# Graph Convolutional Matrix Completion

**Rianne van den Berg**
University of Amsterdam
r.vandenberg@uva.nl

**Thomas N. Kipf**
University of Amsterdam
t.n.kipf@uva.nl

**Max Welling**
University of Amsterdam, CIFAR*
m.welling@uva.nl

## Abstract

In this paper we revisit matrix completion for recommender systems from the point of view of link prediction on graphs. Interaction data such as movie ratings can be represented by a bipartite user-item graph with labeled edges representing observed ratings. Building on recent progress in deep learning on graph-structured data, we propose a graph auto-encoder framework based on differentiable message passing on the bipartite interaction graph. This framework can be viewed as an important first step towards end-to-end learning in settings where the interaction data is integrated into larger graphs such as social networks or knowledge graphs, circumventing the need for multistage frameworks. Our model achieves competitive performance on standard collaborative filtering benchmarks, significantly outperforming related methods in a recommendation task with side information.

## 1 Introduction

With the explosive growth of e-commerce and social media platforms, recommendation algorithms have become indispensable tools for many businesses. Two main branches of recommender algorithms are often distinguished. On the one hand, content-based recommender systems use available information of users and items, such as their respective occupation and genre, to predict the next purchase of a user or rating of an item. On the other hand, collaborative filtering models take into account the collective interaction data to predict future ratings or purchases, for instance by considering user-user or item-item similarities based on their rating patterns.

Collaborative filtering algorithms can be subdivided further into memory-based models and learning-based models. The former class computes recommendations based on user-user similarity or item-item similarity measures by acting directly on the rating data, such as the Pearson correlation coefficient or cosine-based similarities. Learning-based models try to learn user preferences or item features based on the collective interaction data. One particular successful class of learning-based collaborative filtering algorithms is that of matrix factorization models [13], such as probabilistic matrix factorization proposed by Salakhutdinov & Mnih in 2007 [17]. Within these methods, the matrix representing the interactions between users and items, e.g. clicks, purchases or ratings, is approximated by a low rank matrix $U^T V$. The rows of $U$ and $V$ represent user and item latent representations, reflecting for instance user interests and item genres based on the collective interaction data.

In this paper, we revisit the idea of treating matrix completion as a link prediction problem on graphs, as suggested among others by Li & Chen [15]. The interaction data in collaborative filtering can be represented by a bipartite graph between user and item nodes, with observed ratings/purchases represented by links. Here, we focus on datasets containing explicit ratings, and leave implicit interactions such as purchases or clicks for future work. In the case of explicit ratings, the absence of a link can be more easily interpreted as an unobserved rating between a user and item that have not yet interacted. The task of predicting ratings can thus be mapped to predicting links in the bipartite user-item graph, labeled with an associated rating level.

---

*Canadian Institute for Advanced Research

By extending recent progress in graph based methods for semi-supervised classification and link prediction [1, 5, 16, 4, 12, 26, 11], we propose an end-to-end graph-based auto-encoder model for matrix completion. The model produces latent features of user and item nodes through message passing on the bipartite interaction graph, after which these latent user and item representations are used to reconstruct the rating links through a bilinear decoder.

The benefits of formulating matrix completion as link prediction on a bipartite graph become especially apparent when considering the more general situation where the recommender graph is embedded into a much larger graph, such as a social network connecting users or a knowledge graph connecting items with external information. Our graph auto-encoder model serves as an important first step towards enabling end-to-end learning on such large-scale graphs for recommendation.

The paper is structured as follows: in Section 2 we introduce our graph auto-encoder model for matrix completion. Section 3 discusses related work. Experimental results are shown in Section 4, and conclusion and future research directions are discussed in Section 5.

## 2 Matrix completion as link prediction in bipartite graphs

Consider a rating matrix $M$ of shape $N_u \times N_v$, where $N_u$ is the number of users and $N_v$ is the number of items. Entries $M_{ij}$ in this matrix encode either an observed rating (user $i$ rated item $j$) from a set of discrete possible rating values, or the fact that the rating is unobserved (encoded by the value 0). See Figure 1 for an illustration. The task of matrix completion or recommendation can be seen as predicting the value of unobserved entries in $M$.
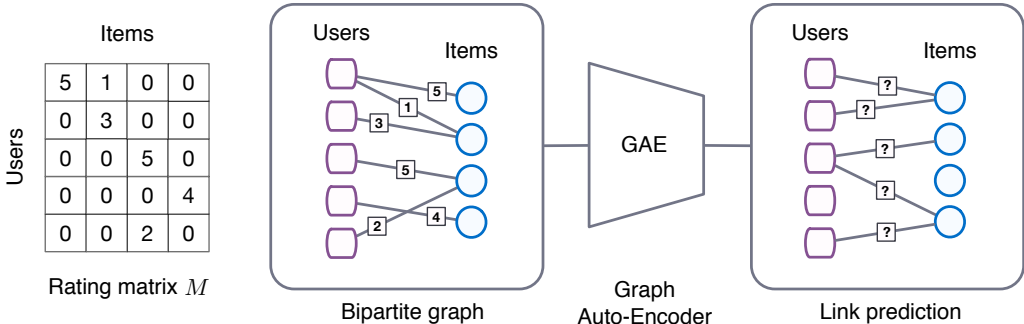


Figure 1: *Left*: Rating matrix $M$ with entries that correspond to user-item interactions (ratings between 1-5) or missing observations (0). *Right*: User-item interaction graph with bipartite structure. Edges correspond to interaction events, numbers on edges denote the rating a user has given to a particular item. The matrix completion task (i.e. predictions for unobserved interactions) can be cast as a link prediction problem and modeled using an end-to-end trainable graph auto-encoder.

In an equivalent picture, matrix completion or recommendation can be cast as a link prediction problem on a bipartite user-item interaction graph. More precisely, the interaction data can be represented by an undirected graph $G = (\mathcal{V}, \mathcal{E}, \mathcal{R})$ with entities consisting of a collection of user nodes $u_i \in \mathcal{U}$ with $i \in \{1, ..., N_u\}$, and item nodes $v_j \in \mathcal{V}$ with $j \in \{1, ..., N_v\}$. The edges $(u_i, r, v_j) \in \mathcal{E}$ carry labels that represent ordinal rating levels, such as $r \in \{1, ..., R\} = \mathcal{R}$. This connection was previously explored in [15] and led to the development of graph-based methods for recommendation.

Previous graph-based approaches for recommender systems (see [15] for an overview) employ a multi-stage pipeline, typically consisting of a graph feature extraction model and a link prediction model, all of which are trained separately. Recent results, however, have shown that results can often be significantly improved by modeling graph-structured data with end-to-end learning techniques [1, 5, 16, 20, 4, 12, 18] and specifically with graph auto-encoders [26, 11] for unsupervised learning and link prediction. In what follows, we introduce a specific variant of graph auto-encoders for the task of recommendation.

## 2.1 Graph auto-encoders

We revisit graph auto-encoders which were originally introduced in [26, 11] as an end-to-end model for unsupervised learning [26] and link prediction [11] on undirected graphs. We specifically consider the setup introduced in [11], as it makes efficient use of (convolutional) weight sharing and allows for inclusion of side information in the form of node features. Graph auto-encoders are comprised of 1) a graph encoder model $Z = f(X, A)$, which take as input an $N \times D$ feature matrix $X$ and a graph adjacency matrix $A$, and produce an $N \times E$ node embedding matrix $Z = [z_1^T, \ldots, z_N^T]^T$, and 2) a pairwise decoder model $\check{A} = g(Z)$, which takes pairs of node embeddings $(z_i, z_j)$ and predicts respective entries $\check{A}_{ij}$ in the adjacency matrix. Note that $N$ denotes the number of nodes, $D$ the number of input features, and $E$ the embedding size.

For bipartite recommender graphs $G = (\mathcal{V}, \mathcal{E}, \mathcal{R})$, we can reformulate the encoder as $[U, V] = f(X, M_1, \ldots, M_R)$, where $M_r \in \{0, 1\}^{N_u \times N_v}$ is the adjacency matrix associated with rating type $r \in \mathcal{R}$, such that $M_r$ contains 1's for those elements for which the original rating matrix $M$ contains observed ratings with value $r$. $U$ and $V$ are now matrices of user and item embeddings with shape $N_u \times E$ and $N_v \times E$, respectively. A single user (item) embedding takes the form of a real-valued vector $U_{i,:}$ ($V_{j,:}$) for user $i$ (item $j$). The specific functional form and parameterization of $f(\cdot)$ and the choice of $X$ are yet to be defined.

Analogously, we can reformulate the decoder as $\check{M} = g(U, V)$, i.e. as a function acting on the user and item embeddings and returning a (reconstructed) rating matrix $\check{M}$ of shape $N_u \times N_v$. We can train this graph auto-encoder by minimizing the reconstruction error between the predicted ratings in $\check{M}$ and the observed ground-truth ratings in $M$. Examples of metrics for the reconstruction error are the root mean square error, or the cross entropy when treating the rating levels as different classes.
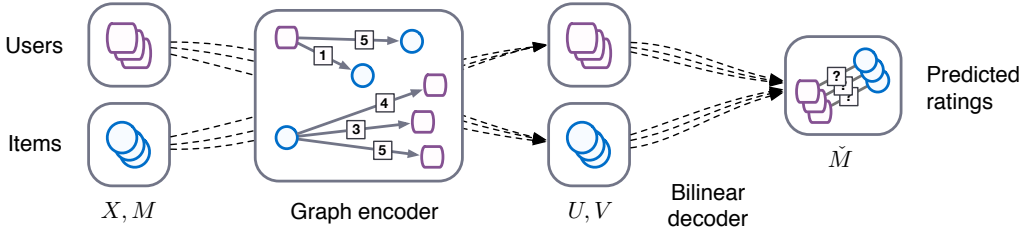


Figure 2: Schematic of a forward-pass through the MC-GC model, which is comprised of a graph convolutional encoder $[U, V] = f(X, M_1, \ldots, M_R)$ that passes and transforms messages from user to item nodes, and vice versa, followed by a bilinear decoder model that predicts entries of the (reconstructed) rating matrix $\check{M} = g(U, V)$, based on pairs of user and item embeddings.

We shall note at this point that several recent state-of-the-art models for matrix completion [14, 23, 6, 28] can be cast into this framework and understood as a special case of our model. An overview of these models is provided in Section 3.

## 2.2 Graph convolutional encoder

In what follows, we propose a particular choice of encoder model that makes efficient use of weight sharing across locations in the graph and that assigns separate processing channels for each edge type (or rating type) $r \in \mathcal{R}$. The form of weight sharing is inspired by a recent class of convolutional neural networks that operate directly on graph-structured data [1, 5, 4, 12], in the sense that our model performs local operations that only take the first-order neighborhood of a node into account, whereby the same transformation is applied across all locations in the graph.

This type of local graph convolution can be seen as a form of message passing [3, 7], where vector-valued messages are being passed and transformed across edges of the graph. In our case, we can assign a specific transformation for each rating level, resulting in edge-type specific messages $\mu_{j \to i, r}$ from items $j$ to users $i$ of the following form:

$$\mu_{j \to i, r} = \frac{1}{c_{ij}} W_r x_j .\tag{1}$$

Here, $c_{ij}$ is a normalization constant, which we choose to either be $1/|\mathcal{N}_i|$ (left normalization) or $1/\sqrt{|\mathcal{N}_i||\mathcal{N}_j|}$ (symmetric normalization) with $\mathcal{N}_i$ denoting the set of neighbors of node $i$. $W_r$ is an edge-type specific parameter matrix and $x_j$ is the (initial) feature vector of node $j$. Messages $\mu_{i \to j, r}$ from users to items are processed in an analogous way. After the message passing step, we accumulate incoming messages at every node by summing over all neighbors $\mathcal{N}_{i,r}$ under a specific edge-type $r$, and by subsequently accumulating them into a single vector representation:

$$h_i = \sigma\left[\operatorname{accum}\left(\sum_{j \in \mathcal{N}_{i,1}} \mu_{j \to i, 1}, \ldots, \sum_{j \in \mathcal{N}_{i,R}} \mu_{j \to i, R}\right)\right], \tag{2}$$

where $\operatorname{accum}(\cdot)$ denotes an accumulation operation, such as $\operatorname{stack}(\cdot)$, i.e. a concatenation of vectors (or matrices along their first dimension), or $\operatorname{sum}(\cdot)$, i.e. summation of all messages. $\sigma(\cdot)$ denotes an element-wise activation function such as the $\operatorname{ReLU}(\cdot) = \max(0, \cdot)$. To arrive at the final embedding of user node $i$, we transform the intermediate output $h_i$ as follows:

$$u_i = \sigma(W h_i). \tag{3}$$

The item embedding $v_i$ is calculated analogously with the same parameter matrix $W$. We will refer to (2) as a *graph convolution* layer and to (3) as a *dense* layer. Note that deeper models can be built by chaining multiple of these layers (in arbitrary combinations) with appropriate activation functions.

It is worth mentioning that the model demonstrated here is only one particular possible, yet comparatively simple choice of an encoder, and other variations are potentially worth exploring. Instead of a simple linear message transformation, one could explore variations where $\mu_{j \to i, r} = nn(x_i, x_j, r)$ is a neural network in itself. Instead of choosing a specific normalization constant for individual messages, such as done here, one could employ some form of attention mechanism, where the individual contribution of each message is learned and determined by the model.

## 2.3 Bilinear decoder

For the reconstruction of the links in the bipartite interaction graph we consider a bilinear decoder, and treat each rating level as a separate class. Indicating the reconstructed rating between user $i$ and item $j$ with $\check{M}_{ij}$, the decoder produces a probability distribution over possible rating levels through a bilinear operation followed by the application of a $\operatorname{softmax}$ function:

$$p(\check{M}_{ij} = r) = \frac{e^{u_i^T Q_r v_j}}{\sum_{s \in R} e^{u_i^T V_s v_j}}, \tag{4}$$

where $Q_r$ is a trainable parameter matrix of shape $E \times E$, where $E$ is the dimensionality of hidden user (item) representations $u_i$ ($v_j$). The predicted rating is computed as

$$\check{M}_{ij} = g(u_i, v_j) = \mathbb{E}_{p(\check{M}_{ij}=r)}[r] = \sum_{r \in R} r \, p(\check{M}_{ij} = r). \tag{5}$$

## 2.4 Model training

**Loss function**   During model training, we minimize the following negative log likelihood of the predicted ratings $\check{M}_{ij}$:

$$\mathcal{L} = -\sum_{i,j; \boldsymbol{\Omega}_{ij}=1} \sum_{r=1}^{R} I_{r, M_{ij}} \log p(\check{M}_{ij} = r), \tag{6}$$

with $I_{k,l} = 1$ when $k = l$ and $I_{k,l} = 0$ otherwise. The matrix $\boldsymbol{\Omega} \in \{0,1\}^{N_u \times N_i}$ serves as a mask for unobserved ratings, such that ones occur for elements corresponding to observed ratings in $M$, and zeros for unobserved ratings. Hence, we only optimize over observed ratings.

**Node dropout**   In order for the model to generalize well to unobserved ratings, it is trained in a denoising setup by randomly dropping out all outgoing messages of a particular node, with a probability $p_{\mathrm{dropout}}$, which we will refer to as *node dropout*. Note that messages are rescaled after dropout as in [24]. In contrast to independently dropping out individual outgoing messages, which makes embeddings more robust against the presence or absence of single edges, this form of dropout also causes embeddings to be more independent of particular user or item influences. We furthermore also apply regular dropout [24] to hidden layer units (3).

**Mini-batching**   We introduce mini-batching by sampling contributions to the loss function in Eq. (6) from different observed ratings. That is, we sample only a fixed number of contributions from the sum over user and item pairs. By only considering a fixed number of contributions to the loss function, we can remove respective rows of users and items in $M_1, ..., M_R$ in Eq. (7) that do not appear in the current batch. This serves both as an effective means of regularization, and reduces the memory requirement to train the model, which is necessary to fit the full model for MovieLens-10M into GPU memory.

## 2.5   Vectorized implementation

In practice, we can use efficient sparse matrix multiplications, with complexity linear in the number of edges, i.e. $\mathcal{O}(|\mathcal{E}|)$, to implement the graph auto-encoder model. The graph convolutional encoder (Eq. 3), for example, can be vectorized as follows:

$$\begin{bmatrix} U \\ V \end{bmatrix} = f(X, M_1, \dots, M_R) = \sigma\left(\begin{bmatrix} H_u \\ H_v \end{bmatrix} W^T\right) , \tag{7}$$

$$\text{with} \quad \begin{bmatrix} H_u \\ H_v \end{bmatrix} = \sigma\left(\text{accum}\left(D^{-1}\mathcal{M}_1 X W_1^T, \dots, D^{-1}\mathcal{M}_R X W_R^T\right)\right), \tag{8}$$

$$\text{and} \quad \mathcal{M}_r = \begin{pmatrix} 0 & M_r \\ M_r^T & 0 \end{pmatrix} . \tag{9}$$

Here, $D$ denotes the diagonal node degree matrix with nonzero elements $D_{ii} = |\mathcal{N}_i|^{-1}$. Vectorization of the bilinear decoder follows in an analogous, straightforward manner. Note that it is only necessary to evaluate observed elements in $\tilde{M}$, given by the mask $\mathbf{\Omega}$ in Eq. 6.

## 2.6   Input feature representation and side information

Since the graph encoder preserves the identity of every node, side information can be injected into node representations both at the input-level (i.e. in the form of an input feature matrix $X$) or at later stages in the model, e.g. in the dense transformation layer (3). We generally had most success with including side information in the form of additional user and item feature vectors $x_i^f$ (for node $i$) via a separate processing channel directly into the the dense hidden layer:

$$u_i = \sigma(W h_i + W_2^f f_i) \ , \ \text{with} \quad f_i = \sigma(W_1^f x_i^f), \tag{10}$$

where $W_1^f$ and $W_2^f$ are trainable weight matrices. We represent the input feature matrix $X = [x_1^T, \dots, x_N^T]^T$ as an identity matrix, which assigns a unique one-hot vector to every node in the graph.

In [25], Strub et al. propose to include content information along similar lines, although in their case the proposed model is strictly user- or item-based, and thus only supports side information for either users or items.

Note that side information does not necessarily need to come in the form of per-node feature vectors, but can also be provided in the form of, e.g., graph-structured, natural language, or image data. In this case, the dense layer in (10) is replaced by an appropriate differentiable module, such as a recurrent neural network, a convolutional neural network, or another graph convolutional network.

## 2.7   Ordinal weight sharing

In the collaborative filtering setting with one-hot vectors as input, the columns of the weight matrices $W_r$ play the role of latent factors for each separate node for one specific rating value $r$. These latent factors are passed onto connected user or item nodes through message passing. However, not all users and items necessarily have an equal number of ratings for each rating level. This results in certain columns of $W_r$ to be optimized significantly less frequently than others. Therefore, some form of weight sharing between the matrices $W_r$ for different $r$ is desirable to alleviate this optimization problem. Following [28], we therefore implement the following weight sharing setup:

$$W_r = \sum_{s=1}^{r} T_s . \tag{11}$$

We will refer to this type of weight sharing as ordinal weight sharing due to the increasing number of weight matrices included for higher rating levels.

As an effective means of regularization of the pairwise bilinear decoder, we resort to weight sharing in the form of a linear combination of a set of basis weight matrices $P_s$:

$$Q_r = \sum_{s=1}^{n_b} a_{rs} P_s \, , \tag{12}$$

with $s \in (1, ..., n_b)$ and $n_b$ being the number of basis weight matrices. Here, $a_{rs}$ are the learnable coefficients that determine the linear combination for each decoder weight matrix $Q_r$. Note that in order to avoid overfitting and to reduce the number of parameters, the number of basis weight matrices $n_b$ should naturally be lower than the number of rating levels.

## 3 Related work

**Auto-encoders** User- or item-based auto-encoders [23, 28, 25] are a recent class of state-of-the-art collaborative filtering models that can be seen as a special case of our graph auto-encoder model, where only either user or item embeddings are considered in the encoder. AutoRec by Sedhain et al. [23] is the first such model, where the user's (or item's) partially observed rating vector is projected onto a latent space through an encoder layer, and reconstructed using a decoder layer with mean squared reconstruction error loss.

The CF-NADE algorithm by Zheng et al. [28] can be considered as a special case of the above auto-encoder architecture. In the user-based setting, messages are only passed from items to users, and in the item-based case the reverse holds. Note that in contrast to our model, unrated items are assigned a default rating of 3 in the encoder, thereby creating fully-connected interaction graph. CF-NADE imposes a random ordering on nodes, and splits incoming messages into two sets via a random cut, only one of which is kept. This model can therefore be seen as a denoising auto-encoder, where part of the input space is dropped out at random in every iteration.

**Factorization models** Many of the most popular collaborative filtering algorithms fall into the class of matrix factorization (MF) models. Methods of this sort assume the rating matrix to be well approximated by a low rank matrix: $M \approx UV^T$, with $U \in \mathbb{R}^{N_u \times k}$ and $V \in \mathbb{R}^{N_i \times k}$, with $k \ll N_u, N_i$. The rows of $U$ and $V$ can be seen as latent feature representations of users and items, representing an encoding for their interests through their rating pattern. Probabilistic matrix factorization (PMF) by Salakhutdinov et al. [17] assumes that the ratings contained in $M$ are independent stochastic variables with Gaussian noise. Optimization of the maximum likelihood then leads one to minimize the mean squared error between the observed entries in $M$ and the reconstructed ratings in $UV^T$. BiasedMF by Koren et al. [13] improves upon PMF by incorporating a user and item specific bias, as well as a global bias. Neural network matrix factorization (NNMF) [6] extends the MF approach by passing the latent user and item features through a feed forward neural network. Local low rank matrix approximation by Lee et al. [14], introduces the idea of reconstructing rating matrix entries using different (entry dependent) combinations of low rank approximations.

**Matrix completion with side information** In matrix completion (MC) [2], the objective is to approximate the rating matrix with a low-rank rating matrix. Rank minimization, however, is an intractable problem, and Candes & Recht [2] replaced the rank minimization with a minimization of the nuclear norm (the sum of the singular values of a matrix), turning the objective function into a tractable convex one. Inductive matrix completion (IMC) by Jain & Dhillon, 2013 and Xu et al., 2013 incorporates content information of users and items in feature vectors and approximates the observed elements of the rating matrix as $M_{ij} = x_i^T UV^T y_j$, with $x_i$ and $y_j$ representing the feature vector of user $i$ and item $j$ respectively.

The geometric matrix completion (GMC) model proposed by Kalofolias et al. in 2014 [9] introduces a regularization of the MC model by adding side information in the form of user and item graphs. In [21], a more efficient alternating least squares optimization optimization method (GRALS) is introduced to the graph-regularized matrix completion problem. Most recently, Monti et al. [19] suggested to incorporate graph-based side information in matrix completion via the use of convolutional neural networks on graphs, combined with a recurrent neural network to model the dynamic rating generation

process. Their work is different from ours, in that we model the rating graph directly using a graph convolutional encoder/decoder approach that predicts unseen ratings in a single, non-iterative step.

## 4 Experiments

We evaluate our model on the three standard collaborative filtering datasets[2]: MovieLens 100K, 1M, and 10M. The datasets consist of user ratings for movies collected from the MovieLens website. Dataset statistics are summarized in Table 1. For all experiments, we choose from the following set-

| Dataset | Users | Items | Ratings | Density | Rating levels |
|---|---|---|---|---|---|
| MovieLens 100K (ML-100K) | 943 | 1,682 | 100,000 | 0.0630 | 1, 2, ..., 5 |
| MovieLens 1M (ML-1M) | 6,040 | 3,706 | 1,000,209 | 0.0447 | 1, 2, ..., 5 |
| MovieLens 10M (ML-10M) | 69,878 | 10,677 | 10,000,054 | 0.0134 | 0.5, 1, ..., 5 |

Table 1: Number of users, items and ratings for each of the MovieLens datasets used in our experiments. We further indicate rating density and rating levels.

tings based on validation performance: accumulation function ($\mathrm{stack}$ vs. $\mathrm{sum}$), whether to use ordinal weight sharing, left vs. symmetric normalization, and dropout rate $p_{\mathrm{dropout}} \in \{0.3, 0.4, 0.5, 0.6, 0.7\}$. Unless otherwise noted, we use a Adam [10] with a learning rate of $0.01$, weight sharing in the decoder with 2 basis weight matrices, and layer sizes of $500$ and $75$ for the graph convolution (with $\mathrm{ReLU}$) and dense layer (no activation function), respectively. We evaluate our model on the held out test sets using an exponential moving average of the learned model parameters with a decay factor set to $0.995$.

**MovieLens 100K** For this task, we compare against matrix completion baselines that make use of side information in the form of user/item features[3]. Side information is present both for users (e.g. age, gender, and occupation) and movies (genres). Following Rao et al. [21], we map the additional information onto feature vectors for users and movies, and compare the performance of our model with (GC-MC+Feat) and without (GC-MC) the inclusion of these features. Note that GMC [9], GRALS [21] and sRGCNN [19] represent user/item features via a k-nearest-neighbor graph. GC-MC+Feat uses 10 hidden units for the dense side information layer (with $\mathrm{ReLU}$ activation) as described in Eq. 10. We train both models for 1,000 full-batch epochs. We report RMSE scores averaged over 5 runs with random initializations[4]. Results are summarized in Table 2.

**MovieLens 1M and 10M** We compare against current state-of-the-art collaborative filtering algorithms, such as AutoRec [23], LLorma [14], and CF-NADE [28]. Results are reported as averages over the same five 90/10 training/test set splits as in [28] and summarized in Table 3. Model choices are validated on an internal 95/5 split of the training set. As ML-10M has twice the number of rating classes, we use twice the number of basis function matrices in the decoder. We train for 3,500 full-batch epochs, and 18,000 mini-batch iterations (20 epochs with batch size 10,000) on the ML-1M and ML-10M dataset, respectively.

**Discussion** On the ML-100K task with side information, our model outperforms related methods by a significant margin. Remarkably, this is even the case without the use of side information. Most related to our method is sRGCNN by Monti et al. [19] that uses graph convolutions on the nearest-neighbor graphs of users and items, and learns representations in an iterative manner using recurrent neural networks. Our results demonstrate that a direct estimation of the rating matrix from learned user/item representations using a simple decoder model can be more effective, while being computationally more efficient.

Our results on ML-1M and ML-10M demonstrate that it is possible to scale our method to larger datasets, putting it into the vicinity of recent state-of-the-art user- or item-based methods in terms

---

[2]https://grouplens.org/datasets/movielens/

[3]For ML-100K, we report performance on the canonical u1.base/u1.test train/test split. Hyperparameters are optimized on a 80/20 train/validation split of the original training set.

[4]Standard error less than 0.001.

| Model | ML-100K + Feat |
|---|---|
| MC [2] | 0.973 |
| IMC [8, 27] | 1.653 |
| GMC [9] | 0.996 |
| GRALS [21] | 0.945 |
| sRGCNN [19] | 0.929 |
| GC-MC (Ours) | 0.905 |
| GC-MC+Feat | **0.901** |

Table 2: RMSE scores for the Movie-Lens 100K task with side information on a canonical 80/20 training/test set split. Side information is either presented as a nearest-neighbor graph in user/item feature space or as raw feature vectors. Baseline numbers are taken from [19].

| Model | ML-1M | ML-10M |
|---|---|---|
| PMF [17] | 0.883 | – |
| I-RBM [22] | 0.854 | 0.825 |
| BiasMF [13] | 0.845 | 0.803 |
| NNMF [6] | 0.843 | – |
| LLORMA-Local [14] | 0.833 | 0.782 |
| I-AUTOREC [23] | 0.831 | 0.782 |
| CF-NADE [28] | **0.829** | **0.771** |
| GC-MC (Ours) | 0.832 | 0.777 |

Table 3: Comparison of average test RMSE scores on five 90/10 training/test set splits (as in [28]) without the use of side information. Baseline scores are taken from [28]. For CF-NADE, we report the best-performing model variant.

of predictive performance. At this point, it is important to note that several techniques introduced in CF-NADE [28], such as layer-specific learning rates, a special ordinal loss function, and the auto-regressive modeling of ratings, can be seen as orthogonal to our approach and can be used in conjunction with our framework.

## 5 Conclusions

In this paper, we have introduced graph convolutional matrix completion (GC-MC), a novel framework for end-to-end learning on bipartite user-item interaction graphs. Our model takes the form of a graph auto-encoder, comprised of a graph convolutional encoder model and a bilinear decoder, and can be seen as a generalisation of previous user- or item-based auto-encoder models to the case where users and items are modeled in a joint representation space.

As our encoder model learns both user and item representations, it is straightforward to include side information for both types of nodes. In this setting, our proposed model outperforms recent related methods that make use of such side information by a large margin, as demonstrated on the MovieLens 100K dataset. We further show that our model can be trained on datasets of larger scale without resorting to approximations or subsampling schemes. In this setting, our graph auto-encoder achieves results comparable to recent state-of-the-art user- or item-based methods.

Our model can be seen as a first step towards modeling recommender systems where the interaction data is integrated into other structured modalities, such as a social network or a knowledge graph. As a next step, it would be interesting to investigate how the differentiable message passing scheme of our encoder model can be extended to such structured data environments. We expect that further approximations, e.g. subsampling of local graph neighborhoods, will be necessary in order to keep requirements in terms of computation and memory in a feasible range.

## References

[1] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013.

[2] Emmanuel Candes and Benjamin Recht. Exact matrix completion via convex optimization. *Communications of the ACM*, 55(6):111–119, 2012.

[3] Hanjun Dai, Bo Dai, and Le Song. Discriminative embeddings of latent variable models for structured data. In *International Conference on Machine Learning (ICML)*, 2016.

[4] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems*, pages 3837–3845, 2016.

[5] David K. Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P. Adams. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in neural information processing systems (NIPS)*, pages 2224–2232, 2015.

[6] Gintare Karolina Dziugaite and Daniel M Roy. Neural network matrix factorization. *arXiv preprint arXiv:1511.06443*, 2015.

[7] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural message passing for quantum chemistry.

[8] Prateek Jain and Inderjit S Dhillon. Provable inductive matrix completion. *arXiv preprint arXiv:1306.0626*, 2013.

[9] Vassilis Kalofolias, Xavier Bresson, Michael Bronstein, and Pierre Vandergheynst. Matrix completion on graphs. *arXiv preprint arXiv:1408.1717*, 2014.

[10] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[11] Thomas N. Kipf and Max Welling. Variational graph auto-encoders. *NIPS Bayesian Deep Learning Workshop*, 2016.

[12] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *ICLR*, 2017.

[13] Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, August 2009.

[14] Joonseok Lee, Seungyeon Kim, Guy Lebanon, and Yoram Singer. Local low-rank matrix approximation. *ICML (2)*, 28:82–90, 2013.

[15] Xin Li and Hsinchun Chen. Recommendation as link prediction in bipartite graphs: A graph kernel-based machine learning approach. *Decision Support Systems*, 54(2):880 – 890, 2013.

[16] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated graph sequence neural networks. *ICLR*, 2016.

[17] Andriy Mnih and Ruslan R Salakhutdinov. Probabilistic matrix factorization. In *Advances in neural information processing systems*, pages 1257–1264, 2008.

[18] Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodolà, Jan Svoboda, and Michael M Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. *arXiv preprint arXiv:1611.08402*, 2016.

[19] Federico Monti, Michael M. Bronstein, and Xavier Bresson. Geometric matrix completion with recurrent multi-graph neural networks. *preprint arXiv:1704.06803*, 2017.

[20] Mathias Niepert, Mohamed Ahmed, and Konstantin Kutzkov. Learning convolutional neural networks for graphs. In *Proceedings of the 33rd annual international conference on machine learning. ACM*, 2016.

[21] Nikhil Rao, Hsiang-Fu Yu, Pradeep K Ravikumar, and Inderjit S Dhillon. Collaborative filtering with graph information: Consistency and scalable methods. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2107–2115. Curran Associates, Inc., 2015.

[22] Ruslan Salakhutdinov, Andriy Mnih, and Geoffrey Hinton. Restricted boltzmann machines for collaborative filtering. In *Proceedings of the 24th international conference on Machine learning*, pages 791–798. ACM, 2007.

[23] 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*, pages 111–112. ACM, 2015.

[24] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014.

[25] 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, pages 11–16, New York, NY, USA, 2016. ACM.

[26] Fei Tian, Bin Gao, Qing Cui, Enhong Chen, and Tie-Yan Liu. Learning deep representations for graph clustering. In *AAAI*, pages 1293–1299, 2014.

[27] Miao Xu, Rong Jin, and Zhi-Hua Zhou. Speedup matrix completion with side information: Application to multi-label learning. In *Advances in Neural Information Processing Systems*, pages 2301–2309, 2013.

[28] Yin Zheng, Bangsheng Tang, Wenkui Ding, and Hanning Zhou. A neural autoregressive approach to collaborative filtering. In *Proceedings of the 33nd International Conference on Machine Learning*, pages 764–773, 2016.