

Neural Collaborative Filtering vs. Matrix Factorization Revisited

Steffen Rendle* Walid Krichene*
Li Zhang* John Anderson*

Abstract

Embedding based models have been the state of the art in collaborative filtering for over a decade. Traditionally, the dot product or higher order equivalents have been used to combine two or more embeddings, e.g., most notably in matrix factorization. In recent years, it was suggested to **replace the dot product with a learned similarity e.g. using a multilayer perceptron** (MLP). This approach is often referred to as *neural collaborative filtering* (NCF). In this work, we revisit the experiments of the NCF paper that popularized learned similarities using MLPs. First, we show that with a proper hyperparameter selection, a simple dot product substantially outperforms the proposed learned similarities. Second, while a MLP can in theory approximate any function, we show that **it is non-trivial to learn a dot product with an MLP**. Finally, we discuss practical issues that arise when applying MLP based similarities and show that MLPs are too costly to use for item recommendation in production environments while dot products allow to apply very efficient retrieval algorithms. We conclude that MLPs should be used with care as embedding combiner and that dot products might be a better default choice.

1 Introduction

Embedding based models have been the state of the art in collaborative filtering for over a decade. A core operation of most of these embedding based models is to combine two or more embeddings. For example, combining a user embedding with an item embedding to obtain a single score that indicates the preference of the user for the item. This can be viewed as a **similarity function** in the embedding space. Traditionally, a dot product or higher order products have been used for the similarity. Recently, it has become popular to learn the similarity function with a neural network. Most commonly, a multilayer perceptron (MLP) is used for the network architecture (e.g. [18, 37, 38, 20, 32, 27]). This approach is often referred to as *neural collaborative filtering* (NCF) [16]. The rationale is

*Google Research, Mountain View, USA. {srendle,walidk,liqzhang,janders}@google.com

that MLPs are general function approximators so that they should be strictly better than a fixed similarity function such as the dot product. This has made NCF the model of choice for comparison in many recommender studies (e.g. [18, 37, 30, 38, 20, 27]).

In this work, we study MLP versus dot product similarities in more detail. We start with revisiting the experiments of the NCF paper [16] that popularized the use of MLPs in recommender systems. We show that a carefully configured dot product baseline largely outperforms the MLP. At first glance, it looks surprising that the MLP, which is a universal function approximator, does not perform at least as well as the dot product. We investigate this issue in a second experiment and show empirically that learning a dot product with high accuracy for a decently large embedding dimension requires a large model capacity as well as many training data. Besides prediction quality, we also discuss the inference cost of dot product versus MLPs, where dot products have a large advantage due to the existence of efficient maximum inner product search algorithms. Finally, we discuss that dot product vs MLP is not a question of whether a deep neural network (DNN) is useful. In fact, many of the most competitive DNN models, such as transformers in natural language processing [9] or resnets for image classification [14], use a dot product similarity in their output layer.

To summarize, this paper argues that MLP-based similarities for combining embeddings should be used with care. While MLPs can approximate any continuous function, their inductive bias might not be well suited for a similarity measure. Unless the dataset is large or the embedding dimension is very small, a dot product is likely a better choice.

2 Definitions

In this section, we formalize the problem and review dot product (esp., matrix factorization) and learned similarity functions (esp., MLP and NeuMF). We denote matrices by upper case letters X , vectors by lowercase bold letters \mathbf{x} , scalars by lowercase letters x . A concatenation of two vectors \mathbf{x}, \mathbf{z} is denoted by $[\mathbf{x}, \mathbf{z}]$.

Our paper studies functions $\phi : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ that combine two d -dimensional embedding vectors $\mathbf{p} \in \mathbb{R}^d$ and $\mathbf{q} \in \mathbb{R}^d$ into a single score. For example \mathbf{p} could be the embedding of a user, \mathbf{q} the embedding of an item, and $\phi(\mathbf{p}, \mathbf{q})$ is the affinity of this user to the item.

The embeddings \mathbf{p} and \mathbf{q} can be model parameters such as in matrix factorization, but they can also be functions of other features, for example the user embedding \mathbf{p} could be the output of a deep neural network taking user features as input. From here on, we focus mainly on the similarity function ϕ but in Section 6.1 we will discuss the embeddings in more detail.

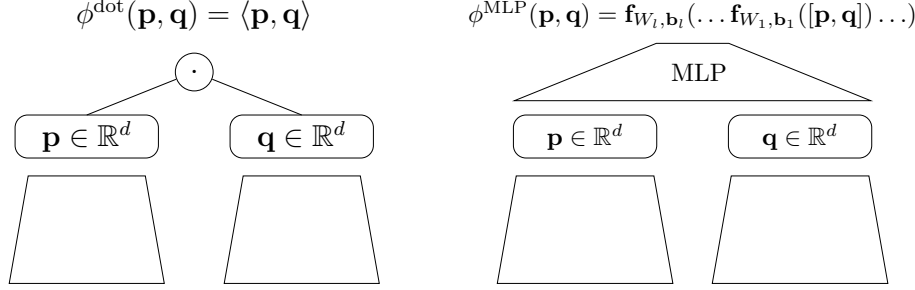


Figure 1: A model with dot product similarity (left) and MLP-based learned similarity (right).

Dot Product The most common combination of two embeddings is the dot product.

$$\phi^{\text{dot}}(\mathbf{p}, \mathbf{q}) := \langle \mathbf{p}, \mathbf{q} \rangle = \mathbf{p}^T \mathbf{q} = \sum_{f=1}^d p_f q_f. \quad (1)$$

If \mathbf{p} and \mathbf{q} are **free model parameters**, then this is equivalent to matrix factorization. **A common trick is to add explicit biases:**

$$\phi^{\text{dot}}(\mathbf{p}, \mathbf{q}) := b + p_1 + q_1 + \langle \mathbf{p}_{[2, \dots, d]}, \mathbf{q}_{[2, \dots, d]} \rangle. \quad (2)$$

This modification does not add expressiveness but has been found to be useful in many studies, **likely because its inductive bias is better suited to the problem** [31, 22].

Learned Similarity Multi layer perceptrons (MLPs) are known to be universal approximators that can approximate any continuous function on a compact set as long as the MLP has enough hidden states [7]. One layer of a multi layer perceptron can be defined as a function $\mathbf{f}: \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$:

$$\mathbf{f}_{W, \mathbf{b}}(\mathbf{x}) = \sigma(W \mathbf{x} + \mathbf{b}), \quad \sigma(\mathbf{z}) = [\sigma(z_1), \dots, \sigma(z_{\text{out}})], \quad (3)$$

which is parameterized by $W \in \mathbb{R}^{\text{in} \times \text{out}}$, $\mathbf{b} \in \mathbb{R}^{\text{out}}$ and an activation function $\sigma: \mathbb{R} \rightarrow \mathbb{R}$. In a multilayer perceptron (MLP), several layers of \mathbf{f} are stacked, e.g., for a three layer MLP, $\mathbf{f}_{W_3, \mathbf{b}_3}(\mathbf{f}_{W_2, \mathbf{b}_2}(\mathbf{f}_{W_1, \mathbf{b}_1}(\mathbf{x})))$.

He et al. [16] propose to replace the dot product with learned similarity functions for collaborative filtering. They suggest to concatenate the two embeddings, \mathbf{p} and \mathbf{q} , and apply an MLP:

$$\phi^{\text{MLP}}(\mathbf{p}, \mathbf{q}) := \mathbf{f}_{W_l, \mathbf{b}_l}(\dots \mathbf{f}_{W_1, \mathbf{b}_1}([\mathbf{p}, \mathbf{q}]) \dots). \quad (4)$$

They further suggest a variation that combines the MLP with a weighted dot product model and name it *neural matrix factorization* (NeuMF):

$$\phi^{\text{NeuMF}}(\mathbf{p}, \mathbf{q}) := \phi^{\text{MLP}}(\mathbf{p}_{[1, \dots, j]}, \mathbf{q}_{[1, \dots, j]}) + \phi^{\text{GMF}}(\mathbf{p}_{[j+1, \dots, d]}, \mathbf{q}_{[j+1, \dots, d]}), \quad (5)$$

where GMF is a ‘generalized’ matrix factorization model:

$$\phi^{\text{GMF}}(\mathbf{p}, \mathbf{q}) := \sigma(\mathbf{w}^T(\mathbf{p} \odot \mathbf{q})) = \sigma(\langle \mathbf{w} \odot \mathbf{p}, \mathbf{q} \rangle) = \sigma\left(\sum_{f=1}^d w_f p_f q_f\right). \quad (6)$$

with learned weights $\mathbf{w} \in \mathbb{R}^d$. For NeuMF, they recommend to use one part of the embedding (here the first j entries) in the MLP and the remaining $d - j$ entries with the GMF.

Fig. 1 illustrates two models with dot product and MLP-based similarity.

3 Revisiting NCF Experiments

In this section, we revisit the experiments of the NCF paper [16] that popularized the use of MLPs as embedding combiners in recommender systems. We show that a simple dot product yields better results.

3.1 Experimental setup

The NCF paper [16] evaluates on an item retrieval task on two datasets: a binarized version of Movielens 1M [13] and a dataset from Pinterest [12]. Both are implicit feedback datasets, i.e. they contain only binary positive tuples between a user and an item. **For each user, the last item is held out and used as the test set, the remaining items of the user are placed into the training set.** For evaluation, each recommender ranks, for each user, a set of 101 items consisting of the withheld test item together with 100 random items. For each user, the position at which the withheld item is ranked by the recommender is recorded, then two metrics are measured: (1) Hit Ratio (i.e. Recall) among the top 10 ranked items – which in this case is 1 if the withheld item is in the top 10 or 0 otherwise. (2) NDCG among the top 10 ranked items – which in this case is $1/\log(r + 1)$ where r is the rank of the withheld item. The average metric over all users is reported. The authors have published the dataset splits and the evaluation code. This allows us to evaluate on exactly the same setting and to compare our results directly with the ones reported in [16].

3.2 Models, loss and training algorithm

We compare three models: MLP-learned similarity models introduced in [16], which use ϕ^{MLP} and ϕ^{NeuMF} respectively, and a simple matrix factorization baseline which uses ϕ^{dot} from Eq. (2). The only difference between these models is the similarity function. In particular, the embeddings \mathbf{p}, \mathbf{q} are free parameters in all models. **We train the matrix factorization baseline by minimizing a logistic loss with L2 regularization, using stochastic gradient descent** (with no batching, no momentum or other variations) with negative sampling, as in the original

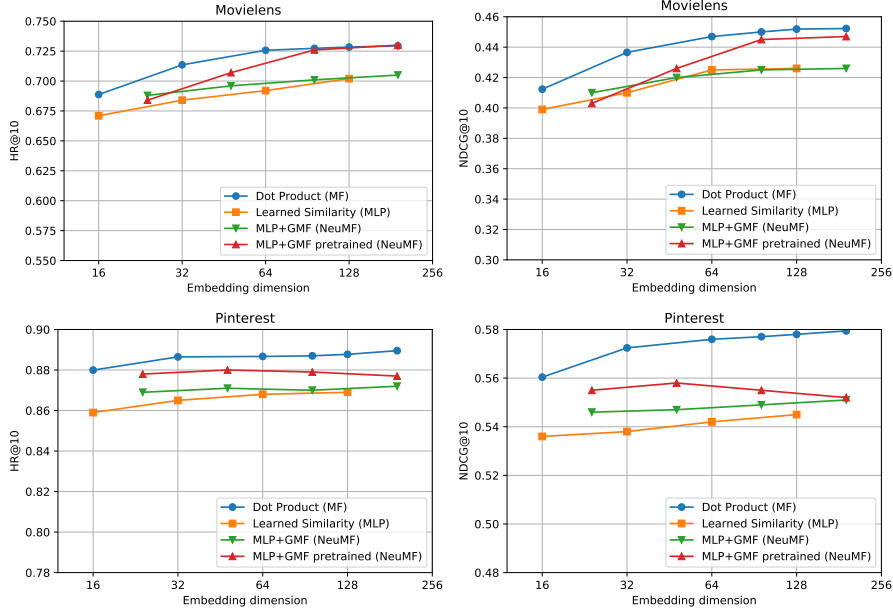


Figure 2: Comparison of learned similarities (MLP, NeuMF) to a dot product: The results for MLP and NeuMF are from [16]. The dot product substantially outperforms the learned similarity measures. Only the pretrained NeuMF is competitive, on one dataset, and for large embedding dimension.

paper [16]¹. More precisely, for each training example (consisting of a user and a positive item), we sample m negative items, uniformly at random. Finally, we vary the embedding dimension $d \in \{16, 32, 64, 96, 128, 192\}$. Additional details about the setup can be found in Appendix A.

3.3 Results

The results are reported in Fig. 2. Contrary to the findings of the NCF paper, the simple matrix factorization model exhibits the best quality over all evaluation metrics, and all embedding dimensions but one.

3.3.1 Matrix Factorization vs MLP

Our main interest is to investigate if the MLP-learned similarity is superior to a simple dot product. As can be seen in Fig. 2, the dot product substantially

¹It is possible that a different loss or a different sampling strategy could lead to an even better performance of our method. However, we wanted to use the same loss and sampling strategy for all competing methods to ensure that this is a meaningful comparison, which will allow us to attribute differences in quality to the choice of similarity functions.

outperforms MLP on all datasets, evaluation metrics and embedding dimensions. With a properly set up matrix factorization model, the experiments do not show any evidence that a MLP is superior. In addition to a lower prediction quality, MLP-learned similarity suffers from other disadvantages compared to dot-product: the model has more model parameters (see Section 4.1), and is more expensive to serve (see Section 5).

3.3.2 Matrix Factorization vs NeuMF

The NCF paper [16] also proposes a combined model where the similarity function is a sum of dot-product and MLP, as in Eq. (5) – this is called NeuMF². The green curve in Fig. 2 shows the performance of this combined model. One can observe only a minor improvement over MLP and overall a much worse quality than MF. The experiments do not support the claim in [16] that a dot product model can be enhanced by feeding some part of its embeddings through an MLP.

A second variant of NeuMF was proposed in [16], that first trains MLP and MF models separately, then fine tunes the combined model. This can be viewed as a form of ensembling. The red curve shows this variant, which performs better than training the combined model directly (in green), but performs worse than the MF baseline overall, except on one datapoint (HR on MovieLens with embedding dimension $d = 192$). Once again, the results do not support the claim that a learned similarity using a MLP is superior to a dot product. The experiment only indicates that ensembling two models can be helpful, a fact that has been observed for a variety of applications, and it is possible that ensembling with different models may yield a similar improvement. The fact remains that using a simple dot product outperforms this ensemble.

3.3.3 On the performance of GMF

Other variants of matrix factorization were considered in [16]. In particular, the GMF model uses a weighted dot product ϕ^{GMF} as described in Eq. (6). Except for the weights in the dot product, this model is very similar to the MF baseline we trained, in particular, both models use the same loss and negative sampling method. Nevertheless, the GMF results reported in [16] are much worse than our MF results. This discrepancy may seem surprising at first glance. We can see two reasons for this difference. First, properly setting up and tuning baseline methods can be difficult in general, as argued in [33], and the reported results may be improved by a more careful setup.

Second, ϕ^{GMF} introduces new model parameters – the vector \mathbf{w} in Eq. (6). While this appears to be an innocuous generalization of the dot product similarity, it can have negative effects. For example, L2 regularization of the embeddings (\mathbf{p} and \mathbf{q}) is meaningless unless \mathbf{w} is regularized as well. More precisely,

²Following [16], the NeuMF uses 2/3rds of the embeddings for the MLP and 1/3rd for the MF. See the discussion about “predictive factors” in Section A.3 for details.

suppose the loss function is of the form

$$L(P, Q, \mathbf{w}, \lambda) = \ell(\{\phi_{\mathbf{w}}^{\text{GMF}}(\mathbf{p}, \mathbf{q}) : \mathbf{p} \in \text{Rows}(P), \mathbf{q} \in \text{Rows}(Q)\}) + \lambda(\|P\|_F^2 + \|Q\|_F^2)$$

where P, Q are embedding matrices, the first term of the loss ℓ depends on the pairwise similarities (i.e. the model output), and the second term is a regularization term, where P, Q are regularized but \mathbf{w} is not. Observe that if we scale the model parameters as $P/a, Q/a, a^2\mathbf{w}$ for some positive scalar a , then the model output is unchanged (given the expression of ϕ^{GMF}), and we have

$$L(P, Q, \mathbf{w}, \lambda) = L\left(\frac{1}{a}P, \frac{1}{a}Q, a^2\mathbf{w}, a^2\lambda\right). \quad (7)$$

It follows that minimizing L with a given λ is equivalent to minimizing L with any other $\tilde{\lambda}$ up to the change of variable $(P/a, Q/a, a^2\mathbf{w})$ with $a = \sqrt{\tilde{\lambda}/\lambda}$, a change of variable which leaves the model output unchanged. The solution is therefore unaffected by regularization. A second consequence is that unless $\lambda = 0$, minimizing the loss L will likely result in embedding matrices P, Q of vanishing norm and a vector of weights \mathbf{w} of diverging norm, leading to numerical instability.

The GMF results in [16] support that the model is indeed not properly regularized because its results do not improve with a higher embedding dimension – unlike in our experiments.

Finally, we observe that GMF does not improve model expressivity compared to a simple dot product, since the weights \mathbf{w} can simply be absorbed into the embedding matrices P and Q . This is another indicator that adding parameters to a simple model is not always a good idea and has to be done carefully.

3.4 Further comparison

As reported in the meta study of [8], the results for NeuMF and MLP in [16] **were cherry-picked in the following sense**: the metrics are reported for the best iteration *selected on the test set*. The NeuMF and MLP numbers we report in Fig. 2 are from the original paper and likely over-estimate the actual test performance of those methods. On the other hand, our MF results in Fig. 2 are not cherry picked, because we select all hyperparameters including the stopping iteration on a validation set – see Appendix A for details. **The fact that our baseline MF outperforms the MLP-learned similarity despite the cherry-picking in the latter strengthens our conclusions.**

In this section, we give an additional comparison using non cherry-picked results produced by [8]. Table 1 includes their results together with our matrix factorization (same as in Fig. 2), with embedding dimension $d = 192$. The results confirm that the simple matrix factorization model substantially outperforms NeuMF on all metrics and datasets. Our results provide further evidence to the conclusion of [8] that simple, well-known baselines outperform NCF. Note that matrix factorization was also one of the baselines in [8] (the iALS method in Table 1), but our experiment shows a much larger margin than was obtained in [8].

Table 1: Comparison from [8] of MLP+GMF (NeuMF) with various baselines and our results. The best results are highlighted in bold, the second best result is underlined.

Method	Movielens		Pinterest		Result from
	HR@10	NDCG@10	HR@10	NDCG@10	
Popularity	0.4535	0.2543	0.2740	0.1409	[8]
SLIM [29, 24]	<u>0.7162</u>	<u>0.4468</u>	0.8679	<u>0.5601</u>	[8]
iALS [19]	0.7111	0.4383	0.8762	0.5590	[8]
MLP+GMF [16]	0.7093	0.4349	<u>0.8777</u>	0.5576	[8]
Matrix Factorization	0.7294	0.4523	0.8895	0.5794	Fig. 2

3.5 Discussion

Following the arguments in [33], it is possible that the studies in [16] and [8] did not properly set up MLP and NeuMF, and that these results could be further improved. It is also possible that the performance of these models is different on other datasets. Nevertheless, at this point, the revised experiments from [16] provide **no evidence supporting the claim that a MLP-learned similarity is superior to a dot product**. This negative result also holds for NeuMF where a GMF is added to the MLP. And it also holds for the pretrained version of NeuMF. Our study treats MLP and NeuMF favorably: (1) we report the results for MLP and NeuMF that were obtained by the original authors, avoiding any bias in improperly running their methods. (2) These cited numbers for MLP and NeuMF are likely too optimistic as they were obtained through cherry picking as identified by [8].

4 Learning a Dot Product with MLP is Hard

An MLP is a universal function approximator: any continuous function on a compact set can be approximated with a large enough MLP [7, 17, 3]. It is tempting to argue that this makes the MLP a more powerful embedding combiner and it should thus perform at least as well or better than a dot product. However, such an argument neglects the difficulty of learning the target function using MLPs: the larger class of functions also implies more parameters needed for representing the function. Hence it would require more data to learn the function and may encounter difficulty in actually learning the desired target function. Indeed, specialized structures, e.g. convolutional, recurrent, and attention structures, are common in neural networks. There is probably no hope to replace them using an MLP though they should all be representable. However, is this also true for the simple “structure” of the dot product? Similar problems turn out to be actively studied subject in machine learning theory [2, 25, 10, 1]. To our knowledge, the best theoretical bound for learning the dot product, a degree two polynomial, requires $O(d^4/\epsilon^2)$ steps for an error bound of ϵ [2]. While the theory gives only a sufficient condition, it does hint that the difficulty scales

polynomially with dimension d and $1/\epsilon$. This motivates us to investigate the question empirically.

4.1 Experimental setup

We set up a synthetic learning task³ where given two embeddings $\mathbf{p}, \mathbf{q} \in \mathbb{R}^d$ and a label $y(\mathbf{p}, \mathbf{q})$, we want to learn a function $\hat{y} : \mathbb{R}^{2d} \rightarrow \mathbb{R}$ that approximates y with $\hat{y}(\mathbf{p}, \mathbf{q})$. We draw the embeddings \mathbf{p}, \mathbf{q} from $\mathcal{N}(0, \sigma_{\text{emb}}^2 I)$ and set the true label as $y(\mathbf{p}, \mathbf{q}) = \langle \mathbf{p}, \mathbf{q} \rangle + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma_{\text{label}}^2)$ models the label noise. From this process we create three datasets each consisting of tuples $(\mathbf{p}, \mathbf{q}, y)$. One of the datasets is used for training and the remaining two for testing. For the training and first test dataset, we first sample M different user embeddings and N different item embeddings, i.e., there are two fixed embedding matrices $P \in \mathbb{R}^{M \times d}$ and $Q \in \mathbb{R}^{N \times d}$. Then we uniformly sample (without replacement) 100 M user-item combinations and put 90% into the training set and 10% into the test set. We create a second test set that consists of *fresh* embeddings that did not appear in the training or test set, i.e., we sample the embeddings for every case from $\mathcal{N}(0, \sigma_{\text{emb}}^2 I)$ instead of picking them from P and Q . The motivation for this setup is to investigate if the learned similarity function generalizes to embeddings that were not seen during training.

We train the MLP on the training dataset and evaluate it on both test datasets. For the architecture of the MLP, we follow the suggestions in the NCF paper: we use an input layer of size $2d$ consisting of the concatenation of the two embeddings, and 3 hidden layers with sizes $[4h, 2h, h]$ where h is a parameter, and use the ReLU as the activation function. The NCF paper suggests to use $h = d/2$, we also experiment with $h = d$ and $h = 2d$. For $h = d$, the number of model parameters are about $18d^2$, so for example for $d=8$: 1,152 or $d=64$: 73,728 or for $d=256$: 1,179,648. For optimization, we also follow the NCF paper and choose the Adam optimizer.

As evaluation metric, we compute the RMSE between the predicted similarity of the MLP and the true similarity y . We also measure the RMSE of a trivial model that predicts always 0 (=average rating in our dataset). In our setup, this RMSE is equal in expectation to $\sqrt{\text{Var}(y)} = \sqrt{\sigma_{\text{label}}^2 + d\sigma_{\text{emb}}^4}$. Secondly, we measure the RMSE of the dot product model, i.e., $\hat{y}(\mathbf{p}, \mathbf{q}) = \langle \mathbf{p}, \mathbf{q} \rangle$. This RMSE is equal in expectation to σ_{label} . We report the approximation error, i.e., the difference between the RMSE of the dot product model and the MLP. Each experiment is repeated 5 times and we report the mean.

We want to choose the experimental parameters σ_{label} and σ_{emb} such that the approximation error gives some indication what values are acceptable. To do this we choose values that are related to well-studied rating prediction tasks. In the Netflix prize, the best models have RMSEs of about 0.85 [21] – for Movielens 10M, the best models have about 0.75 [33]. For these datasets, it is likely that the label noise is close to these values, thus we choose the label

³The code is available at https://github.com/google-research/google-research/tree/master/dot_vs_learned_similarity.

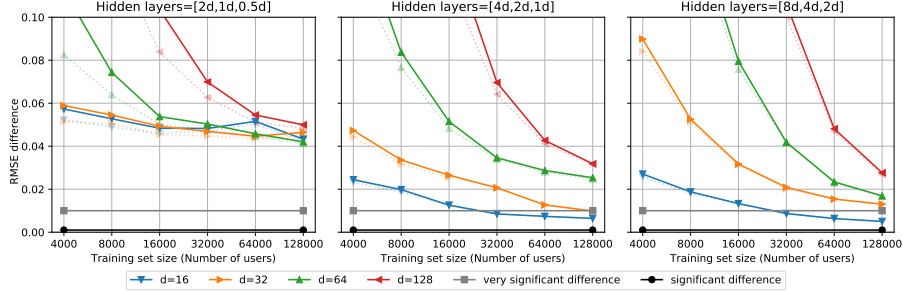


Figure 3: How well a MLP can learn a dot product over embeddings of dimension d . The ground truth is generated from a dot product of Gaussian embeddings plus Gaussian label noise. The graphs show the difference between the RMSE of the dot product and the RMSE of the learned similarity measure; the solid line measures the difference on the fresh set, the dotted on the test set. Noise and scale have been chosen such that 0.01 could indicate a very significant difference and 0.001 a significant difference.

noise $\sigma_{label} = 0.85$. For the Netflix prize, the trivial model that predicts always the average rating has an RMSE of 1.13. Thus we set $\sigma_{emb}^2 = \sqrt{\frac{1.13^2 - 0.85^2}{d}}$. With this setup, the trivial model in our experiment has the same RMSE as the trivial model on Netflix. By aligning both the trivial model and the noise to the Netflix prize, absolute differences in our experiment give some indication of the scale of acceptable errors. In both Netflix and ML 10M, a difference in RMSE of 0.01 is considered very large. For example, for the Netflix prize it took the community about a year⁴ to lower the RMSE from 0.8712 to 0.8616. Similarly, for Movielens 10M, it took about 4 years to lower the RMSE from 0.7815 to 0.7634. Much smaller differences have been published. For example many published increments on Movielens 10M are about 0.001 [33]. We will use these thresholds of 0.01 and 0.001 as an indication whether the approximation errors are acceptable in our experiments. While this is not a perfect comparison, we hope that it can serve as a reasonable indicator.

4.2 Results

Figure 3 shows the approximation error of the MLP for different choices of embedding dimensions and as a function of training data. The figure suggests that with enough training data and wide enough hidden layers, an MLP can approximate a dot product. This holds for embeddings that have been seen in the training data as well as for fresh embeddings. However, consistent with the theory, the number of samples needed scales polynomially with the increasing dimensions and reduced error. Anecdotally, we observe the number of samples needed is about $O(d/\epsilon)^\alpha$ for $1 \leq \alpha \leq 2$. The experiments clearly indicate that

⁴https://www.netflixprize.com/leaderboard_quiz.html

it becomes increasingly difficult for an MLP to fit the dot product function with increasing dimensions. In all cases, the approximation error is well above what is considered a large difference for problems with comparable scale. For example, for the moderate $d = 128$, with 128000 users, the error is still above 0.02, much higher than the *very significant* difference of 0.01.

This experiment shows the difficulty of using an MLP to approximate the dot product, even when explicitly trained to do so. Hence, if the dot product performs well on a given task, there could be a significant price to pay for an MLP to approximate it. We hope this can explain, at least partially, why the dot product model outperforms the MLP model in the experiments of Section 3.3.

5 Applicability of Dot Product Models

Most academic studies focus on training runtime when discussing applicability. However, in industrial applications, the serving runtime is often more important, in particular when the recommendations cannot be precomputed offline but need to be computed at the time of the user’s request. This is the case for most context-aware recommenders in which the recommendation depends on contextual features that are only available at query time. For instance, consider a sequential recommender that recommends items to a user based on the previously selected L items. Here the top scoring items cannot be precomputed for all possible combinations of L items. Instead the recommender would need to retrieve the highest scoring items from the whole item catalogue with a latency of a few milliseconds after the user’s request. Such real time retrieval is a common application in real world recommender systems [6].

Computing a dot product similarity takes $\mathcal{O}(d)$ time while computing an MLP-learned similarity takes $\mathcal{O}(d^2)$ time. If there are n items to score, then the total costs are $\mathcal{O}(dn)$ (for dot) vs $\mathcal{O}(d^2n)$ (for MLP). For large scale applications, n is typically in the range of millions and d is in the hundreds, and while dot has a lower complexity, both are impractical for retrieval applications that require latencies of a few milliseconds. However, for a dot product, the problem of finding the top scoring items can be approximated efficiently. Indeed, given the user embedding \mathbf{p} , the problem is to find items i that maximize $\langle \mathbf{p}, \mathbf{q}_i \rangle$. This is a well-studied problem, known as *approximate nearest neighbor search* [26] or *maximum inner product search* [34]. Efficient sublinear time algorithms exist that makes dot product retrieval feasible in typically a few milliseconds, even with millions of items n [6]. To the best of our knowledge, no such sublinear techniques exist for nearest neighbor retrieval with MLPs.

To summarize, MLP similarity is not applicable for real time top-N recommenders, while the dot product allows fast retrieval using well established nearest neighbor search algorithms.

6 Related Work

6.1 Dot products at the Output Layer of DNNs

At first glance it might appear that our work questions the use of neural networks in recommender systems. This is not the case, and as we will discuss now, many of the most competitive neural networks use a dot product for the output but not an MLP. Consider the general multiclass classification task where (\mathbf{x}, y) is a labeled training example with input \mathbf{x} and label $y \in \{1, \dots, n\}$. A common approach is to define a DNN \mathbf{f} that maps the input \mathbf{x} to a representation (or embedding) $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^d$. At the final stage, this representation is combined with the class labels to produce a vector of scores. Commonly, this is done by multiplying the input representation $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^d$ with a class matrix $Q \in \mathbb{R}^{n \times d}$ to obtain a scalar score for each of the n classes. This vector is then used in the loss function, for example as logits in a softmax cross entropy with the label y . This falls exactly under the family of models discussed in this paper, where $\mathbf{p} = \mathbf{f}(\mathbf{x}) \in \mathbb{R}^d$ and the classes are the items. In fact, the model as described above is a dot product model because at the output $Q\mathbf{f}(\mathbf{x}) = Q\mathbf{p} = [\langle \mathbf{p}, \mathbf{q}_i \rangle]_{i=1}^n$ which means each input-label or user-item combination is a dot product between an input (or user) embedding and label (or item) embedding. This dot product combination of input and class representation is commonly used in sophisticated DNNs for image classification [23, 14] and for natural language processing [4, 28, 9]. This makes our findings that a dot product is a powerful embedding combiner well aligned with the broader DNN community where it is common to apply a dot product at the output for multiclass classification.

6.2 MLPs at the Output Layer of DNNs

NeuMF is very closely related to the previously proposed *neural network matrix factorization* [11]. Neural network matrix factorization also uses a combination of an MLP plus extra embeddings with an explicit dot product like structure as in GMF. A follow up paper [15] proposes to replace the MLP in NCF by an outerproduct and pass this matrix through a convolutional neural network. Finding the dot product with this technique is trivial because the sum of the diagonal in the outerproduct is the dot product. Unfortunately, while written by the same authors as the NCF paper, it evaluates on different data, so our results in Section 3.3 cannot be compared to their numbers and it remains unclear if their work improves over a well tuned baseline with a dot product. Besides prediction quality, this proposal suffers from the same applicability issues as the MLP (see Section 5).

6.3 Specialized Structures inside a DNN

In DNN modeling it is very common to replace an MLP by a more specialized structure that has an inductive bias that represents the problem better. For example, in image classification structures such as convolutional neural networks

are very popular because they represent the spatial structure of the input data. In recurrent neural networks, such parameter sharing is very important too. Another example are attention models, e.g. in Neural Machine Translation [36] and in the Transformer model [35], that contain a matrix product inside the neural network for combining multiple inputs – they can be regarded as the dot product model for combining “internal” embeddings too. All these specialized structures are crucial for advancing the state of the art of deep learning, although in theory they can all be approximated by MLPs.

The inefficiency of MLPs to capture dot and tensor products has been studied by [5] in the context of recommender systems. Here the authors examine how to add context to recurrent neural networks. Similar to our work and Section 4.1, [5] points out that MLPs do not model multiplications and it investigates approximating dot products and tensor products with MLPs empirically. Their study focuses on the model size required to learn a tensor product for embeddings of dimension $d = 1$ and $d = 2$, where the number of distinct embeddings is 100 per mode and the training error is measured.

6.4 Experimental Issues in Recommender Systems

In their meta study, [8] point out issues with evaluation in recommender system research. Their experiments also cover the NCF paper. They show that well studied baselines can get comparable results to (a reproducible value of) NeuMF (see Section 3.4). The goal of our study and [8] is different. While [8] covers a broad set of methods and publications, we are investigating the specific issue of learned similarity functions in more detail. Our work provides apples to apples comparisons of dot product vs MLP, stronger results (outperforming the original NCF results), and a thorough investigation of the reasons and consequences.

7 Conclusion

Our findings indicate that **a dot product might be a better** default choice for combining embeddings than learned similarities using MLP or NeuMF. Shifting the focus in the recommender system research community from learned similarities to dot products might have several positive effects: (1) The research becomes more relevant for the industry because models are applicable (see Section 5). (2) Dot product similarity simplifies modeling and learning (no pretraining, no need for large datasets) which facilitates both experimentation and understanding. (3) Better alignment with other research areas such as natural language processing or image models where the dot product is commonly used.

Finally, our experiments give further evidence that running machine learning methods properly is difficult [33] and one-off studies are prone to drawing wrong conclusions. Introducing shared benchmarks might help to better identify improvements.

References

- [1] ALLEN-ZHU, Z., LI, Y., AND SONG, Z. A convergence theory for deep learning via over-parameterization. In *Proceedings of the 36th International Conference on Machine Learning* (2019), pp. 242–252.
- [2] ANDONI, A., PANIGRAHY, R., VALIANT, G., AND ZHANG, L. Learning polynomials with neural networks. In *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32* (2014), ICML’14, JMLR.org, p. II–1908–II–1916.
- [3] BARRON, A. R. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information theory* 39, 3 (1993), 930–945.
- [4] BENGIO, Y., DUCHARME, R., VINCENT, P., AND JAUVIN, C. A neural probabilistic language model. *Journal of machine learning research* 3, Feb (2003), 1137–1155.
- [5] BEUTEL, A., COVINGTON, P., JAIN, S., XU, C., LI, J., GATTO, V., AND CHI, E. H. Latent cross: Making use of context in recurrent recommender systems. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining* (New York, NY, USA, 2018), WSDM ’18, Association for Computing Machinery, p. 46–54.
- [6] COVINGTON, P., ADAMS, J., AND SARGIN, E. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems* (New York, NY, USA, 2016), RecSys ’16, Association for Computing Machinery, p. 191–198.
- [7] CYBENKO, G. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems* 2, 4 (1989), 303–314.
- [8] DACREMA, M. F., BOGLIO, S., CREMONESI, P., AND JANNACH, D. **A troubling analysis of reproducibility and progress in recommender systems research, 2019.**
- [9] DEVLIN, J., CHANG, M.-W., LEE, K., AND TOUTANOVA, K. Bert: Pre-training of deep bidirectional transformers for language understanding, 2018.
- [10] DU, S., LEE, J., LI, H., WANG, L., AND ZHAI, X. Gradient descent finds global minima of deep neural networks. In *Proceedings of the 36th International Conference on Machine Learning* (2019), pp. 1675–1685.
- [11] DZIUGAITE, G. K., AND ROY, D. M. Neural network matrix factorization, 2015.
- [12] GENG, X., ZHANG, H., BIAN, J., AND CHUA, T. Learning image and user features for recommendation in social networks. In *2015 IEEE International Conference on Computer Vision (ICCV)* (2015), pp. 4274–4282.

- [13] HARPER, F. M., AND KONSTAN, J. A. The movielens datasets: History and context. *ACM Trans. Interact. Intell. Syst.* 5, 4 (Dec. 2015).
- [14] HE, K., ZHANG, X., REN, S., AND SUN, J. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (Jun 2016).
- [15] HE, X., DU, X., WANG, X., TIAN, F., TANG, J., AND CHUA, T.-S. Outer product-based neural collaborative filtering. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18* (7 2018), International Joint Conferences on Artificial Intelligence Organization, pp. 2227–2233.
- [16] HE, X., LIAO, L., ZHANG, H., NIE, L., HU, X., AND CHUA, T.-S. Neural collaborative filtering. In *Proceedings of the 26th International Conference on World Wide Web* (Republic and Canton of Geneva, Switzerland, 2017), WWW ’17, International World Wide Web Conferences Steering Committee, pp. 173–182.
- [17] HORNIK, K., STINCHCOMBE, M., WHITE, H., ET AL. Multilayer feedforward networks are universal approximators. *Neural networks* 2, 5 (1989), 359–366.
- [18] HU, B., SHI, C., ZHAO, W. X., AND YU, P. S. Leveraging meta-path based context for top- n recommendation with a neural co-attention model. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (New York, NY, USA, 2018), KDD ’18, Association for Computing Machinery, p. 1531–1540.
- [19] HU, Y., KOREN, Y., AND VOLINSKY, C. Collaborative filtering for implicit feedback datasets. In *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining* (2008), ICDM ’08, pp. 263–272.
- [20] JAWARNEH, I. M. A., BELLAVISTA, P., CORRADI, A., FOSCHINI, L., MONTANARI, R., BERROCAL, J., AND MURILLO, J. M. A pre-filtering approach for incorporating contextual information into deep learning based recommender systems. *IEEE Access* 8 (2020), 40485–40498.
- [21] KOREN, Y. The bellkor solution to the netflix grand prize, 2009.
- [22] KOREN, Y., AND BELL, R. *Advances in Collaborative Filtering*. Springer US, Boston, MA, 2011, pp. 145–186.
- [23] KRIZHEVSKY, A., SUTSKEVER, I., AND HINTON, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*. 2012, pp. 1097–1105.
- [24] LEVY, M., AND JACK, K. Efficient top-n recommendation by linear regression. In *RecSys Large Scale Recommender Systems Workshop* (2013).

- [25] LI, D., CHEN, C., LIU, W., LU, T., GU, N., AND CHU, S. Mixture-rank matrix approximation for collaborative filtering. In *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017, pp. 477–485.
- [26] LIU, T., MOORE, A. W., GRAY, A., AND YANG, K. An investigation of practical approximate nearest neighbor algorithms. In *Proceedings of the 17th International Conference on Neural Information Processing Systems* (Cambridge, MA, USA, 2004), NIPS’04, MIT Press, p. 825–832.
- [27] MATTSON, P., CHENG, C., COLEMAN, C., DIAMOS, G., MICIKEVICIUS, P., PATTERSON, D., TANG, H., WEI, G.-Y., BAILIS, P., BITTORF, V., BROOKS, D., CHEN, D., DUTTA, D., GUPTA, U., HAZELWOOD, K., HOCK, A., HUANG, X., IKE, A., JIA, B., KANG, D., KANTER, D., KUMAR, N., LIAO, J., MA, G., NARAYANAN, D., OGUNTEBI, T., PEKHIMENKO, G., PENTECOST, L., REDDI, V. J., ROBIE, T., JOHN, T. S., TABARU, T., WU, C.-J., XU, L., YAMAZAKI, M., YOUNG, C., AND ZAHARIA, M. Mlperf training benchmark, 2019.
- [28] MIKOLOV, T., SUTSKEVER, I., CHEN, K., CORRADO, G. S., AND DEAN, J. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems* (2013), pp. 3111–3119.
- [29] NING, X., AND KARYPIS, G. Slim: Sparse linear methods for top-n recommender systems. In *2011 IEEE 11th International Conference on Data Mining* (2011), IEEE, pp. 497–506.
- [30] NIU, W., CAVERLEE, J., AND LU, H. Neural personalized ranking for image recommendation. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining* (New York, NY, USA, 2018), WSDM ’18, Association for Computing Machinery, p. 423–431.
- [31] PATEREK, A. Improving regularized singular value decomposition for collaborative filtering. In *Proceedings of KDD cup and workshop* (2007), vol. 2007, pp. 5–8.
- [32] QIN, J., REN, K., FANG, Y., ZHANG, W., AND YU, Y. Sequential recommendation with dual side neighbor-based collaborative relation modeling. In *Proceedings of the 13th International Conference on Web Search and Data Mining* (New York, NY, USA, 2020), WSDM ’20, Association for Computing Machinery, p. 465–473.
- [33] RENDLE, S., ZHANG, L., AND KOREN, Y. On the difficulty of evaluating baselines: A study on recommender systems. *CoRR abs/1905.01395* (2019).

- [34] SHRIVASTAVA, A., AND LI, P. Asymmetric lsh (alsh) for sublinear time maximum inner product search (mips). In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2* (Cambridge, MA, USA, 2014), NIPS’14, MIT Press, p. 2321–2329.
- [35] VASWANI, A., SHAZEER, N., PARMAR, N., USZKOREIT, J., JONES, L., GOMEZ, A. N., KAISER, L., AND POLOSUKHIN, I. Attention is all you need. In *Advances in neural information processing systems* (2017), pp. 5998–6008.
- [36] WU, Y., SCHUSTER, M., CHEN, Z., LE, Q. V., NOROUZI, M., MACHEREY, W., KRIKUN, M., CAO, Y., GAO, Q., MACHEREY, K., ET AL. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144* (2016).
- [37] ZAMANI, H., AND CROFT, W. B. Learning a joint search and recommendation model from user-item interactions. In *Proceedings of the 13th International Conference on Web Search and Data Mining* (New York, NY, USA, 2020), WSDM ’20, Association for Computing Machinery, p. 717–725.
- [38] ZHAO, X., ZHU, Z., ZHANG, Y., AND CAVERLEE, J. Improving the estimation of tail ratings in recommender system with multi-latent representations. In *Proceedings of the 13th International Conference on Web Search and Data Mining* (New York, NY, USA, 2020), WSDM ’20, Association for Computing Machinery, p. 762–770.

A Experiments from NCF Paper

This section provides details about our setup for establishing a dot product baseline for the experiments of the NCF paper (Section 3.3). The code and datasets of NCF were provided by its authors⁵. We provide code for our implementation of matrix factorization and the script to generate the tuning split⁶.

A.1 Model and Optimization

We implemented a matrix factorization with bias (see Eq. 2). The parameters of this model are the embeddings $P \in \mathbb{R}^{M \times d}$ for M users and $Q \in \mathbb{R}^{N \times d}$ for N items. Following the NCF paper, for training we cast the implicit data, which contains only positive observations, into a binary two class classification problem and sample m negative items for each tuple in the implicit data. In each epoch a new set of negatives is drawn – the sampling distribution is uniform.

⁵https://github.com/hexiangnan/neural_collaborative_filtering

⁶https://github.com/google-research/google-research/tree/master/dot_vs_learned_similarity

We minimize the **binary logistic loss with L2 regularization**. For each training example (u, i, y) where $y \in \{0, 1\}$ is the binary label, the regularized loss is

$$l(u, i, y) = -y \ln \sigma(\phi(\mathbf{p}_u, \mathbf{q}_i)) - (1 - y) \ln(1 - \sigma(\phi(\mathbf{p}_u, \mathbf{q}_i))) + \lambda \|\mathbf{p}_u\| + \lambda \|\mathbf{q}_i\| \quad (8)$$

with the regularization constant $\lambda \in \mathbb{R}^+$. The loss is optimized with stochastic gradient descent with learning rate η , with the update rules:

$$p_{u,1} \leftarrow p_{u,1} - \eta[(\sigma(\phi(\mathbf{p}_u, \mathbf{q}_i)) - y) + \lambda p_{u,1}] \quad (9)$$

$$q_{i,1} \leftarrow q_{i,1} - \eta[(\sigma(\phi(\mathbf{p}_u, \mathbf{q}_i)) - y) + \lambda q_{i,1}] \quad (10)$$

$$\mathbf{p}_{u,[2,\dots,d]} \leftarrow \mathbf{p}_{u,[2,\dots,d]} - \eta[(\sigma(\phi(\mathbf{p}_u, \mathbf{q}_i)) - y) \mathbf{q}_{i,[2,\dots,d]} + \lambda \mathbf{p}_{u,[2,\dots,d]}] \quad (11)$$

$$\mathbf{q}_{i,[2,\dots,d]} \leftarrow \mathbf{q}_{i,[2,\dots,d]} - \eta[(\sigma(\phi(\mathbf{p}_u, \mathbf{q}_i)) - y) \mathbf{p}_{u,[2,\dots,d]} + \lambda \mathbf{q}_{i,[2,\dots,d]}] \quad (12)$$

The embeddings are initialized from a normal distribution. This configuration shares the same loss, regularization, negative sampling approach, and initialization procedure with MLP and NeuMF as proposed in [16].

The hyperparameters of the dot product model are: embedding dimension d , regularization λ , learning rate η , number of negative samples m , number of training epochs, standard deviation for initialization. Analogously to the NCF paper, we report results for $d \in \{16, 32, 64, 96, 128, 192\}$.

A.2 Hyperparameter Tuning

We create a tuning dataset that follows the same splitting protocol as the final training/test split. In particular, we remove the last feedback from each user from the training set and place it in a test set for tuning and keep the remaining training cases in a training set for tuning. We then train models on the training set for tuning and evaluate the model on the test set for tuning. We choose all hyperparameters including the number of training epochs on this tuning set. Note that both the training set for tuning and test set for tuning contain no information about the final test set.

From our past experience with matrix factorization models, if the other hyperparameters are chosen properly, then the larger the embedding dimension the better the quality – our experiments Figure 2 confirm this. For the other hyperparameters: learning rate and number of training epochs influence the convergence curves. Usually, the lower the learning rate, the better the quality but also the more epochs are needed. We set a computational budget of up to 256 epochs and search for the learning rate within this setting. In the first hyperparameter pass, we search a coarse grid of learning rates $\eta \in \{0.001, 0.003, 0.01\}$ and number of negatives $m = \{4, 8, 16\}$ while fixing the regularization to $\lambda = 0$. Then we did a search for regularization in $\{0.001, 0.003, 0.01\}$ around the promising candidates. To speed up the search, these first coarse passes were done with 128 epochs and a fixed dimension of $d = 64$ (Movielens) and $d = 128$ (Pinterest). We did further refinements around the most promising values of learning rate, number of negatives and regularization using $d = 128$ and 256 epochs.

Throughout the experiments we initialize embeddings from a Gaussian distribution with standard deviation of 0.1; we tested some variation of the standard deviation but did not see much effect.

The final hyperparameters for Movielens are: learning rate $\eta = 0.002$, number of negatives $m = 8$, regularization $\lambda = 0.005$, number of epochs 256. For Pinterest: learning rate $\eta = 0.007$, number of negative samples $m = 10$, regularization $\lambda = 0.01$, number of epochs 256.

After hyperparameter selection, we trained on the full dataset with these hyperparameters and evaluated according to the protocol in [16]. We repeated the final training and evaluation 8 times and report the mean of the metrics.

A.3 MLP and NeuMF Results

We report the results for MLP and NeuMF from the original NCF paper [16]. As we share the same evaluation protocol and splits, the numbers are comparable. We report the results for NeuMF from Table 2 in [16] and the results for MLP from Tables 3,4 in [16] using the ‘MLP-3’ setting.

It should be noted that in [16], the tables and plots use “predictive factor” instead of embedding dimension. The predictive factor is defined as the size of the last hidden layer of the MLP, and as described in [16], for the 3-layer MLP a predictive factor of k operates on two input embeddings, each of dimension $d = 2k$. For the NeuMF model, a predictive factor of k operates on embeddings of dimension $d = 3k$ because it consists of an independent MLP with predictive factor of k (embedding size $d = 2k$) and a GMF with embedding size $d = k$. This definition of *predictive factor* is arbitrary, in fact it can be made arbitrarily small by adding layers to the MLP without changing anything else in the model. We think it is more meaningful to compare models with a fixed embedding dimension. In particular, we want to investigate the prediction quality of an MLP or a dot product over two embeddings of the same size d . We recast all results from the NCF paper in terms of embedding dimension, by multiplying the predictive factor by 3 for NeuMF results and by 2 for MLP results. This allows us to do an apples to apples comparison of different similarity functions over an embedding space of dimension d .