambda_3 \rightarrow \infty$ and expect to get the same result as our objective function in Eq. (2)? In practice, this will not happen for several reasons. First, as we set $\lambda_2 = \lambda_3 \rightarrow \infty$, the Hessian of the objective function in Eq. (1) becomes ill-conditioned, and the contours will have a ”banana” shape rather than ”elliptical” shape. The reason is that some eigenvalues of the Hessian will be in the order of ∞ , while the other ones will be small. More information can be found in Chapter 17.1 of [17]. As a result, first-order methods iterate zigzag and slowly approach toward the optimal solution.

Second, note that previous autoencoder-based methods use alternating optimization over the parameters. Consider the step of optimizing over \mathbf{U} and \mathbf{V} for a fixed θ :

$$\min_{\mathbf{U}, \mathbf{V}} \lim_{\lambda_2 = \lambda_3 \rightarrow \infty} \lambda_1 \sum_{j,k} \mathbb{1}(R_{jk} > 0) \|R_{jk} - \mathbf{U}_{j,:}^T \mathbf{V}_{k,:}\|^2 + \lambda_2 \|\mathbf{U} - \mathbf{g}^u(\mathbf{R}, \mathbf{X})\|^2 + \lambda_3 \|\mathbf{V} - \mathbf{g}^i(\mathbf{R}, \mathbf{Y})\|^2.$$

Note that even a small mismatch between the neural’s and MF’s representations, the error in the second and the third terms, makes the objective error ∞ . So the solution to the above optimization problem will be $\mathbf{U} = \mathbf{g}^u(\mathbf{R}, \mathbf{X})$ and $\mathbf{V} = \mathbf{g}^i(\mathbf{R}, \mathbf{Y})$. In other words, the first term, which is responsible for rating prediction, will always be ignored and the performance degrades significantly.

Finally, note that we have proved the equivalency of Eq. (1) and (2) in terms of their optimal solutions. In practice, the objective functions are highly non-convex and we can expect the methods to end up close to a local solution. As we optimize the neural network using SGD, our method starts from some initial point, traverses a path, and stops as soon as overfitting happens. *From*

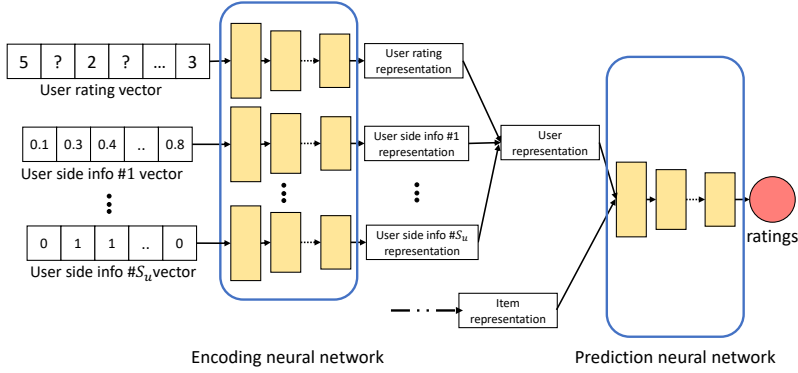


Figure 2: The diagram of our direct neural network structure. Encoding neural network maps the multiple sources of the input to a low-dimensional user representation. Item representation is created by the same structure. Prediction neural network takes the joint (user,item) representation and predicts the rating.

start to end, the solution of our method remains inside the feasible set, while the final solution of the autoencoder-based methods can be anywhere in the space.

Size of the feasible set. The constraints of the Eq. 3 determines the feasible set. Assuming d -dimensional representations, the size of the feasible set is determined by the number of d -dimensional vectors that can be generated by the user and item encoders, which depends on the complexity of the model, i.e., increasing the complexity of the encoders leads to a larger feasible set. In our visualization in Fig. 1, as the feasible set gets larger, the optimal solution of our method gets closer to the optimal solution of $Q(\cdot)$. As discussed before, the optimal solution of the previous autoencoder-based methods always stays somewhere in the middle of the path.

3.2 NRP with a direct structure

We define the direct structure as a neural network structure without the decoders, which predict the ratings without reconstructing the user/item ratings and side information. Our main goal of designing the direct structure and combining it with NRP framework is to use it as a baseline. The comparison between the direct structure and the autoencoder-based methods let us know the effectiveness of the reconstruction based methods in hybrid recommendation systems.

Our direct structure is achieved by making two modifications to our autoencoder structure. First, we remove the decoders from the structure, which leads to saving around 50% of memory and faster optimization. Second, we use a set of fully connected layers to predict the final rating, instead of the dot product. This makes our model more expressive. Note that these fully connected layers add a small amount of memory to our model since they connect the low-dimensional user and item representations to the final rating. On the other hand, the decoders are usually huge, since they reconstruct the input ratings with thousands to millions of dimensions.

Fig. 2 shows the overall architecture of our direct structure. We define three cooperating sub-networks: user encoding network, item encoding network, and prediction network. The user and item encoding networks generate the user and item representations, respectively, which are combined to get the joint (user,item) representation. Finally, the prediction network takes the joint representation and predicts the rating.

Learning the user and item representations. We focus on the user representation, as the item representation can be learned in the same way. We have access to multiple sources of information for each user, including ratings of the user on a subset of items, and S_u sources of side information. The user encoding network maps each source to a low-dimensional space and then combines them to get the final user representation.

One source of information is the ratings of the users on items. For the j th user, $\mathbf{R}_{j,:}$ contains the ratings of this user on all items. We give this as the input to a set of fully connected layers and get the rating representation $\mathbf{g}_{\text{rating}}^u(j) \in \mathbb{R}^{1 \times d_{\text{rating}}}$ for the j th user.

We also learn one representation for each source of side information. Each source will be given as the input to a neural network, which maps it to a low-dimensional representation. The representation for the s_u th source is shown by $\mathbf{g}_{s_u}^u(j) \in \mathbb{R}^{1 \times d_{s_u}}$. The structure of the network for different sources could be different. For example, it could be a convolutional network for the images, a long-short term memory (LSTM) network for the text descriptions and titles, and fully connected layers for the other general type of inputs.

To get the final representation for the j th user, we concatenate the representations of all the sources. We show this by $\mathbf{z}_j = [\mathbf{g}_{\text{rating}}^u(j), \{\mathbf{g}_{s_u}^u(j)\}_{s_u=1}^{S_u}] \in \mathbb{R}^{1 \times d_u}$, where $[\]$ is used for concatenation of the vectors.

We can get the k th item representation \mathbf{z}_k with minor changes to the user representation explained above. For the rating representation, the input should be a column of the rating matrix \mathbf{R} , instead of the row. We concatenate the representations of the different sources of information to get the final item representation \mathbf{z}_k .

Learning the rating. We first concatenate the user and item representations to get the joint representation \mathbf{z}_{jk} . This joint representation is the input to another set of fully connected networks that try to predict the rating of the users on the items. We show this last neural network by $h(\cdot)$. We define the objective function as the mean squared error between the output of our neural network and the true rating:

$$\frac{1}{mn} \min_{\boldsymbol{\theta}} \sum_{j=1}^m \sum_{k=1}^n \mathbb{1}(R_{jk} > 0) \|R_{jk} - h(\mathbf{z}_{jk})\|^2 + \lambda_1 \|\boldsymbol{\theta}\|^2$$

s.t. $\mathbf{z}_{jk} = [\mathbf{z}_j, \mathbf{z}_k]$, $\mathbf{z}_j = [\mathbf{g}_{\text{rating}}^u(j), \{\mathbf{g}_{s_u}^u(j)\}_{s_u=1}^{S_u}]$, $\mathbf{z}_k = [\mathbf{g}_{\text{rating}}^i(k), \{\mathbf{g}_{s_i}^i(k)\}_{s_i=1}^{S_i}]$ (6)

We jointly optimize all parameters of the neural network using stochastic gradient descent. Note that our method works as long as we have at least one source of information for users and one source of information for items.

Direct structure in the literature. We noticed that different variants of the direct structure have been used in other areas of the recommendation systems, such as implicit feedback prediction from user-item interactions [4, 10] and modeling high-order user-item interactions [2, 6]. To the best of our knowledge, the direct structure has been used once in hybrid recommender systems by Shi et al. [22] (denoted by ACCM) to address the cold-start issue. There are three main differences between our direct structure and ACCM: 1) our approach uses interaction vector as the input, while ACCM uses user/item ID as the input, 2) our structure uses MLPs to map user-item interaction to the rating, while ACCM uses the dot product, and 3) our approach uses concatenation to combine user and item representation, while ACCM uses the weighted sum. We will show in our experiments that ACCM achieves comparable results to the autoencoder-based methods, while our direct structure outperforms them.

ID of the users and items as another source of information It is easy to adjust our model and objective function to include other sources, such as IDs. In our experiments, we found that adding IDs as another source does not improve the performance. Since IDs might be useful in some other datasets, we briefly explain how to add this source to our model.

Let’s focus on the user ID first. To learn the low-dimensional representation for the user IDs, we use an embedding layer. The embedding layer works as a look-up table; it takes the user ID and returns one row, which contains the embedding vector of the user. We denote the embedding vector of the j th user by $\mathbf{g}_{\text{ID}}^u(j)$ and add it to the \mathbf{z}_j in Eq. 6. The embedding vector of the k th item, $\mathbf{g}_{\text{ID}}^i(k)$, can be added to the \mathbf{z}_k in the same way. The rest of our objective function does not change.

Let’s see how it works mathematically for the user embedding. We encode each user ID to a one-hot vector of size m , where m is the total number of users. For the j th user, $\mathbf{I}_{j,:}$ shows this encoding, where \mathbf{I} is an identity matrix of size $m \times m$. This one-hot vector is connected to a layer with the weight matrix $\mathbf{E}^{(u)} \in \mathbb{R}^{m \times d_e}$. The output of the embedding layer is achieved by matrix multiplication $\mathbf{g}_{\text{ID}}^u(j) = \mathbf{I}_{j,:} \mathbf{E}^{(u)} \in \mathbb{R}^{1 \times d_e}$, which returns the j th row of $\mathbf{E}^{(u)}$ containing the d_e dimensional embedding of the j th user.

4 Experiments

For each dataset, we randomly select 80% of the ratings as the training set, 10% as the validation set, and 10% as the test set. We repeat the process three times to create three training/validation/test sets and report the mean and standard deviation of each method on these three sets. Unless otherwise stated, we use bag of words (BoW) to represent the item and user side information.

Datasets. We use four datasets in our experiments. Table 1 lists the number of users and items, and the amount of sparsity of each dataset.

1. ml100k [7]. The dataset contains 100 000 ratings (1 to 5) from around 1 000 users on 1 600 movies. The user side information contains age, gender, occupation, and zip code; the size of the feature vector is 879. The item side information contains the movie title and the genre; the size of the feature vector is 2 479.
2. ml1m [7]. It contains 1 million ratings of around 6 000 users on 4 000 movies. The content of the user and item side information are the same as ml100k. The dimensions of the user and item side information are 92 and 4 606, respectively.
3. Amazon review data: Grocery and Gourmet Food [8]. The dataset contains 508 800 Amazon’s Grocery and Gourmet Food product review from 86 400 users on 108 500 items. There is no user side information. The item side information contains the title and the category, and is represented by a 36 258 dimensional vector.
4. Ichiba. It contains 1.5 million Rakuten Ichiba¹’s product review from 324 000 users on 294 000 items. The user side information contains age and gender represented by 121 dimensional vector. The item side information contains the category with 455 dimensions.

¹https://rit.rakuten.co.jp/data_release/

Table 1: Summary of the four datasets.

Dataset	# of users	# of items	sparsity
ml100k	1 000	1 600	94%
ml1m	6 000	4 000	96%
Amazon R.	86 400	108 500	99.994%
Ichiba	324 000	294 000	99.84%

Evaluation metrics. We use the root mean square error (RMSE) and precision to evaluate the prediction performance. Let us assume set T contains all the ratings in the test set, where $R_{jk} \in T$ is the actual rating of the user j on item k . We define \hat{R}_{jk} as the predicted rating, generated by a recommender system. Then, RMSE is defined as follows:

$$\text{RMSE} = \sqrt{\frac{1}{|T|} \sum_{R_{jk} \in T} (R_{jk} - \hat{R}_{jk})^2}. \quad (7)$$

To report precision, we need to define the set of retrieved items and the set of relevant items per user. We define set S_j as the set of items rated by the user j . To define the relevant items (groundtruth) for user j , we sort the items in S_j based on their rating and pick the top $p\%$. At the test time, for user j , we predict the ratings of the items in S_j and consider the top $p\%$ of the items (with the highest predicted ratings) as the retrieved set. We define the precision for user j as:

$$\text{precision} = \frac{|\text{relevant items}| \cap |\text{retrieved items}|}{|\text{retrieved items}|} \quad (8)$$

We report the average precision of all the users in test set.

Experimental setting. We implement our method using Keras with TensorFlow 1.12.0 backend. We ran all the experiments on a 12GB GPU. For each method, we tried a set of activation functions (relu, selu, and tanh), a range of learning rates and regularization parameters from 10^{-1} to 10^{-5} , a set of optimizers (Adam, SGD, and RMSprop), and picked the one that works best. For a fair comparison, all autoencoder methods have the same structure (# of layers, neurons, etc.). Supplementary material at the end of this script contains the details of the experimental setting.

Neural representations are better in prediction than regularization. In Table 2, we compare our proposed method with matrix factorization (MF), the autoencoder-based methods, DHA [15] and aSDAE [3], and the attention-based direct structure, ACCM [22], on MovieLens datasets. We also replaced the dot product of ACCM with MLP, denoted by ACCM_{MLP} in Table 2. Our framework applied to the same encoder-decoder structure as DHA and aSDAE are called NRP_{DHA} and NRP_{aSDAE}, respectively. Our framework applied to the direct structure is called NRP_{direct}.

MF has the worst performance (the largest RMSE), which means that MF formulation with the L2 norm of weights as the regularization is not enough for the rating prediction. By setting the hyper-parameters carefully, both DHA and aSDAE outperform SVD, which suggests that the neural network’s representations are better regularizers than the L2 norm of the weights. Our methods NRP_{DHA} and NRP_{aSDAE} outperform the original DHA and aSDAE, respectively, in

Table 2: Our proposed NRP framework, trained with the autoencoder and direct structures, versus MF, autoencoder-based methods, and ACCM (the direct structure of Shi et al. [22]) on ml100k and ml1m datasets. We report mean and standard deviation of the methods. The first/second number in the fourth column is the number of parameters involved in learning the neural network’s/MF’s representations. Time refers to the training time per epoch. Our framework combined with the direct structure achieves better prediction results, faster training, and less memory usage compared to the autoencoder-based methods.

.....ml100k.....				
method	RMSE	precision	# params.	time
MF	0.940 ± 0.003	$68.4\% \pm 0.5$	(0, 0.26M)	15s
DHA	0.939 ± 0.002	$68.2\% \pm 0.6$	(8.6M, 0.26M)	85s
NRP_{DHA}	0.926 ± 0.004	$68.4\% \pm 0.6$	(8.6M, 0)	68s
aSDAE	0.946 ± 0.005	$68.0\% \pm 1.1$	(13M, 0.26M)	98s
NRP_{aSDAE}	0.910 ± 0.008	$69.0\% \pm 0.2$	(13M, 0)	70s
ACCM	0.928 ± 0.004	$68.2\% \pm 0.1$	(3.1M, 0)	37s
ACCM _{MLP}	0.925 ± 0.005	$67.7\% \pm 0.3$	(3.3M, 0)	37s
NRP_{direct}	0.899 ± 0.006	$70.2\% \pm 0.4$	(4.7M, 0)	42s

.....ml1m.....				
method	RMSE	precision	# params.	time
MF	0.892 ± 0.004	$68.2\% \pm 0.3$	(0, 1M)	45s
DHA	0.865 ± 0.001	$69.3\% \pm 0.2$	(44M, 1M)	1 097s
NRP_{DHA}	0.855 ± 0.002	$69.6\% \pm 0.2$	(44M, 0)	1 027s
aSDAE	0.879 ± 0.005	$69.0\% \pm 0.1$	(66M, 1M)	1 155s
NRP_{aSDAE}	0.877 ± 0.008	$68.5\% \pm 0.4$	(66M, 0)	1 055s
ACCM	0.856 ± 0.002	$69.5\% \pm 0.3$	(11.5M, 0)	450s
ACCM _{MLP}	0.865 ± 0.002	$68.9\% \pm 0.2$	(11.8M, 0)	470s
NRP_{direct}	0.851 ± 0.001	$70.0\% \pm 0.1$	(22M, 0)	640s

Table 3: We report the precision by creating the relevant and retrieved sets using top 10% and 25% of the items. We put "OM" in the tables whenever we get an out-of-memory error.

method	ml1m		Amazon review	
	top 10%	top 25%	top 10%	top 25%
MF	55.6% \pm 0.16	68.05% \pm 0.45	64.9% \pm 0.04	71.5% \pm 0.67
Autorec	57.6% \pm 0.26	69.5% \pm 0.43	62.6% \pm 0.98	69.8% \pm 0.62
NeuMF	56.8% \pm 0.12	68.9% \pm 0.48	66.8% \pm 0.30	72.6% \pm 0.09
DSSM	54.7% \pm 0.35	67.2% \pm 0.30	NA	NA
DHA	57.4% \pm 0.78	69.3% \pm 0.23	OM	OM
NRP _{DHA}	57.3% \pm 0.17	69.5% \pm 0.32	66.6% \pm 0.60	72.9% \pm 0.63
aSDAE	56.4% \pm 0.39	68.7% \pm 0.42	OM	OM
NRP _{aSDAE}	57.1% \pm 0.3	69.0% \pm 0.41	64.5% \pm 0.43	71.2% \pm 0.68
HIRE	57.4% \pm 0.08	69.4% \pm 0.55	OM	OM
NRP _{direct}	58.1% \pm 0.16	69.9% \pm 0.42	67.3% \pm 0.38	73.1% \pm 0.24

terms of RMSE and precision. This improvement comes from removing the MF terms, relying on neural representations, and staying inside the feasible set of neural network’s output. ACCM, the direct structure of Shi et al. [22], outperforms MF, DHA, and aSDAE, but shows similar performance to the NRP_{DHA}. So by using the direct structure of ACCM we cannot conclude that the decoders are unnecessary to get the best results. Finally, NRP_{direct} has the best performance and outperforms all autoencoder-based methods. This shows that removing the decoders, making the neural network free of reconstructing the inputs, and replacing the dot product with MLPs lead to learning a better model.

Note that NRP_{direct} outperforms ACCM because of 1) using interaction vectors as the input instead of IDs and 2) using MLPs, instead of the dot product, to map the joint representation to the rating. To show this, we have replaced the dot product of ACCM with MLPs, denoted by ACCM_{MLP} in Table 2. We can see that the result of ACCM_{MLP} is slightly better than ACCM in ml100k and slightly worse in ml1m. NRP_{direct} outperforms ACCM_{MLP} in both datasets.

The last two columns of Table 2 compare the number of learnable parameters and training time per epoch of each method. DHA and aSDAE have the largest memory usage and training time, as they optimize the two representations alternatively. NRP_{DHA} and NRP_{aSDAE} achieve better training time and memory than DHA and aSDAE, respectively. Among the neural network-based methods, the direct structures, ACCM and NRP_{direct}, have the fastest training and lowest memory usage because of their simple structure. Note that NRP_{direct} has a larger number of parameters than ACCM. This is because NRP_{direct} uses MLPs to map the interaction vectors to the low-dimensional representations, while ACCM uses embedding layers.

Comparison with other hybrid and collaborative filtering methods. We compare our method with several baselines and recent works in Tables 3 and 4. MF [13] is a baseline collaborative filtering method that uses dot product of the user and item representations to predict the

Table 4: RMSE of our NRP framework compared with the hybrid and collaborative filtering methods. Our approach outperforms the rest of the methods.

method	ml100k	ml1m	Amazon	Ichiba
MF	0.940 ± 0.003	0.892 ± 0.0004	1.153 ± 0.003	1.00 ± 0.104
Autorec	0.921 ± 0.002	0.889 ± 0.0003	2.19 ± 0.01	2.47 ± 0.059
NeuMF	0.948 ± 0.005	0.886 ± 0.001	1.140 ± 0.004	0.900 ± 0.004
DSSM	0.934 ± 0.002	0.941 ± 0.0004	NA	0.913 ± 0.003
DHA	0.939 ± 0.002	0.865 ± 0.001	OM	OM
NRP _{DHA}	0.926 ± 0.004	0.855 ± 0.002	1.135 ± 0.002	OM
aSDAE	0.946 ± 0.005	0.879 ± 0.005	OM	OM
NRP _{aSDAE}	0.910 ± 0.008	0.877 ± 0.008	1.24 ± 0.004	OM
HIRE	0.930 ± 0.006	0.861 ± 0.004	OM	OM
NRP _{direct}	0.897 ± 0.003	0.851 ± 0.001	1.135 ± 0.002	0.889 ± 0.002

Table 5: RMSE of our method NRP_{direct}, which is trained with and without user and item side information, on ml100k and Ichiba datasets. Side information helps in achieving better prediction performance.

	ml100k		Ichiba	
	precision	RMSE	precision	RMSE
no side info	0.901 ± 0.006	$69.9\% \pm 0.43$	0.895 ± 0.003	$78.9\% \pm 0.002$
side info	0.897 ± 0.003	$70.2\% \pm 0.42$	0.889 ± 0.002	$80.1\% \pm 0.004$

rating. Autorec [21] is another collaborative filtering method which uses autoencoders to reconstruct the ratings. NeuMF [10] combines deep and shallow networks and uses the user/item ID as the input to predict the ratings. NeuMF was originally proposed for the prediction of implicit feedback, but it can easily be modified to work with explicit feedback. DSSM [11] is a content-based recommender method, which uses deep neural networks to learn the representations. We modify DSSM to make it applicable to the explicit feedback prediction by connecting the user/item representations to a MLP with a mean squared error loss. DSSM is not applicable in the Amazon dataset, as there is no user side information. HIRE [16] is a hybrid method that considers the hierarchical user and item side information. DHA [15] and aSDAE [3] are autoencoder-based methods, which use neural representations as the regularizer.

In Table 3 we report precision and in Table 4 we report RMSE of methods on four datasets. We can see that our method achieves the best results on different datasets.

Importance of side information. We train our model with and without the user/item side information to verify that the side information can improve the prediction results. Table 5 lists the RMSE and precision results on different datasets. We can see that the side information helps

to achieve better results.

5 Conclusion

The current autoencoder-based hybrid recommender systems learn two types of representations. One comes from the matrix factorization and is used for prediction. The other comes from neural networks and is used for regularization. In this paper, we proposed a new framework that uses the neural networks' representation directly for the prediction task. We showed that by applying our approach to the same autoencoder structure as previous works, we can achieve faster training and better performance. We also proposed a simpler network structure by removing the decoders and replacing dot product with MLP in autoencoders. Our approach combined with the new proposed framework outperformed the previous works. It also had a fast training and small memory usage compared to the autoencoder-based methods.

References

- [1] Robin Burke. Hybrid web recommender systems. In *The Adaptive Web: Methods and Strategies of Web Personalization*, pages 377–408. Springer Berlin Heidelberg, 2007.
- [2] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishikesh Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, Rohan Anil, Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, and Hemal Shah. Wide & deep learning for recommender systems. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems (DLRS)*, 2016.
- [3] Xin Dong, Lei Yu, Zhonghuo Wu, Yuxia Sun, Lingfeng Yuan, and Fangxi Zhang. A hybrid collaborative filtering model with deep structure for recommender systems. In *Proceedings of the 31st AAAI Conference on Artificial Intelligence*, 2017.
- [4] Zhi-Hong Dong, Ling Huang, Chang-Dong Wang, Jian-Huang Lai, and Philip S Yu. DeepCF: a unified framework of representation learning and matching function learning in recommender system. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*, 2019.
- [5] Ali Mamdouh Elkahky, Yang Song, and Xiaodong He. A multi-view deep learning approach for cross domain user modeling in recommendation systems. In *Proceedings of the 24th International Conference on World Wide Web (WWW)*, 2015.
- [6] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. DeepFM: a factorization-machine based neural network for CTR prediction. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI)*, 2017.
- [7] F. Maxwell Harper and Joseph A. Konstan. The MovieLens datasets. *ACM Transactions on Interactive Intelligent Systems*, 5(4):1–19, 2015.
- [8] Ruining He and Julian McAuley. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *Proceedings of the 25th International Conference on World Wide Web (WWW)*, 2016.
- [9] Xiangnan He and Tat-Seng Chua. Neural factorization machines for sparse predictive analytics. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, 2017.

- [10] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. Neural collaborative filtering. In *Proceedings of the 26th International Conference on World Wide Web (WWW)*, 2017.
- [11] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM International on Conference on Information and Knowledge Management (CIKM)*, 2013.
- [12] Yehuda Koren. Factorization meets the neighborhood: a multifaceted collaborative filtering model. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 2008.
- [13] Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.
- [14] Sheng Li, Jaya Kawale, and Yun Fu. Deep collaborative filtering via marginalized denoising auto-encoder. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management (CIKM)*, 2015.
- [15] Tianyu Li, Yukun Ma, Jiu Xu, Bjorn Stenger, Chen Liu, and Yu Hirate. Deep heterogeneous autoencoders for collaborative filtering. In *Proceedings of the 18th IEEE International Conference Data Mining (ICDM)*, 2018.
- [16] Tianqiao Liu, Zhiwei Wang, Jiliang Tang, Songfan Yang, Gale Yan Huang, and Zitao Liu. Recommender systems with heterogeneous side information. In *Proceedings of the 28th International Conference on World Wide Web (WWW)*, 2019.
- [17] Jorge Nocedal and Stephen Wright. *Numerical Optimization*. Springer, 2006. ISBN 978-0-387-40065-5.
- [18] Steffen Rendle. Factorization machines. In *Proceedings of the 10th IEEE International Conference Data Mining (ICDM)*, 2010.
- [19] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. BPR: bayesian personalized ranking from implicit feedback. In *Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence (UAI)*, 2009.
- [20] Francesco Ricci, Lior Rokach, Bracha Shapira, and Paul B. Kantor, editors. *Recommender Systems Handbook*. Springer US, 2011.
- [21] Suvash Sedhain, Aditya Krishna Menon, Scott Sanner, and Lexing Xie. AutoRec: Autoencoders meet collaborative filtering. In *Proceedings of the 24th International Conference on World Wide Web (WWW)*, 2015.
- [22] Shaoyun Shi, Min Zhang, Yiqun Liu, and Shaoping Ma. Attention-based adaptive model to unify warm and cold starts recommendation. In *Proceedings of the 27th ACM International on Conference on Information and Knowledge Management (CIKM)*, 2018.
- [23] Florian Strub, Romaric Gaudel, and Jérémie Mary. Hybrid recommender system based on autoencoders. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems (DLRS)*, 2016.

- [24] Gábor Takács, István Pilászy, Bottyán Németh, and Domonkos Tikk. Matrix factorization and neighbor based algorithms for the netflix prize problem. In *Proceedings of the 2nd ACM Conference on Recommender Systems (RecSys)*, 2008.
- [25] Hao Wang, Naiyan Wang, and Dit-Yan Yeung. Collaborative deep learning for recommender systems. In *Proceedings of the 21st ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 2015.
- [26] Hong-Jian Xue, Xinyu Dai, Jianbing Zhang, Shujian Huang, and Jiajun Chen. Deep matrix factorization models for recommender systems. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI)*, 2017.
- [27] Shuai Zhang, Lina Yao, and Xiwei Xu. AutoSVD++: An efficient hybrid collaborative filtering model via contractive auto-encoders. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, 2017.

Supplementary material: experimental settings

We implement our method using Keras with TensorFlow 1.12.0 backend. We ran all the experiments on a 12GB GPU. For each method, we tried a set of activation functions (relu, selu, and tanh), a range of learning rates and regularization parameters from 10^{-1} to 10^{-5} , a set of optimizers (Adam, SGD, and RMSprop), and picked the one that works best. For a fair comparison, all autoencoder methods have the same structure (# of layers, neurons, etc.). Table 6 lists these details for all the methods.

In the following, we give more information about the structure of the neural network and other hyper-parameters of each method. The notation $[a, b, c]$ denotes a network with three fully connected layers, where a, b , and c are the number of neurons in each layer.

- **NRP_{direct}** . The encoding networks are $[500, 200, 100]$, $[1000, 500, 300, 100]$, $[500, 300, 100]$, and $[500, 300, 100]$ in ml100k, ml1m, Amazon, and Ichiba datasets. The prediction network is $[500, 200, 100, 50, 1]$ in all datasets.
- **DHA [15] and aSDAE [3]** . The encoders and decoders in these two methods are the same as NRP_{direct}, which makes the comparisons fair. These two methods have a large number of hyper-parameters and we tried a large range of values to get the best results. We have set $\lambda_2 = \lambda_3 = 0.5$ and $\lambda_1 = 1$. We use SGD with learning rate of 0.1 for the optimization over the variables \mathbf{U} and \mathbf{V} .
- **NRP_{DHA}** and **NRP_{aSDAE}** . The encoders and decoders in these two methods are the same as NRP_{direct} , DHA, and aSDAE, which make the comparisons fair.
- **HIRE [16]**. We used the code provided by the authors without changing the hyper-parameters.
- **NeuMF [10]**. This method has one deep and one shallow branches. The structure of the deep branch is the same as the encoding networks of the NRP_{direct}. The shallow branch, which uses the IDs as the input, computes the element-wise product of the representations, as explained in the original paper.
- **Autorec [21]** . We implemented the I-Autorec, which reconstructs the ratings of the items. This method overfits to the training data fast, even using small networks. The autoencoder is $[m, 500, m]$ in ml100k, as suggested by the original paper, and $[m, 500, 300, 100, 300, 500, m]$ in ml1m, Amazon, and Ichiba datasets, where m is the number of users.
- **DSSM [11]** Each branch has the same structure as the encoder network in our NRP_{direct}.
- **ACCM [22]**. It uses an embedding layer to map user and item IDs to low-dimensional representations. To map side information to the low-dimensional representations, it uses the same structure as the encoder network of the NRP_{direct}.
- **ACCM_{MLP} [22]**. Same as the ACCM above. To map the joint user and item representation to the rating, it uses the same structure as the prediction network of the NRP_{direct}.

Table 6: Implementation details: optimization method (optimizer), learning rate (lr), regularization (regu.), activation function (activ.).

ml100kml1m			
	optimizer	lr	regu.	activ.	optimizer	lr	regu.	activ.
NRP _{direct}	RMSprop	0.001	0	selu	SGD	0.0005	0	selu
aSDAE	SGD	0.1	10 ⁻⁵	tanh	RMSprop	0.001	10 ⁻⁵	tanh
DHA	SGD	0.1	10 ⁻⁵	tanh	RMSprop	0.001	10 ⁻⁵	tanh
NRP _{aSDAE}	RMSprop	0.001	0.03	tanh	RMSprop	0.001	0	tanh
NRP _{DHA}	RMSprop	0.001	0.03	tanh	RMSprop	0.001	10 ⁻⁴	tanh
Autorec	Adam	0.001	0.005	tanh	Adam	0.001	0.001	tanh
NeuMF	RMSprop	0.001	0.01	selu	RMSprop	0.001	0.01	selu
DSSM	RMSprop	0.001	0.001	selu	SGD	0.0005	0	selu
MF	SGD	0.1	0.0001	NA	SGD	0.1	0.0001	NA
ACCM	RMSprop	0.001	0	tanh	SGD	0.001	0	tanh

 Amazon Ichiba			
	optimizer	lr	regu.	activ.	optimizer	lr	regu.	activ.
NRP _{direct}	RMSprop	0.001	0.001	selu	SGD	0.0005	0.001	selu
aSDAE	OM	OM	OM	OM	OM	OM	OM	OM
DHA	OM	OM	OM	OM	OM	OM	OM	OM
NRP _{aSDAE}	RMSprop	0.001	0.001	selu	OM	OM	OM	OM
NRP _{DHA}	RMSprop	0.001	0.001	tanh	OM	OM	OM	OM
Autorec	RMSprop	0.001	0	selu	Adam	0.001	0.001	tanh
NeuMF	RMSprop	0.00001	0.001	selu	SGD	0.0001	0.01	selu
DSSM	NA	NA	NA	NA	SGD	0.0005	0.001	selu
MF	SGD	0.1	0.0001	NA	SGD	0.1	0.0001	NA