CONSTRAINING EMBEDDING LEARNING WITH SELF-MATRIX FACTORIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

We focus on the problem of learning object representations from solely association data, that is observed associations between objects of two different types, e.g. movies rated by users. We aim to obtain embeddings encoding object attributes that were not part of the learning process, e.g. movie genres. It has been shown that meaningful representations can be obtained by constraining the learning with manually curated object similarities. We propose Self-Matrix Factorization (SMF), a method that learns object representations and object similarities from observed associations, with the latter constraining the learned representations.^{AA} In our extensive evaluation across three real-world datasets, we compared SMF with SLIM, HCCF and NMF obtaining better performance at predicting missing associations as measured by RMSE and precision at top-K. We also show that SMF outperforms the competitors at encoding object attributes as measured by the embedding distances between objects divided into attribute-driven groups.

023 024 025

026 027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

This paper focuses on the problem of learning object representations from observed associations between objects of two different types. We assume that no data is available other than the associations between the objects, and our aim is to obtain representations that reflect object attributes and properties that, in general, are unknown.

The observed associations between two groups of n and m objects respectively can be represented by a data matrix $X \in \mathbb{R}^{n \times m}$, where the association between object i and object j is stored in $X_{i,j}$. We consider the case in which this value can be either a binary or a positive number, and a value of zero indicates no known association between objects i and j.

Learning meaningful object representations has been shown to be relevant for several tasks including 037 recovering missing associations and clustering objects into meaningful groups, and several methods 038 have been proposed. Matrix Factorization (MF) methods, for instance, assume that the association matrix is low rank, allowing X to be decomposed into lower-dimensional matrix factors, containing the representations of the objects. Techniques like principal component analysis (PCA) (Hotelling, 040 1933), singular value decomposition (SVD)(Eckart & Young, 1936) and non-negative matrix fac-041 torization (NMF) (Lee & Seung, 1999) have been successfully applied to tasks such as recovering 042 associations (Sarwar et al., 2002; Vozalis & Margaritis, 2007; Luo et al., 2014) and clustering objects 043 (Yang & Seoighe, 2016; Yeung & Ruzzo, 2001). 044

- Deep Learning is another widely used technique for learning object embeddings, particularly through graph neural networks (GNNs). An association data matrix can be thought of as a bipartite graph where the nodes represent the objects in the data and the links represent the associations between them. GNNs leverage this network structure to extract insights from the encoded graphs. While deep learning methods have been shown to be particularly effective at incorporating prior known object properties (Wu et al., 2022), a number of techniques have also been developed that can use solely association data such as, for example, LightGCN (He et al., 2020), SEAL (Zhang et al., 2021) and HCCF (Xia et al., 2022).
- Recently, learning strategies that rely on manually curated similarities between objects have been proposed to constrain embedding learning somehow. For example, Neo-GNNs (Yun et al., 2021)

and BUDDY (Chamberlain et al., 2023) are GNN methods relying on higher-order interactions in the graph. These interactions function as additional node similarity features and were used to enhance link-prediction performance. However, selecting such similarities is not trivial.

057 In this paper, we argue that object similarities can be learned directly from the data matrix. We rely 058 on the fact that^{AA} the objects lie on multiple linear low-dimensional manifolds embedded in a high-059 dimensional space (Elhamifar & Vidal, 2013). Our matrix decomposition approach, Self-Matrix 060 Factorization (SMF), learns distributed representations while constraining them using learned object 061 similarities. These similarities depend on the manifold structures implicit in the association matrix 062 X and are learned together with the embeddings. In other words, the object similarities, determined 063 by their positions in the manifolds, naturally constrain the object embeddings during the learning. Our method is the first to explore this idea in a matrix factorization model^{AA}. In our extensive 064 evaluation across three distinct benchmark datasets, we show that SMF consistently outperforms the 065 competitors at encoding object attributes as measured by the embedding distances between objects 066 divided into attribute-driven groups. We also performed experiments to recover missing values on 067 the different association matrices and show that SMF obtains comparable or better predictions than 068 its competitors. 069

- 070
- 071
- 072 073

2 RELATED WORKS

074 075

076 MF and GNN techniques encompass numerous methods for learning object representations from 077 association data (Koren et al., 2021; Wu et al., 2022). MF techniques decompose the association 078 matrix X into two or more matrix factors, where the object representations are encoded as rows 079 or columns of these matrix factors, mapping objects to a shared latent space of lower dimension-080 ality (Aggarwal et al., 2016). Several methods for link prediction have been proposed, including 081 SVD (Koren et al., 2009), SVD++ (Koren, 2008) and probabilistic matrix factorization (Yang et al., 082 2014). NMF (Lee & Seung, 1999) and its variations have been used across fields ranging from 083 medicine to engineering (Hamamoto et al., 2022; Sturluson et al., 2021). Graph-regularized NMF (Cai et al., 2010), symmetric NMF (Luo et al., 2021) and robust NMF (Peng et al., 2021) have 084 been successfully used for object clustering and community detection. Additionally, NMF with 11, 085 12 or elastic net regularization has been applied successfully across diverse applications, including 086 precision medicine (Hamamoto et al., 2022), gene-expression analysis (Sweeney et al., 2023) and 087 recommender systems (Rendle et al., 2020), showing state-of-the-art performance.^{AA} 088

GNNs have gained popularity for their strong capabilities in graph representation learning. These 089 methods can effectively learn node representations that are well-suited for link prediction tasks 090 (Zhang et al., 2021). One advantage of GNNs is their ability to incorporate external object fea-091 tures, which can significantly enhance prediction performance (Wu et al., 2022). Some approaches, 092 like graph-regularized NMF (Cai et al., 2010), BUDDY (Chamberlain et al., 2023), and Neo-GNNs 093 (Yun et al., 2021), leverage similarity measures to improve object clustering and link prediction per-094 formance. HCCF, a specialized GNN technique, learns hyper-edges between objects, enabling it to 095 simultaneously learn embeddings and refine object similarities for improved representation learn-096 ing.^{AA}

Manually curated similarities have proven useful for embedding learning, stemming from the fact 098 that these similarities can themselves be used in recommender systems (Aggarwal et al., 2016). Sparse Linear Models (SLIM) (Ning & Karypis, 2011) are state-of-the-art recommender systems 100 (Ferrari Dacrema et al., 2019) that rely on learning object similarities rather than embeddings. SLIM 101 learns coefficients such that each object can be represented as a linear combination of other objects. 102 This means that a new link between objects i and j is predicted only if objects similar to i were 103 originally linked with j. The coefficients used to reconstruct objects depend on the linear manifolds 104 present in the data matrix X. In this way, new links are recommended to an object based on the 105 links other objects belonging to the same linear manifold have. Although these similarities have demonstrated predictive power, they have not yet been used to inform embedding learning. In this 106 work, we address this gap by proposing a framework that jointly learns object embeddings and object 107 similarities, where the latter constrains the embedding space, resulting in richer representations.^{AA}

3 Self-Matrix Factorization

108

110

111

112

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

SMF learns two non-negative matrices $W \in \mathbb{R}^{n \times k}$ and $H \in \mathbb{R}^{k \times m}$, with $k \ll (m \times n)$. Each matrix contains distinct low dimensional object embeddings, such that their product approximates the low-rank interaction data matrix $X \in \mathbb{R}^{n \times m}$:

$$X \simeq WH. \tag{1}$$

While this model is not new, its novelty^{AA} resides in the learning of the embeddings in W to encode linear manifold information implicitly contained in the association data itself. Relying on the above mentioned assumption that objects lie on multiple linear low-dimensional manifolds embedded in high-dimensional space (Elhamifar & Vidal, 2013), let us consider the situation depicted in Figure 1a in which we have points in the 3-D space that are approximately localized onto 3 distinct linear manifolds. Rows of X are represented as squares, triangles and circles, with triangles and squares lying on one-dimensional sub-space (red and brown lines) and circles lying on a two-dimensional sub-space (green plane). Let us focus on the three blue points of which i and p lie on the plane and qon the red line.^{AA}. We assume that objects that belong to the same subspace, are more similar to each other than objects that reside in different subspaces. We would like these similarities to constrain the learning of the embeddings – that is, we would like the embedding for two objects that belong to the same subspace, to be more similar to each other than the embeddings of objects that reside in different subspaces. Thus, in the embedding space (2-dimensional, in Figure 1b), object i should be closer to object p than to object q, mimicking their behavior in the high-dimensional space. Figure 1b demonstrates the expected behavior of SMF-learned object embeddings. Points that belong to the same linear manifold in the high-dimensional space are projected into a lower-dimensional space, where they closely approximate one another.^{AA}



Figure 1: SMF explicit constraint. In this example, the association matrix X contains only 3 columns. X is decomposed into the product WH, where W have 2 columns. (a) Positions of X rows in the 3-dimensional space. Points represented as dots, triangles and squares belong to different subspaces. (b) Positions of the 2-dimensional rows of W in the space, SMF uses the similarities established by the linear manifolds to constrain W such that a pair of object embeddings are likely to have a high dot product if they belong to the same linear manifold in the 3-dimensional space.

159 160 161

153

154

155

156

157

158

We propose the following loss function for learning a model with these properties:

165

166

167 168

183

185

169 where \circ represents the element-wise product and $\|\cdot\|_F^2$ indicates the Frobenius norm. The first term 170 of Equation 2 is the Euclidean distance between the non-negative matrix X and the product of two 171 non-negative matrices W and H. Minimizing this distance results in projecting high-dimensional 172 data into a low-dimensional representation. The non-negativity constraint of matrices W and H is 173 a crucial factor for the interpretability of the representations (Lee & Seung, 1999). This constraint 174 naturally encourages any pair of vectors $W_{i,:}$ and $H_{i,:i}$ to exhibit a significant overlap in their highvalued components if the objects i and j share an observed association in X. Conversely, objects 175 that do not meet this condition can be understood as lacking common attributes, making them less 176 likely to interact. 177

 $+\lambda_1 \|W\|_1 + \lambda_1 \|H\|_1 + \frac{\lambda_2}{2} \|W\|_F^2 + \frac{\lambda_2}{2} \|H\|_F^2$

(2)

subject to $W, H \ge 0$.

 $\min_{W,H} \mathcal{L}_{SMF}(W,H) = \frac{1}{2} \|X - WH\|_F^2 + \frac{\lambda_{se}}{4} \|X - [T \circ (WW')]X\|_F^2$

178While parts of Equation 2 resemble the loss function of NMF, its second term introduces a funda-179mental novelty. AA It is designed to preserve the linear manifold information implicit in the matrix180X. Matrix T is populated with ones, except for the diagonal where elements are set to zero. By181reducing the distance between the matrix $[T \circ (WW')]X$ and the original association matrix, we are182reconstructing each row of X using other rows of X:

$$X_{i,:} = \sum_{j} T_{i,j} [WW']_{i,j} X_{j,j}$$

Note that, since we are learning W, the presence of the matrix T is necessary to avoid the trivial so-186 lution in which WW' becomes the identity matrix. The last 4 terms apply elastic-net regularization 187 to the matrices W and H to promote sparsity and mitigate overfitting. Therefore, since minimizing 188 the loss function in Equation 2 is attempting to reconstruct a row using only a few other rows, the 189 learning will favour reconstructing each row using only rows representing the objects in the same 190 subspace. For instance, in Figure 1a one can reconstruct each point in the plane by using only points 191 in the plane, without the need to use points from different subspaces. Let us also note that matrix 192 $[T \circ (WW')]$ contains the coefficients for the reconstruction of the rows of X. Therefore, matrix 193 WW' is attempting to encode the inherent similarities between the objects established by the linear 194 manifolds. During the learning of W, this amount to promote higher values for the dot product 195 between the lower-dimensional representations of objects within the same subspace than for the dot 196 product between representations of objects in different subspaces.

197 By minimizing the loss function in equation 2, we approximate each interaction $X_{i,j}$ as $(W_{i,:})$. 198 $H_{i,j}$ (first term) as well as $\sum_{s} T_{i,s}(W_{i,i} \cdot W'_{s,i}) X_{s,j}$ (second term). The first term enforces shared 199 latent features between the rows and column objects, while the second term incorporates an explicit 200 constraint for all the embeddings of the objects in the row of X. This second constraint is directly 201 related to the similarity between object embeddings in W, so that the dot product between any pair $W_{i,i}$ and $W_{p,i}$ is informed by the linear manifolds in which objects i and p lies. Notably, SMF does 202 not require prior knowledge of these manifolds; instead, it simultaneously learns the embeddings 203 and the manifold structure, making it the first method to integrate these two processes.^{AA} 204

Similarly to NMF (Lee & Seung, 2000), we derived a multiplicative update rule to minimize the function in Equation 2:^{AA}

207 208 209

$$W_{i,j} \leftarrow W_{i,j} \times \frac{[XH' + \lambda_{se}((XX') \circ T)W]_{i,j}}{[WHH' + \lambda_{se}(((T \circ (WW'))XX') \circ T)W + \lambda_2W + \lambda_1 \operatorname{sgn}(W)]_{i,j}}$$
(3)

210 211

212 213

$$H_{i,j} \leftarrow H_{i,j} \times \frac{[W'X]_{i,j}}{[W'WH + \lambda_2 H + \lambda_1 \operatorname{sgn}(H)]_{i,j}}$$
(4)

· · - - - • •

where $sgn(\cdot)$ is the sign function.^{AA} W and H were initialized with non-negative values to ensure that the proposed multiplicative update rules process results in non-negative embeddings after each iteration.

224 225

226 227 228

229 230

241

242

243 244

245 246

Finally, since in our data matrix X the zeros denote our lack of knowledge about a possible association, it is often convenient to modulate the importance of the zeros during the learning. This has been done by other authors before us (e.g. (Galeano et al., 2020; Blondel et al., 2008)) and is normally achieved by weighting the contribution of the zero values by a factor $\alpha << 1$ in the loss function. In this case, our loss function becomes:

$$\min_{W,H} \mathcal{L}_{WSMF}(W,H) = \frac{1}{2} \|P \circ (X - WH)\|_F^2 + \frac{\lambda_{se}}{4} \|P \circ (X - [T \circ (WW')]X)\|_F^2
+ \lambda_1 \|W\|_1 + \lambda_1 \|H\|_1 + \frac{\lambda_2}{2} \|W\|_F^2 + \frac{\lambda_2}{2} \|H\|_F^2$$
(5)
subject to $W, H > 0$.

where the matrix $P \in \mathbb{R}^{n \times m}$ is defined as $P_{i,j} = 1$ when $X_{i,j} > 0$ and $P_{i,j} = \alpha$ otherwise.

Equations 3 and 4 are then modified accordingly to optimize Equation 5:

$$W_{i,j} \leftarrow W_{i,j} \times \frac{[XH' + \lambda_{se}((XX') \circ T)W]_{i,j}}{[((P \circ P) \circ (WH))H' + \lambda_{se}((((P \circ P) \circ ((T \circ (WW'))X))X') \circ T)W + \lambda_2W + \lambda_1 \operatorname{sgn}(W)]_{i,j}}$$
(6)

$$H_{i,j} \leftarrow H_{i,j} \times \frac{[W'X]_{i,j}}{[W'((P \circ P) \circ (WH)) + \lambda_2 H + \lambda_1 \operatorname{sgn}(H)]_{i,j}}$$
(7)

From equations 5, 6, and 7, and assuming that $k \ll n, m$, it follows that the time complexity of each iteration of SMF optimization algorithm is $O(n^2 \cdot m)$. We discuss the complexity of SMF in the Appendix A.5, together with details about computational time and number of iterations.^{AA}

4 EXPERIMENTAL RESULTS

To test SMF, we compared its performance against three different models, namely NMF with elastic-247 net regularization (Pauca et al., 2006), Hypergraph Contrastive Collaborative Filtering (HCCF) (Xia 248 et al., 2022) and Sparse Linear models (SLIM) (Ning & Karypis, 2011). We compared the perfor-249 mance of the models at predicting associations and the quality of the embeddings through clustering 250 analysis. NMF with elastic-net regularization and HCCF were chosen as representative of matrix de-251 composition techniques and GNN-based methods, respectively. NMF has been shown to be effective 252 at recovering missing associations and encoding object attributes. HCCF is a state-of-the-art GNN 253 model for link prediction. Like SMF, HCCF also learns similarities between objects of the same 254 type, avoiding dependence on manually curated similarities. However, it enhances embeddings by 255 learning hyper-edges (connections involving more than two objects) that contribute to the embedding construction. For the task of predicting associations, we also compared SMF with SLIM (Ning 256 & Karypis, 2011), a state-of-the-art approach for predicting missing associations (Ferrari Dacrema 257 et al., 2019) that has been shown to be competitive with deep learning models. SLIM does not 258 learn object embeddings and for this reason, we could not perform any clustering analysis. Size 259 embeddings for the models are given in the last three columns of Table 1, other algorithm details 260 and their implementation are available in the Appendix, while the code we used in our experiments 261 is available in the Supplementary Material. 262

We run our experiments on three datasets namely, Movielens, Drug-SE and ModCloth. These datasets were chosen among those that have been used for the task of predicting associations because they also included attribute information of the objects. Overall details are given in Table 1 and descriptions of the three datasets are given below.

Movielens: This dataset describes ratings ranging from 1 to 5 that users gave to movies. It is
 a smaller version of Group lens that is made available for educational and development purposes
 (Harper & Konstan, 2015). The one used in this work includes object attributes for both users and
 movies. Each movie is associated with its respective genre (18 genres in total), and each user is

associated with its gender (2 genders in total). The data matrix contains 943 users and 1682 movies. It contains 100000 non-zero elements representing the known ratings, resulting in an association data matrix with a density equal to 6.3%.

Frequencies of Drug Side Effects (Drug-SE): Galeano et al. obtained a data matrix containing the
frequencies in which certain drugs produce specific side effects(Galeano et al., 2020) by filtering the
frequencies obtained from the Side Effect Resource Database (SIDER) (Kuhn et al., 2016). Integers
between 1 and 5 represent side-effect frequency terms for 'very rare', 'rare', 'infrequent', 'frequent', and 'very frequent' respectively.

Drugs can be grouped by clinical activity using their main Anatomical, Therapeutic and Chemical (ATC) class levels. ATC is a hierarchical organization of terms maintained by the World Health Organization. A term at a lower level indicates a more specific descriptor of clinical activity. Each drug in the matrix *R* is associated with its respective ATC-category term in three different levels. The drugs in this dataset belong to all 14 groups at the more general Anatomical level, 70 out of 94 groups at the intermediate Therapeutic level and 147 out of 262 groups at the more specific Chemical level. The data matrix contains 759 drugs and 994 side effects. It contains 37441 non-zero elements representing known frequencies, resulting in an association matrix with a density of 5%.

ModCloth: This dataset contains ratings that users gave to different clothing items (Misra et al., 2018). Originally, a rating in this dataset could be 2, 4, 6, 8, or 10. However, we divided all the values by 2, resulting in ratings ranging from 1 to 5. Due to the low density of known ratings in the association matrix, we eliminated those users and clothes with less than 10 associations from the data matrix. The resulting matrix had higher density but still had objects with less than 10 associations, resulting in a 0.05% density. Each clothing item belongs to 1 of 66 different categories.

293 294 295

Table 1: Datasets and embedding sizes

Datasets	rows	columns	density	NMF	HCCF	SMF
Movielens Drug-SE ModCloth	$943 \\759 \\5419$	$1682 \\ 994 \\ 32089$	${6.3\%} \atop {5\%} 0.05\%$	$10 \\ 10 \\ 30$	$32 \\ 32 \\ 32 \\ 32$	$10 \\ 10 \\ 30$

- 300 301
- 302 303

304

4.1 PERFORMANCE EVALUATION AT PREDICTING ASSOCIATION

SMF achieves scores closer to the real values. We evaluated SMF by assessing the model's per formance at recovering the different levels of associations. In our experiments, we set 10% of the
 known associations to zero and then we compared the performance of the different models at recov ering them.

We used Root Mean Square Error (RMSE) to assess the reconstruction of the association matrix. The outcomes of the evaluations for NMF, SLIM and SMF are reported in Tables 2 as the mean RMSE across 30 runs of each model, along with the corresponding variance. HCCF is not included because it cannot predict the ratings nor the frequencies, only the presence of a link.

The RMSE is a comparison between known associations and the scores predicted by the models, where lower RMSE values indicate that the predicted scores are closer to the actual associations. We can see that SMF scores remain consistently closer to the real values than those produced by NMF and SLIM. SMF achieves a 15% lower RMSE than NMF in the sparser ModCloth dataset and at least 65% lower RMSE than SLIM for all datasets.

318 SMF achieves better performance at top-K predictions. It is important to measure a system's ability to predict the existence of associations between objects, independently of their specific values. For our datasets, this amounts to predicting which movies is a user more likely to watch, which side effects is a drug likely to cause, and which clothes is a user more likely to rent. To measure this, we generated three binary datasets. These new datasets were built by replacing the non-zero elements on all the datasets. In practice, we are often interested in predicting only a small number of associations with high accuracy. This is commonly referred to as the top-K recommendation task

22/

332

352

353

354

355

356

357

359

360

361

362

364

365

367

371

524				
325		Table 2: Root	t Mean Square Error	
326	MODELS	Movielens	Drug-SE	ModCloth
327			8	
328	NMF	$0.9777 \pm 3e{-5}$	$0.6558 \pm 1\mathrm{e}{-4}$	$1.5759 \pm 3e{-4}$
329	SLIM	$2.9480\pm3\mathrm{e}{-6}$	$1.8622\pm5\mathrm{e}{-6}$	$3.7790 \pm 4\mathrm{e}{-7}$
330	SMF	$0.9352\pm1\mathrm{e}{-5}$	$0.6455\pm5\mathrm{e}{-5}$	$1.3258 \pm 1\mathrm{e}{-5}$
331				

333 (Cremonesi et al., 2010) where a system's performance is measured using precision at top-K. For 334 this purpose, we ranked the scores to retrieve the K higher elements. These were predicted as new associations between objects and compared with the test set to obtain the precision at top-K, which 335 is the ratio of known associations within the predicted associations. To have fair measurements, 336 the true positives for the analysis are the ones on the test set, and all the unknown elements on the 337 original datasets are considered as true negatives (Krichene & Rendle, 2020). The outcomes of the 338 evaluations are reported in Figure 2, where SMF outperforms the competitors in almost every set-339 ting. ModCloth datasets results are not shown due to the low association density, all the models only 340 manage to predict a few true positives at the top 1000. 341

We can see that SMF predictions achieve the top precision in 7 out of 10 settings. HCCF achieves a 342 better precision at the top 10 and top 20 for the Drug-SE dataset, indicating that it contains more true 343 associations in the top predictions. However, the precision drops as K increments, resulting in the 344 worst precision of all the models for the tops 50, 100 and 150 in the same Drug-SE dataset. SLIM, 345 a state-of-the-art predictor for top-K recommendations (Ferrari Dacrema et al., 2019), achieves the 346 best precision for the top 150 in the Drug-SE dataset. 347

The Area under the receiver operating characteristic curve (AUROC) and the Area under the 348 precision-recall curve (AUPRC) are also useful metrics to evaluate the overall distribution of the 349 true positives. The AUROC, AUPRC and correlation outcomes for all datasets are shown in the 350 Appendix. 351



366 Figure 2: Precision at top-K: Bar plot of the precision of NMF, SLIM, HCCF and SMF for different values of K while predicting missing links in the interaction data. The error bars indicate the 368 variance of the precision for 30 different runs of the models. (a) Precision for the Movielens dataset 369 while predicting links between users and movies, the negative to positive ratio in the test set is ap-370 proximately 1600. (b) Precision for the Drug-SE dataset while predicting links between side effects and drugs, the negative to positive ratio in the test set is approximately 200. 372

373 SMF sensibility to hyperparameters settings: The SMF loss function proposed in Eq. 5 contains 374 five hyperparameters. SMF demonstrates stable performance across a wide range of hyperparameter 375 values, indicating that its practical application does not require extensive hyperparameter tuning. The parameter λ_{se} controls the importance of the self-expressive term and we set it to 1 in all exper-376 iments in this paper. Figure 4 in the Appendix A.4 explores the effect other hyperparameters have 377 on embedding learning by assessing the RMSE and AUPRC on the validation set using the Movie378 Lens dataset. SMF is robust to the choice of the object embedding dimension k, achieving good 379 performance even for low values of k. As it was also shown by other authors (Galeano et al., 2020), 380 α value depends on the task. α should be set to a low value (closer to zero) when the objective is 381 to accurately retrieve the numerical values of the associations, as in tasks focused on minimizing 382 RMSE. Conversely, α should be set to a high value (closer to 1) when correctly identifying the associations themselves is more critical, as in tasks that optimize AUPRC. Additionally, this experiment 383 shows that SMF is resilient to different values of the λ_1 and λ_2 regularization weights. Finally, per-384 formance remains consistent across the explored search space, with the only significant variations 385 arising predictably from changes in α .^{AA} 386

387 388

389

4.2 EMBEDDING EVALUATION THROUGH CLUSTER ANALYSIS

390 To demonstrate that the SMF-derived embeddings offer a more meaningful encoding of previously unseen latent object attributes, we analyzed these low-dimensional representations and compared 391 them with the ones learned by NMF and HCCF. Our aim is to assess their capacity for encapsulating 392 inherent data characteristics, which were not part of the training process but may play crucial roles 393 in establishing connections between objects. Our embedding analysis was conducted on two levels: 394 first, to verify whether SMF effectively clusters objects into meaningful groups within the low-395 dimensional space; and second, to assess whether SMF achieves superior class separation of objects 396 compared to NMF and HCCF. 397

We took advantage of this information and grouped the embeddings in W into disjoint sets based on the classes to which their corresponding objects belong. Subsequently, we calculate a similarity matrix, $W_{sim} \in \mathbb{R}^{n \times n}$, containing the cosine similarity between all the embeddings. Finally, we employ a two-sample *t*-test to assess whether the similarities between objects within the same class (intra-similarities) differ significantly from the similarities between objects in different classes (intersimilarities). An illustration of this procedure is provided in Figure 3a. This process was repeated 30 times across different runs of the models.

For the Movielens dataset, we organized users based on their gender, and movies by their respective 405 genres. In our genre analysis, we refined the dataset to include only movies with a single genre, 406 enabling the classification of movies into disjoint categories. In the case of the Drug-SE dataset, we 407 categorized drugs according to their various levels within the ATC hierarchy¹. The lower levels of 408 the ATC hierarchy provide more specific terms for drug classification. For this study, we compared 409 the similarity of drug embeddings across three levels of the hierarchy: anatomical, therapeutical, 410 and chemical. Finally, the clothing items of the ModCloth dataset can be separated into different 411 groups depending on which type of clothes they are (dresses, jeans, blazers, etc). 412

SMF consistently clusters objects in the low-dimensional space. In our analysis of the distributions of intra- and inter-similarities for the clothing types included in the ModCloth dataset, all NMF, HCCF and SMF achieved significant separation in 100% of the runs. However, for the Movielens dataset, clustering movies by gender in the embedding space proved challenging for HCCF. NMF, HCCF and SMF attain significant separation 100%, 67% and 93% of the runs, respectively.

When considering the various levels of the ATC drug classification hierarchy, SMF and HCCF achieved significant distribution separation 100% of the runs for every level. On the other hand, NMF struggled to maintain consistent separation, achieving statistical significance in only 3%, 13%, and 13% of runs for the 1^{st} , 2^{nd} , and 3^{rd} levels, respectively. The consistent achievement of statistical significance in these experiments indicates effective clustering of objects in the low-dimensional space. This provides compelling evidence that the SMF-learned embeddings reliably encode meaningful information about the fundamental attributes of the objects.

SMF achieves superior class separation. To assess the efficacy of each method in achieving class separation, we employ the Z-score difference between the means of the intra-class and inter-class similarity distributions:

$$Z = \frac{\mu_{in} - \mu_{out}}{\sqrt{\frac{\sigma_{in}^2}{n_{in}} + \frac{\sigma_{out}^2}{n_{out}}}},$$

¹ATC categories were obtained from the ATC codes WHO 2018 release.

where μ_{in} is the average embedding similarity of object pairs in the same group. μ_{out} is the average embedding similarity of object pairs in different groups. σ_{in} and σ_{out} are the corresponding standard deviations, and n_{in} and n_{out} are the corresponding number of object-pairs.

We can interpret the z-score as a normalized distance that measures how different two distributions are by adjusting the difference between the means according to their standard deviation.

Our results for this experiment are summarized in Figures 3b, 3.c and 3.d. We can see that SMFlearned embeddings effectively group objects into more meaningful clusters than those learned by both NMF and HCCF, across all the datasets and diverse groups. Notably, in Figure 3d, we observe that the separation between groups in the ATC levels increases as we delve from the first to the second level of the hierarchy. This reflects the fact that the drug clinical activity becomes more similar as we move to more specific levels.

444

5 CONCLUSION AND DISCUSSION

445 446

447 Many machine learning approaches rely on learning distributed representations able to reflect relevant object attributes. A common strategy to enrich these embeddings is by directly constraining 448 them to follow similarities extracted from side information (Aggarwal et al., 2016). Similarly, one 449 can directly rely on the similarities in the association matrix to guide the embedding learning to 450 better uncover patterns in the data. In this work, we introduced Self-Matrix Decomposition (SMF), 451 a constrained matrix decomposition approach that learns low-dimensional representations by cons-452 training them to rely on object similarities. These similarities depend on linear manifolds implicit in 453 the association data and are learned with the representations. 454

SMF can decompose a low-rank matrix while preserving its inherent similarities in WW', leveraging the relationships between rows of X that are revealed thanks to the second term in Equation 2. This *Self-Expressive* term learns a coefficient matrix that allows X to be reconstructed by itself, similarly to a Self-Expressive model (Elhamifar & Vidal, 2013). The loss function in Equation 5 facilitates the embeddings to learn better representations and capture the latent attributes, as they allow the embeddings to glean information from objects residing within the same subspace (as depicted in Figure 1).

We conducted experiments to assess whether a set of known object properties could be effectively 462 encoded within the object embeddings. Prior research has also delved into similar investigations, 463 revealing, for instance, that various NMF variants can learn embeddings encoding movie genres 464 (Gomez-Uribe & Hunt, 2015) and drug ATC categories (Galeano et al., 2020). We conducted an 465 in-depth analysis of the similarities among embeddings generated by NMF, HCCF and SMF across 466 multiple runs and diverse groupings. This analysis aimed to ascertain whether objects belonging to 467 the same group consistently clustered together in the low-dimensional space. SMF offers signifi-468 cantly higher stability in learning well-separated embeddings compared to NMF and HCCF. This 469 is evident from the fact that in multiple runs, SMF achieves statistical significance approximately 470 99% of the time, whereas NMF and HCCF accomplish this feat in only 41% and 87% of the runs, respectively. Furthermore, the experimental results demonstrate that SMF consistently achieves su-471 perior class separation in all conducted experiments (depicted in Figure 3). Consequently, SMF can 472 be used to cluster objects in meaningful groups, and an analysis of these groups may help reveal 473 hidden object attributes. 474

475 The experimental results in the supervised setting indicate that SMF attains overall better RMSE 476 values compared to NMF and SLIM suggesting that the subspaces encoded in the matrix coefficient of WW' contribute to the learning of more descriptive embeddings for recommendations. While 477 AUROC and AUPRC are commonly used to evaluate classification tasks by showing how true pos-478 itive samples are ranked when predicting associations, practical effectiveness is better measured by 479 the precision at top-K metric. This metric indicates how close the true positives are to the top of the 480 ranking, which is more relevant for suggesting new associations. In this regard, SMF consistently 481 outperforms NMF, SLIM and HCCF, reinforcing the notion that SMF embeddings adeptly capture 482 the distinctive patterns typically learned by reconstructing the data matrix X directly with W and H483 and indirectly by constraining W. 484

485 Although SMF aims at learning from association data, it can be applied to any data represented as a nonnegative matrix. Furthermore, we expect it to be able to integrate extra information. For



Figure 3: Embedding Analysis: Box plots showing the distributions of the Z-score differences across 521 30 different runs of NMF, HCCF and SMF. (a) Pipeline explaining the experiment, first, we have the 522 embedding matrix W, arranged into three different groups, g_1, g_2 , and g_3 . Next, W_{sim} contains the 523 similarities between all the embeddings. Lastly, we calculate if there is a statistical significance in the 524 difference between the intra-class and inter-class similarities from W_{sim} (b) Movielens experiments. The left plot shows the separation between the distributions of the similarities while grouping the 526 users by their gender (Male and Female for this dataset). The right plot shows the separation between the distribution of the similarities while grouping the movies by their genres (18 distinct groups). (c) 527 ModCloth experiments. Separation of the distribution of the similarities while grouping the clothes 528 by their type. (d) ATC-category experiments. From left to right, the grouping advances in the ATC 529 hierarchy, 1^{st} : anatomical, 2^{nd} : therapeutic, and 3^{rd} : pharmacological. 530

instance, to include additional measures of similarity between objects, one could add a term in the
 cost function that penalizes the difference between the learned similarity and the additional similarity
 measure.^{AA}

Like most machine learning models, SMF does not address biases in the training dataset. As a result, objects with more associations are likely to have higher prediction scores than those with fewer associations. Another limitation of SMF is its inefficiency in handling dynamic datasets.
When association data changes, the model must be retrained to ensure that the embeddings reflect the updated state of the objects.^{AA}

540	ACKNOWLEDGMENTS
541	

542 REFERENCES

543

547

582

583

- 544 Charu C Aggarwal et al. *Recommender systems*, volume 1. Springer, 2016.
- Vincent D Blondel, Ngoc-Diep Ho, and Paul Van Dooren. Weighted nonnegative matrix factoriza tion and face feature extraction. *Image and Vision Computing*, pp. 1–17, 2008.
- Deng Cai, Xiaofei He, Jiawei Han, and Thomas S Huang. Graph regularized nonnegative matrix factorization for data representation. *IEEE transactions on pattern analysis and machine intelligence*, 33(8):1548–1560, 2010.
- Benjamin Paul Chamberlain, Sergey Shirobokov, Emanuele Rossi, Fabrizio Frasca, Thomas Markovich, Nils Hammerla, Michael M Bronstein, and Max Hansmire. Graph neural networks for link prediction with subgraph sketching. *The Eleventh International Conference on Learning Representations*, 2023.
- Paolo Cremonesi, Yehuda Koren, and Roberto Turrin. Performance of recommender algorithms on top-n recommendation tasks. In *Proceedings of the fourth ACM conference on Recommender systems*, pp. 39–46, 2010.
- Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936.
- Ehsan Elhamifar and René Vidal. Sparse subspace clustering: Algorithm, theory, and applications.
 IEEE transactions on pattern analysis and machine intelligence, 35(11):2765–2781, 2013.
- Maurizio Ferrari Dacrema, Paolo Cremonesi, and Dietmar Jannach. Are we really making much
 progress? a worrying analysis of recent neural recommendation approaches. In *Proceedings of the 13th ACM conference on recommender systems*, pp. 101–109, 2019.
- 567 Diego Galeano and Alberto Paccanaro. Machine learning prediction of side effects for drugs in clinical trials. *Cell Reports Methods*, 2(12), 2022.
 569
- Diego Galeano, Shantao Li, Mark Gerstein, and Alberto Paccanaro. Predicting the frequencies of drug side effects. *Nature communications*, 11(1):4575, 2020.
- Jiangzhang Gan, Tong Liu, Li Li, and Jilian Zhang. Non-negative matrix factorization: a survey. *The Computer Journal*, 64(7):1080–1092, 2021.
- Carlos A Gomez-Uribe and Neil Hunt. The netflix recommender system: Algorithms, business value, and innovation. ACM Transactions on Management Information Systems (TMIS), 6(4): 1–19, 2015.
- ⁵⁷⁸ Ryuji Hamamoto, Ken Takasawa, Hidenori Machino, Kazuma Kobayashi, Satoshi Takahashi, Am⁵⁷⁹ ina Bolatkan, Norio Shinkai, Akira Sakai, Rina Aoyama, Masayoshi Yamada, et al. Application of
 ⁵⁸⁰ non-negative matrix factorization in oncology: one approach for establishing precision medicine.
 ⁵⁸¹ Briefings in bioinformatics, 23(4):bbac246, 2022.
 - F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. Acm transactions on interactive intelligent systems (tiis), 5(4):1–19, 2015.
- Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. Lightgcn:
 Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*,
 pp. 639–648, 2020.
- Harold Hotelling. Analysis of a complex of statistical variables into principal components. *Journal of educational psychology*, 24(6):417, 1933.
- 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*, pp. 426–434, 2008.

632

- Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.
- Yehuda Koren, Steffen Rendle, and Robert Bell. Advances in collaborative filtering. *Recommender systems handbook*, pp. 91–142, 2021.
- Walid Krichene and Steffen Rendle. On sampled metrics for item recommendation. In *Proceedings* of the 26th ACM SIGKDD international conference on knowledge discovery & data mining, pp. 1748–1757, 2020.
- Michael Kuhn, Ivica Letunic, Lars Juhl Jensen, and Peer Bork. The sider database of drugs and side
 effects. *Nucleic acids research*, 44(D1):D1075–D1079, 2016.
- Daniel Lee and H Sebastian Seung. Algorithms for non-negative matrix factorization. Advances in neural information processing systems, 13, 2000.
- Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.
- Kin Luo, Mengchu Zhou, Yunni Xia, and Qingsheng Zhu. An efficient non-negative matrix factorization-based approach to collaborative filtering for recommender systems. *IEEE Trans- actions on Industrial Informatics*, 10(2):1273–1284, 2014.
- Kin Luo, Zhigang Liu, Long Jin, Yue Zhou, and MengChu Zhou. Symmetric nonnegative matrix factorization-based community detection models and their convergence analysis. *IEEE Transactions on Neural Networks and Learning Systems*, 33(3):1203–1215, 2021.
- Rishabh Misra, Mengting Wan, and Julian McAuley. Decomposing fit semantics for product size
 recommendation in metric spaces. In *Proceedings of the 12th ACM Conference on Recommender Systems*, pp. 422–426, 2018.
- Kia Ning and George Karypis. Slim: Sparse linear methods for top-n recommender systems. In 2011 IEEE 11th international conference on data mining, pp. 497–506. IEEE, 2011.
- V Paul Pauca, Jon Piper, and Robert J Plemmons. Nonnegative matrix factorization for spectral data analysis. *Linear algebra and its applications*, 416(1):29–47, 2006.
- Siyuan Peng, Wee Ser, Badong Chen, and Zhiping Lin. Robust semi-supervised nonnegative matrix factorization for image clustering. *Pattern Recognition*, 111:107683, 2021.
- Steffen Rendle, Walid Krichene, Li Zhang, and John Anderson. Neural collaborative filtering vs.
 matrix factorization revisited. In *Proceedings of the 14th ACM Conference on Recommender Systems*, pp. 240–248, 2020.
- Badrul Sarwar, George Karypis, Joseph Konstan, and John Riedl. Incremental singular value decomposition algorithms for highly scalable recommender systems. In *Fifth international conference on computer and information science*, volume 1, pp. 27–8. Citeseer, 2002.
- Arni Sturluson, Ali Raza, Grant D McConachie, Daniel W Siderius, Xiaoli Z Fern, and Cory M Si mon. Recommendation system to predict missing adsorption properties of nanoporous materials. *Chemistry of Materials*, 33(18):7203–7216, 2021.
- Michael D Sweeney, Luke A Torre-Healy, Virginia L Ma, Margaret A Hall, Lucie Chrastecka, Alisa
 Yurovsky, and Richard A Moffitt. Fastanmf: a fast and stable non-negative matrix factorization
 for gene expression. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 2023.
- Manolis G Vozalis and Konstantinos G Margaritis. A recommender system using principal component analysis. In *Published in 11th panhellenic conference in informatics*, pp. 271–283, 2007.
- 647 Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. Graph neural networks in recommender systems: a survey. *ACM Computing Surveys*, 55(5):1–37, 2022.

- Lianghao Xia, Chao Huang, Yong Xu, Jiashu Zhao, Dawei Yin, and Jimmy Huang. Hypergraph contrastive collaborative filtering. In *Proceedings of the 45th International ACM SIGIR conference* on research and development in information retrieval, pp. 70–79, 2022.
- Haixuan Yang and Cathal Seoighe. Impact of the choice of normalization method on molecular cancer class discovery using nonnegative matrix factorization. *PloS one*, 11(10):e0164880, 2016.
- Wei Feng Yang, Min Wang, and Zhou Chen. Fast probabilistic matrix factorization for recommender
 system. In 2014 IEEE International Conference on Mechatronics and Automation, pp. 1889–
 1894. IEEE, 2014.
 - Ka Yee Yeung and Walter L. Ruzzo. Principal component analysis for clustering gene expression data. *Bioinformatics*, 17(9):763–774, 2001.
 - Seongjun Yun, Seoyoon Kim, Junhyun Lee, Jaewoo Kang, and Hyunwoo J Kim. Neo-gnns: Neighborhood overlap-aware graph neural networks for link prediction. Advances in Neural Information Processing Systems, 34:13683–13694, 2021.
 - Muhan Zhang, Pan Li, Yinglong Xia, Kai Wang, and Long Jin. Labeling trick: A theory of using graph neural networks for multi-node representation learning. *Advances in Neural Information Processing Systems*, 34:9061–9073, 2021.

A APPENDIX

A.1 MODELS LEARNING

SMF: For our model, we optimized the function shown in Equation 5, using the iterative multiplicative update rule described in Equations 6 and 7 for all the datasets.

674 NMF: Here we opted to train an elastic-net regularized NMF (Pauca et al., 2006), where we modify
675 the loss function similarly as with data-driven regularized NMF (Galeano et al., 2020):

$$\min_{W,H} \mathcal{L}_{WNMF}(W,H) = \frac{1}{2} \|P \circ (X - WH)\|_F^2$$
subject to $W, H \ge 0.$
(8)

that can be optimized with the following multiplicative update rules:

$$W_{i,j} \leftarrow W_{i,j} \times \frac{[(P^2 \circ X)H']_{i,j}}{[(P^2 \circ (WH))H' + \lambda_2 W + \lambda_1 \operatorname{sgn}(W)]_{ij}}$$
(9)

$$H_{i,j} \leftarrow Hi, j \times \frac{[W'(P^2 \circ X)]_{ij}}{[W'(P^2(WH)) + \lambda_2 H + \lambda_1 \operatorname{sgn}(H)]_{ij}}$$
(10)

SLIM: here we modified the original loss by replacing the diagonal constraint by adding the trace of M(tr(M)) multiplied by a large number γ as a new term in the loss function (Galeano & Paccanaro, 2022). The loss is also modified similarly as in Equation 5 to tune the importance of the zeros during the learning:

$$\min_{M} \mathcal{L}_{\text{WSLIM}}(M) = \frac{1}{2} \|P \circ (X - MX)\|_{F}^{2} + \gamma tr(M) + \frac{\lambda_{2}}{2} \|M\|_{F}^{2} + \lambda_{1} \|M\|_{1}$$
subject to $M > 0$
(11)

that can be optimized with the following multiplicative update rules:

$$M_{i,j} \leftarrow M_{i,j} \frac{[(P^2 \circ X)X]_{i,j}}{[(P^2(MX))X' + \gamma I + \lambda_2 M + \lambda_1 \operatorname{sgn}(M)]_{i,j}}$$
(12)

702 All learnable parameters for NMF, SLIM and SMF were initialized by sampling from a uniform 703 distribution between 0 and 0.1. It can be proven that in the optimization proposed for Equation 8 and 704 11, both loss functions are non-increasing at their respective parameters converge in a local minimum 705 for NMF and global minimum for SLIM. For Equation 5, we observed that the loss function always 706 increases for the first iteration, then, for the second iteration forward, the loss function turns out to be non-increasing, and it also manages to achieve convergence for every run of the model, what seems to indicate that this version of adaptive gradient descent is well suited for the problem under 708 study in this work. The algorithm for the optimization of Equation 5 was implemented in Matlab 709 R2023a, and the code is included with this submission. The training is stopped by satisfying the 710 stopping criteria $\delta \leq 1e-3$ for the NMF, SLIM and SMF models, and the maximum relative change 711 δ is defined as: 712

713 714

715

716

$$\delta = \frac{\max(\|W_{i,j}^{\text{old}} - W_{i,j}^{\text{new}}\|)}{\max(\|W_{i,j}^{\text{old}}\|)}$$
(13)

where W^{old} and W^{new} are the values of the matrix W after each iteration, clearly the same formula is also applied for H and M.

HCCF: This is a GNN model relying in two message-passing mechanisms. One message-passing occurs between the representation of graph nodes, and the other occurs between hyper-edge representations of the nodes. Both mechanisms are connected by incorporating message and node embeddings generated by the other message-passing, and by the loss function including a contrastive term to connect the two different types of message embeddings:

726 727

729

730

$$\mathcal{L}_{s}^{(u)} = \sum_{i=0}^{I} \sum_{l=0}^{L} -log \frac{exp(s(z_{i,l}^{(u)}, \Gamma_{i,l}^{(u)})/\tau)}{\sum_{i'=0}^{I} exp(s(z_{i,l}^{(u)}, \Gamma_{i',l}^{(u)})/\tau)}$$

728

where, $z_{i,l}^{(u)}$ and $\Gamma_{i,l}^{(u)}$ are the message embeddings for both processes related to the user *i* for the l^{th}) message passing layer. The function s(.) represents a similarity between both embeddings and τ is a temperature constant that tunes the softmax. The overall loss is:

738 739

$$\mathcal{L} = \mathcal{L}_r + \lambda_c (\mathcal{L}_s^{(u)} + \mathcal{L}_s^{(v)}) + \lambda_{wd} ||\Theta||_F^2$$
(14)

where \mathcal{L}_r is the marginal loss, $\mathcal{L}_s^{(u)}$ and $\mathcal{L}_s^{(v)}$ are the contrastive loss for users and items respectively and the last term is weight decay for the learnable parameters. λ_c tunes the importance of the contrastive term in the learning. The final predictions of HCCF are given by:

$$Pr_{i,j} = \Psi_i^{(u)T} \Psi_j^{(v)}$$

where $Psi_i^{(u)}$ and $Psi_j^{(v)}$ are the final embeddings of user *i* and item *j*. We used the code provided by the authors (Xia et al., 2022) in https://github.com/akaxlh/HCCF to train the HCCF model for our experiments.

744 745 A.2 IMPLEMENTATION DETAILS

Experiments were conducted on a machine equipped with two NVIDIA Quadro RTX 6000 GPUs (each with 24 GB of VRAM), an Intel Xeon Gold 6230 processor, and 192 GB of RAM. SMF, NMF, and SLIM were implemented in MATLAB R2023a, while HCCF was implemented in Python 3.6.12 using TensorFlow 1.14.0. CUDA version 11.8 was utilized to leverage GPU acceleration for HCCF training and evaluation.^{AA}

752 A.3 AUROC, AUPRC AND CORRELATION

753

754 Tables 3, 4 and 5 below show extra evaluation metrics used to evaluate the performance of SMF 755 against its competitors. AUROC and AUPRC measure the distribution of true positives in the ranked scores and correlation indicates how close are the predicted scores to the actual levels of association.

756							
757			Table 3: Results for Movielens dataset				
758		MODELS	CORRELATION	AUROC	AUPRC		
759							
760		NMF	$0.5432\pm3{-5}$	$0.9441 \pm 1\mathrm{e}{-7}$	$0.1402\pm5\mathrm{e}{-6}$		
761		SEM	$0.3531\pm3{-7}$	$0.9436\pm4\mathrm{e}{-9}$	$0.1457 \pm 4{-8}$		
762		HCCF	(-)	$0.8686 \pm 4e{-4}$	$0.0611\pm2\mathrm{e}{-4}$		
763		SMF	$0.5714 \pm 6 - 6$	$0.9436 \pm 2e{-7}$	$0.1387 \pm 5 - 6$		
764							
765							
766			Table 4: Results	for Drug-SE datase	t		
767		MODEL S	CODDEL ATION		AUDDC		
768		MODELS	CORRELATION	AUKOC	AUIKC		
769		NMF	$0.7406 \pm 7e - 5$	$0.8819 \pm 4e - 10$	$0.0879 \pm 4e - 10$		
770		SEM	$0.4181 \pm 4e-7$	$0.9268 \pm 6e-8$	$0.1582 \pm 7e-8$		
771		HCCF	(-)	$0.8113 \pm 3e{-5}$	$0.0616 \pm 2e - 5$		
772		SMF	$0.7432 \pm 1\mathrm{e}{-5}$	$0.8582\pm2\mathrm{e}{-8}$	$0.1016\pm4\mathrm{e}{-9}$		
773							
774							
775							
776			Table 5: Results I	for ModCloth datase	et		
777		MODEI	LS CORRELATIO	N AUROC	AUPRC		
778					- 101 1		
779		NMF	$0.0836 \pm 5e-6$	$0.6904 \pm 3e -$	5 1.64e - 4		
780		SEM	$0.0418 \pm 2e-7$	$0.6230 \pm 1e - 1$	10 2.50e-4		
781			(-)	$0.0781 \pm 4e -$	5 1.97e-4		
782		5171L	$01010 \pm 10-0$	J U. 7919 ± 16-	- J 2.04 0-4		
783							
784							
785	A.4	HYPERPARAMETE	R TUNING				

A.4 HYPERPARAMETER TUNING

A validation set with 10% of the interactions for each dataset was used to select an appropriate 787 set of hyperparameters, the decision was based on RMSE and AUPRC measures. The final set of 788 hyperparameters used to perform the experiments are detailed in Tables 6 and 7. 789

790 SLIM does not have a value set for k. It is important to note that the value of α the Movielens and 791 ModCloth datasets are only used for the link prediction task (when evaluating precision, AUROC 792 and AUPRC). For rating prediction (when evaluating RMSE and correlation), α is set to zero, reflecting the fact that there are no true zeros in the dataset. Assuming that in an ideal scenario, where 793 all the users assign a rating to all the movies and clothing items, those values should be between 1 794 and 5. λ_{se} in Equation 5 was always set to one. 795

796 where λ_{wd} , λ_c and τ are shown in Equation 14 and drop is dropout.HCCF model training was 797 run with the predefined parameters in https://github.com/akaxlh/HCCF, except for the parameters in Table 7. The set of parameters was selected after multiple rounds of testing against the same 798 validation set used for hyperparameter tuning of NMF, SLIM and SMF. 799

800 Figure 4 shows the RMSE (orange) and AUPRC (blue) for different hyperparameter values. The red 801 and grey lines divide the plot into different regions where k (ranging from 4 to 16) and α (ranging 802 from 0 to 1) are constant respectively. Note that within a region where k is constant, there are five 803 regions of constant α . Within each of these regions, there are multiple values of the regularization 804 weights in which λ_2 change after each consecutive point and λ_1 remains constant for 4 straight points. Both regularization weights range from zero to one and notably the lowest values of AUPRC 805 correspond to regularization weights set to zero. Therefore, SMF generalizes better with regularized 806 embeddings. Here, we illustrated SMF robustness across a wide range of hyperparameter values.^{AA} 807

808

786

A.5 COMPUTATIONAL COMPLEXITY AND SCALABILITY

10									
11		Table	e 6: Hyperpa	rameters for I	Movielens, I	Drug-SE, a	nd Mod	Cloth	
12	T 7 B		Movielens		Drug-SE			ModClot	h CD (T
13	Values	NMF	SLIM S	$\frac{SMF}{1}$ NMI	SLIM	SMF		SLIM	SMF
14	λ_1	0.5	0.5		0.5	0.5	0.01	0.01	0.01
15	λ_2	0.0	0.5 0.224 (0 + 0.5 0.224 + 0.002	0.0 5 0.224	0.0	0.0	$1 \\ 0.05$	0.5
16	ά	0.224	0.224	0.002	0.224	0.0020	0.05	0.05	0.05
17									
18			Ta	ble 7: Hyperp	arameters for	or HCCF			
19			Values	Movielens	Drug-SE	ModCl	oth		
20									
21			λ_{wd}	1e - 3	1e - 2	1e -	2		
22			λ_c	1e-6	1e-7	1e -	7		
23			au	0.1	0.1	0.1			
24			drop	0.5	0	0			
25									
26									
</td <td>We derived</td> <td>the comp</td> <td>outational tim</td> <td>e complexity</td> <td>of SMF opt</td> <td>imization a</td> <td>algorith</td> <td>m from eq</td> <td>uations 5,</td>	We derived	the comp	outational tim	e complexity	of SMF opt	imization a	algorith	m from eq	uations 5,
20 20	and 7 that a	re used to	o check the c	onvergence c	riterion and	update W	and H ,	respective	ely. As eac
29	equation has	a fixed i	number of m	atrix operation	ns, it suffice	s to derive	the asy	mptotic co	omplexity of
24	the most exp	bensive o	one, which in (WW) and	volves matrix	multiplicati	ion. Specif	ically, t	he compu	tation of the
20	data matrix	$X \subset \mathbb{R}^{n \times n}$	$\langle m \rangle$ Other mu	A runs in une	$O(n \cdot m),$ such as WF	where n and \overline{I} run in the	$\int \frac{d}{dt} \frac{d}{dt}$	(k,m) w	$\frac{1510115}{bere} \frac{b}{k} \frac{1}{15} 1$
3∠ 22	dimension o	$A \in \mathbb{R}$ f the emb	edding space	Assuming t	hat $k < < n$	m we have	ve that t	he time cc	molexity of
24	each iteratio	n of the S	SMF algorith	m is $O(n^2 \cdot n)$	n). ^{AA}	, <i>m</i> , we na	ve that t		mpienty
34 25				1	CENC	(0)	\\ 1 ."		
38	SMF iteratio	ons have $(O(m))$	the same tin	ne complexity	as SEM's	$(O(n^2 \cdot m))$	()), while a state of the product	three own	erations ai
37	obtained the	mean if	$(\kappa \cdot m)$. To eration time	(in seconds)	in two datas	unning uni sets with d	e of the ifferent	sizes (Mo	wielens an
38	ModCloth).	This is s	hown in the f	irst two colun	in two datas	8. As SMI	F's mult	inlicative	update rule
39	involve mor	e matrix	multiplicatio	ons than those	of NMF ar	nd SEM, it	is expe	ected to ha	ave a high
40	execution tin	ne per ite	eration. AA			· · · · · ·	1		0
41	A fundamer	- tal variai	ble for analy	zing the over	all compute	tional time	of the	SME aloo	rithm is th
42	total number	r of itera	tions until co	nvergence A	s it is not n	ossible to a	derive it	analytica	illy here w
43	show only th	e empiri	cal total num	ber of iteration	ns. The last	two colum	is of Tal	ble 8 show	the number
44	of iterations	that were	e necessary t	o achieve con	vergence for	r SMF, SM	F, and S	EM. ^{AA}	
45	The sector:	: f . l		4h	- 41 41 4			41	
46	nutational st	ny of the	Like NME	SME faces kn	nuy neu to t	ons when a	r, as oc	o large so	ale dataset
47	where the co	mnutati	onal demand	s can hinder e	officiency (C	an et al	2021) 7	Chis reflec	ts a broade
48	challenge in	scaling r	natrix factori	zation technic	ues to acco	mmodate i	ncreasir	ngly larger	· datasets. ^A
49	A 11 /1	4 1		1	1 1 1		•		
50	All three me	thous rec	juire addition	al memory the	at exceeds the (n_1, m_2) where (n_2, m_2)	ne size of th	ie input	. Asympto	bath SEM
51	and SMF red	uire space	complexity a $O(n^2)$ to a	store the simi	$(n \cdot m)$ will arity matrix	AA	. when	m >> m	, bour ser
52		iune spa		store the shift					
53									
54			Table 8: M	ean iteration	time and nur	mber of ite	rations		
55	Ν	Iodels	Movielens ((s) ModClo	th (s) Ma	ovielens (it) Mo	dCloth (it	t)
56									
57	Ν	MF	0.01531	1.423	87	1868.20		467	
58	S	EM	0.4338	7.39	48	244.20		292	
59	S	MF	0.02737	7.46	93	1427.73		1032	
60									
61									



Figure 4: Sensibility to hyperparameter setting: Different points are the RMSE (orange) and AUPRC (blue) of SMF trained with different hyperparameters. Vertical red dotted lines divided the plot into regions where k is constant. Vertical dotted gray lines divided the plot into regions where α is constant ($\alpha_1 = 0, \alpha_2 = 0.0025, \alpha_3 = 0.05, \alpha_4 = 0.22$ and $\alpha_5 = 1$). Within each of these regions, both regularization weights change. λ_2 is different for each consecutive point, and is set to 0, 0.01, 0.5 and 1 in that order. λ_1 takes the same values, but it remains the same for 4 straight points before being changed to the next value. The first points at the left of each region correspond to models with both regularization weights set to zero.