

Understanding and Improving Relational Matrix Factorization in Recommender Systems

Li Pu

Artificial Intelligence Laboratory, EPFL
INR240, EPFL-IC-LIA, Station 14
CH-1015, Lausanne, Switzerland
li.pu@epfl.ch

Boi Faltings

Artificial Intelligence Laboratory, EPFL
INR230, EPFL-IC-LIA, Station 14
CH-1015, Lausanne, Switzerland
boi.faltings@epfl.ch

ABSTRACT

Matrix factorization techniques such as the singular value decomposition (SVD) have had great success in recommender systems. We present a new perspective of SVD for constructing a latent space from the training data, which is justified by the theory of hypergraph model. We show that the vectors representing the items in the latent space can be grouped into (approximately) orthogonal clusters which correspond to the vertex clusters in the co-rating hypergraph, and the lengths of the vectors are indicators of the representativeness of the items. These properties are used for making top- N recommendations in a two-phase algorithm. In this work, we provide a new explanation for the significantly better performance of the asymmetric SVD approaches and a novel algorithm for better diversity in top- N recommendations.

Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*information filtering*

Keywords

recommender system, matrix factorization, hypergraph

1. INTRODUCTION

In a recommender system, there are two types of entities, i.e. users and items. The users give weighted connections, or ratings, to the items. These connections can be modeled as a rating matrix \mathbf{R} where the rows represent the users and the columns represent the items. The entry $\mathbf{R}(i, k)$ is the rating given by user i to item k (if the rating exists). The goal of the recommender system is to predict unseen items that might be rated with high scores. A common assumption for making predictions is that certain latent factors are associated with each user and each item, and a user gives ratings according to her latent factors and the latent factors of the items.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
RecSys '13, October 12–16, 2013, Hong Kong, China.
Copyright 2013 ACM 978-1-4503-2409-0/13/10 ...\$15.00.
<http://dx.doi.org/10.1145/2507157.2507178>.

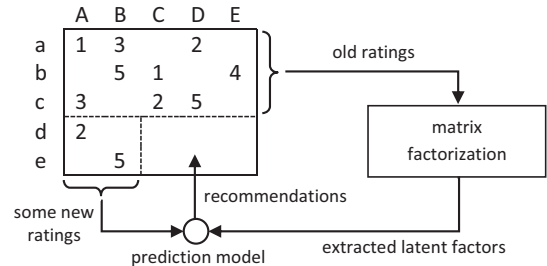


Figure 1: Recommender system using matrix factorizations.

For example, in the *Lastfm* dataset, the items are the artists in a music website. The latent factors could be the genres of the artists. The recommender system first obtains the genres of the artists, and the users' interests on different genres. Then the system makes recommendations by aligning the user's interests to the genres of the artists.

For ratings modeled as the matrix \mathbf{R} , various matrix factorization techniques have been applied to extract the latent factors and make recommendations. Figure 1 shows a typical recommender system using matrix factorization. The old ratings given by existing users are first decomposed into multiple matrices to extract the latent factors. When new users come into the system, their ratings are combined with the extracted latent factors to make recommendations.

Sarwar et al. [17] proposed one of the early works that use the low-rank Singular Value Decomposition (SVD) approximation to replace the missing values in the original rating matrix, i.e. $\mathbf{R} \approx \mathbf{U}_l \mathbf{S}_l \mathbf{V}_l^\top$ where the sizes of the matrices are $\mathbf{R}: m \times n$, $\mathbf{U}_l: m \times l$, $\mathbf{S}_l: l \times l$, and $\mathbf{V}_l: n \times l$. The diagonal matrix \mathbf{S}_l contains the l biggest singular values of \mathbf{R} , and this low-rank approximation is also called the *truncated SVD*. This approach relies on the fact that the rank- l truncated SVD approximation $\hat{\mathbf{R}} = \mathbf{U}_l \mathbf{S}_l \mathbf{V}_l^\top$ minimizes the Frobenius norm $\|\hat{\mathbf{R}} - \mathbf{R}\|_F$ among all the rank- l matrices, and we hope that the matrix $\hat{\mathbf{R}}$ would generalize the correct ratings to the missing entries in \mathbf{R} . Note that imputation of missing values is implicitly required when computing the truncated SVD. We can represent user i with the vector $\bar{\mathbf{u}}_i = \mathbf{S}_l^{-1/2} \mathbf{U}_l(i, \cdot)^\top$ and item k with the vector $\bar{\mathbf{v}}_k = \mathbf{S}_l^{-1/2} \mathbf{V}_l(\cdot, k)^\top$. The predicted rating from user i to item k is then $\hat{\mathbf{R}}(i, k) = \bar{\mathbf{u}}_i^\top \bar{\mathbf{v}}_k$, where $\bar{\mathbf{u}}_i$ and $\bar{\mathbf{v}}_k$ can be considered as latent vectors.

During the Netflix prize competition, Funk [5] proposed an iterative way to compute the low-rank SVD approxima-

tion on a very big matrix that minimizes only the Frobenius norm on the non-empty entries in the original rating matrix. Based on Funk’s work, Paterek [14] suggested a new model to reduce the number of parameters that need to be stored in the recommender system. In Paterek’s approach, a truncated SVD is computed as before, but only the right singular vectors in \mathbf{V}_i are taken into the second step. In the second step, item k is represented by \mathbf{v}_k which is the k -th column of \mathbf{V}_i^\top , while user i ’s latent vector is computed by a linear combination of the \mathbf{v}_k ’s whose corresponding items are rated by user i . Finally the rating prediction is the inner product of a user’s latent vector and an item’s latent vector. Koren [9] pointed out that this “asymmetric SVD” which takes only the right singular vectors actually works better in prediction accuracy. However, besides the empirical studies, the reason of such improvement for asymmetric SVD is not clear.

The **main contributions** of this paper include (1) a new explanation of the (approximate) orthogonal structure of the latent space constructed from the asymmetric SVD, which is used as an intermediate step in a top- N recommender system, and (2) a novel algorithm based on a hypergraph formulation that brings more diversity in the recommendations by applying normalizations in the asymmetric SVD.

In the remainder of the paper, we start with the hypergraph representation and the normalized hypergraph cut. We show that the asymmetric SVD approach follows the hypergraph learning scheme, and there is an (approximate) orthogonal structure associated with the right singular vectors. Based on the lengths of the latent vectors obtained from the right singular vectors, the items can be categorized into anchor items and non-anchor items. Finally we present a top- N recommendation algorithm that utilizes the above structures, and conduct experiments to compare our algorithm with the state-of-the-art approaches.

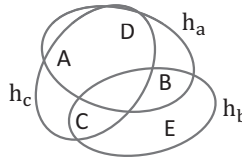
2. PRELIMINARIES

Suppose we have two sets of entities: the user set $Y = \{y_1, y_2, \dots, y_m\}$ of size m , and the item set $Z = \{z_1, z_2, \dots, z_n\}$ of size n . The items in set Z belong to s clusters $\{C_1, C_2, \dots, C_s\}$. Each cluster C_j is a subset of Z and $C_j \cap C_{j'} = \emptyset$ for $j \neq j'$.

Taking the `Lastfm` dataset as an example again, let the set Y be the set of users and the set Z be the set of artists. Clusters or partitions over Z can be established according to the genres of the artists. For example, C_1 includes the rock music artists, while all the country music artists are in C_2 . Each user in the set Y , on the other hand, is assumed to have a preference over music genres. If the user likes rock music, she would listen to the artists from C_1 with high probability. If the user does not like country music, she would rarely listen to the artists from C_2 .

The mutually exclusion between clusters of items is the basic assumption in our analysis. Each cluster can be considered as a latent factor by which the ratings are generated. In some cases, however, an artist could play in multiple genres. This does not pose a problem to our basic assumption, because we can create a cluster (latent factor) that corresponds to a combination of genres. In fact, if the intuition behind the latent factors follows some categorical attributes, e.g. genres, languages, professions, etc., or a combination of categorical attributes, our basic assumption always holds.

A hypergraph from the old ratings.
There are 3 hyperedges $\{h_a, h_b, h_c\}$.



The induced graph from the hypergraph

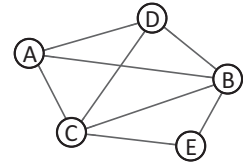


Figure 2: The hypergraph representation of the ratings in Figure 1 and the induced graph.

2.1 From Ratings to Relations

The matrix \mathbf{R} contains the ratings between the users and the items, which can be seen as relations between two sets. For simplicity, we assume that all the ratings are bigger than zero, thus the empty entries in the $m \times n$ rating matrix \mathbf{R} can be represented by zeros. If we ignore the actual ratings in \mathbf{R} , we can obtain a binary matrix $\mathbf{X} = (\mathbf{R} > 0)$ that represents only the relations, i.e. $\mathbf{X}(i, k) = 1$ if $\mathbf{R}(i, k) > 0$, otherwise $\mathbf{X}(i, k) = 0$, and the resulting matrix \mathbf{X} can be precisely represented by a hypergraph.

A *hypergraph* is an extension of a graph for modeling relations. In a graph, an edge connects exactly two vertices, while in a hypergraph a *hyperedge* connects any number of vertices. For a binary relation matrix \mathbf{X} , the items are considered as the vertices of a hypergraph, and a user is represented by a hyperedge which contains all the items (vertices) that have been rated by the user. Following the tradition of visualizing a hypergraph in literatures, Figure 2 shows an example of the hypergraph representation. The matrix \mathbf{X} is also called the *incident matrix* of the hypergraph. The *vertex degree* (of an item) is the sum of the k -th column of \mathbf{X} , i.e. $\text{deg}(z_k) = \sum_i \mathbf{X}(i, k)$. The *hyperedge degree* (of a user) is the sum of the i -th row of \mathbf{X} , i.e. $\text{deg}(y_i) = \sum_k \mathbf{X}(i, k)$.

Recall that we assume there are some clusters of items, and these clusters correspond to the latent factors. Thus extracting the latent factors is equivalent to finding clusters of items (vertices) in the hypergraph. There is a wide range of studies for clustering the vertices in a hypergraph [8, 1, 22]. In this paper, we focus on one particular line of studies that convert a hypergraph into an approximated graph so that the graph clustering algorithms can be applied to the hypergraph. Zhou et al. proposed a transformation to convert a hypergraph into a graph based on the *Normalized Hypergraph Cut* (NHC) [22] (others include the clique-expansion, the start-expansion [1], and the hyperedge-expansion [15]). We denote the transformation proposed in [22] as “NHC”, and use it to find clusters in the hypergraph.

By applying the NHC transformation, the vertices of the hypergraph are copied into a new graph. Then each hyperedge is converted into a clique of weighted edges in the new graph that connect all the vertices in the original hyperedge. The weights of the edges in the clique are normalized by the hyperedge degree so that the contribution of each hyperedge to the new graph is the same. With a linear combination of all the cliques from all the hyperedges, the new graph can be seen as an approximation of the hypergraph in the sense that any two vertices connected by some hyperedges are also connected in the new graph. We call this new graph the *induced graph* (see Figure 2).

One can show that the adjacency matrix of the induced graph is $\mathbf{X}^\top \mathbf{D}_e^{-1} \mathbf{X}$, where \mathbf{X}^\top is the transpose of \mathbf{X} and $\mathbf{D}_e = \text{diag}(\mathbf{X}\mathbf{1})$ is the hyperedge degree matrix. \mathbf{D}_e is actually the diagonal matrix of row sums of \mathbf{X} , and $\mathbf{1}$ is an all-ones vector of proper length. With the NHC transformation, it is proposed to use the normalized Laplacian of the induced graph to find clusters of vertices [22], where the normalized Laplacian matrix \mathbf{L} is defined as

$$\mathbf{L} = \mathbf{I} - \mathbf{D}_v^{-1/2} \mathbf{X}^\top \mathbf{D}_e^{-1} \mathbf{X} \mathbf{D}_v^{-1/2}. \quad (1)$$

Matrix \mathbf{I} is an identity matrix, and $\mathbf{D}_v = \text{diag}(\mathbf{1}^\top \mathbf{X})$ is the vertex degree matrix.

As in standard spectral graph theory [4], the clustering algorithm involves computing the leading (smallest) l eigenvalues $\{\lambda_1, \lambda_2, \dots, \lambda_l\}$ of \mathbf{L} and the corresponding eigenvectors $[\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_l] = \mathbf{F}$. Then a vertex (or an item z_k) can be represented by a vector α_k that is the k -th column of \mathbf{F}^\top (note that \mathbf{F} is of size $n \times l$). We call the l -dimensional vector α_k the *latent vector*, and the l -dimensional space the *latent space*. Once the latent vectors are computed, one can apply any clustering algorithm (such as k-means) in the latent space to find clusters of vertices.

This scheme of spectral technique (the list of eigenvalues are often called the *spectrum* of the matrix) has been widely used in many areas. It has close connections to the Principal Component Analysis (PCA), SVD, and many other matrix factorization techniques. In the next section, we show that the NHC transformation of a hypergraph is equivalent to a SVD, and under certain conditions there is a structure with the resulting latent vectors. Note that although the actual ratings are ignored in the hypergraph representation, they will be used again in the top- N recommendation algorithm presented in Section 4.

3. HYPERGRAPH LEARNING USING SVD

The normalized Laplacian \mathbf{L} is a symmetric matrix. Let

$$\bar{\mathbf{X}} = \mathbf{D}_e^{-1/2} \mathbf{X} \mathbf{D}_v^{-1/2}. \quad (2)$$

We can rewrite $\mathbf{L} = \mathbf{I} - \bar{\mathbf{X}}^\top \bar{\mathbf{X}}$. The biggest l eigenvalues of $\bar{\mathbf{X}}^\top \bar{\mathbf{X}}$ are exactly $\{1 - \lambda_1, 1 - \lambda_2, \dots, 1 - \lambda_l\}$, while the corresponding eigenvectors are still $[\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_l] = \mathbf{F}$. Then we decompose the matrix $\bar{\mathbf{X}}$ by the full SVD $\bar{\mathbf{X}} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$, where \mathbf{U} and \mathbf{V} are unitary matrices, and \mathbf{S} is a rectangular diagonal matrix. The biggest l singular values in \mathbf{S} are exactly $\{\sqrt{1 - \lambda_1}, \sqrt{1 - \lambda_2}, \dots, \sqrt{1 - \lambda_l}\}$, and the columns of \mathbf{V} (right-singular vectors) are the eigenvectors of $\bar{\mathbf{X}}^\top \bar{\mathbf{X}}$. Therefore, instead of computing the eigen-decomposition of \mathbf{L} , we can obtain \mathbf{F} from the SVD of $\bar{\mathbf{X}}$, i.e. \mathbf{F} is the submatrix of the first l columns of \mathbf{V} .

Since we are looking for the right singular vectors of $\bar{\mathbf{X}}$ associated with the largest l singular values, computing \mathbf{F} can be done by the truncated SVD $\bar{\mathbf{X}} \approx \mathbf{U}_l \mathbf{S}_l \mathbf{V}_l^\top = \mathbf{U}_l \mathbf{S}_l \mathbf{F}^\top$. Thus the NHC transformation is equivalent to the asymmetric truncated SVD of $\bar{\mathbf{X}}$. We call the column vectors in $\bar{\mathbf{X}}$ the “profiles” of the items.

There are many advantages of using SVD to compute the α_k ’s (latent vectors) rather than the eigen-decomposition of \mathbf{L} . Firstly, the matrix $\bar{\mathbf{X}}$ is usually sparse, but $\bar{\mathbf{X}}^\top \bar{\mathbf{X}}$ might be non-sparse. When n is large, the computational cost and storage cost of the eigen-decomposition might be impractical. Secondly, there are existing approaches to implement SVD incrementally, which allows us to just compute the minor changes when modifying $\bar{\mathbf{X}}$.

3.1 A Special Case

We first consider a special case where all the items in one cluster are rated by the same set of users. In other words, each cluster C_j is associated with a subset of users and these users have rated all the items in C_j . We call this special case the “full concentrated case” because it can be modeled with a beta distribution with a small concentration parameter. If there are s clusters, the matrix $\bar{\mathbf{X}}$ can be written as

$$\bar{\mathbf{X}} = \begin{bmatrix} \underbrace{\bar{x}_1 \cdots \bar{x}_1}_{|C_1| \text{ vectors}} & \underbrace{\bar{x}_2 \cdots \bar{x}_2}_{|C_2| \text{ vectors}} & \cdots & \underbrace{\bar{x}_s \cdots \bar{x}_s}_{|C_s| \text{ vectors}} \end{bmatrix}, \quad (3)$$

where the column vectors in one cluster are all the same. We call the items in the s clusters the *anchor items* because they define the cluster centers and remain clustered in the latent space.

In spectral graph theory, if the graph contains s disconnected components, one can show that the latent vectors obtained from the first s eigenvectors of the graph Laplacian are aligned with the axes of the s -dimensional latent space [13]. When the disconnected components are connected by some weak links, the latent vectors are still aligned with s approximately orthogonal lines [20]. However, the induced graph of a hypergraph usually does not consist of more than one connected component. Figure 3 shows an example of a hypergraph in full concentrated case. The adjacency matrix of the induced graph is not a diagonal block matrix (not even approximately), so we need to develop a new analysis to support the clustering algorithm for hypergraphs.

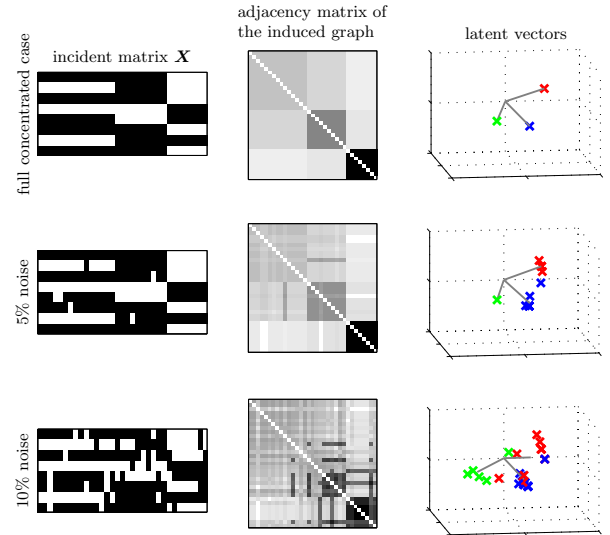


Figure 3: Examples of the incident matrix and the latent vectors with the full concentrated case and noisy cases. The noises are added by randomly flipping some entries of \mathbf{X} . Approximately orthogonal structures can be observed in the noisy cases. In the charts of latent vectors, the gray lines are from the origin to the cluster centers.

If the rank of $\bar{\mathbf{X}}$ is s , i.e. the anchor item profiles \bar{x}_j ’s are linearly independent, the truncated SVD $\bar{\mathbf{X}} = \mathbf{U}_l \mathbf{S}_l \mathbf{V}_l^\top$ is an exact decomposition when $l = s$. Recall that the latent vectors in \mathbf{F} can be obtained from $\mathbf{F} = \mathbf{V}_l$. The following statement shows the structure of \mathbf{F} .

PROPOSITION 1. If $\bar{\mathbf{X}}$ is in the form of Eq. (3) and $\bar{\mathbf{x}}_j$'s are linearly independent, the latent vectors $\mathbf{F} = [\alpha_1 \cdots \alpha_n]^\top$ given by the truncated SVD $\bar{\mathbf{X}} = \mathbf{U}_s \mathbf{S}_s \mathbf{F}^\top$ (columns of \mathbf{F} are normalized to 1) can be grouped by the clusters $\alpha_k = \beta_j$ for $\forall z_k \in C_j$. The α_k 's from the same cluster are identical (denoted as β_j). Furthermore, we have $\beta_j^\top \beta_{j'} = 0$ for $\forall j \neq j'$, and $\beta_j^\top \beta_j = |C_j|^{-1}$ for $\forall j \in \{1, 2, \dots, s\}$. (See appendix for proof.)

If all the items are anchor items, i.e. in the form of Eq. (3) and having linearly independent $\bar{\mathbf{x}}_j$'s, the latent vectors would be in exactly s clusters whose centers are orthogonal to each other (see Figure 3). The length of the latent vector is the reciprocal of the cluster size. Unlike the graph case, the latent vectors computed from a hypergraph are usually not aligned with the axes of the latent space, but up to a rotation that depends on the cluster sizes and the distribution of non-zero entries in \mathbf{X} . Nevertheless, the orthogonal structure of the latent vectors still provides us a similar theoretical foundation for clustering vertices as in the standard spectral graph theory.

3.2 More General Case

The assumption that all the items are anchor items is too strict. We usually have some item profiles which are combinations of some other item profiles. For example, there are pure ‘‘SciFi’’ books and pure ‘‘romance’’ books that attract distinct sets of readers (the anchor items). But one book might have the ‘‘SciFi’’ element and the ‘‘romance’’ element at the same time, so it could attract readers from both sets. The item profile of this book is apparently a combination of two profiles.

Besides the anchor items, we now consider the *non-anchor items* that are combinations of anchor items. We assume that the weights of the combinations are always non-negative, which implies that only additions of profiles are allowed, but not subtractions. If there are n' non-anchor items and s clusters of anchor items, $\bar{\mathbf{X}}$ can be rewritten as

$$\bar{\mathbf{X}} = [\underbrace{\bar{\mathbf{x}}_1 \cdots \bar{\mathbf{x}}_{n'}}_{n' \text{ vectors}} \underbrace{\bar{\mathbf{x}}_1 \cdots \bar{\mathbf{x}}_1}_{|C_1| \text{ vectors}} \underbrace{\bar{\mathbf{x}}_2 \cdots \bar{\mathbf{x}}_2}_{|C_2| \text{ vectors}} \cdots \underbrace{\bar{\mathbf{x}}_s \cdots \bar{\mathbf{x}}_s}_{|C_s| \text{ vectors}}], \quad (4)$$

where the profile of non-anchor item z'_k is denoted as $\bar{\mathbf{x}}'_k$ ($k \in \{1, 2, \dots, n'\}$), and the profile of an anchor item from cluster C_j is still denoted as $\bar{\mathbf{x}}_j$. Note that $\bar{\mathbf{x}}'_k$ is not an exact linear combination of the anchor items, because when combining the profiles of the anchor items, there might be redundant entries. For example, if $\bar{\mathbf{x}}_1 = [1 \ 1 \ 0 \ 0]^\top$, $\bar{\mathbf{x}}_2 = [0 \ 1 \ 1 \ 0]^\top$ and $\bar{\mathbf{x}}'_k$ is a combination of $\bar{\mathbf{x}}_1$ and $\bar{\mathbf{x}}_2$, $\bar{\mathbf{x}}'_k$ would be $[1 \ 1 \ 1 \ 0]^\top$ instead of $[1 \ 2 \ 1 \ 0]^\top$ because the incident matrix \mathbf{X} is always a binary matrix (normalizations are ignored in this example for simplicity). Thus we need to subtract the redundant entries when summing the anchor item profiles:

$$\bar{\mathbf{x}}'_k = \sum_{j \in I(z'_k)} \sqrt{\frac{\deg(z_j)}{\deg(z'_k)}} \bar{\mathbf{x}}_j - \mathbf{r}, \quad (5)$$

where $I(z'_k)$ is the set of anchor items from which z'_k is constructed, and \mathbf{r} denotes the column vector that consists of the redundant entries in the profiles of the anchor items. Obviously all entries of \mathbf{r} are non-negative.

Denote the latent vector of an anchor item from cluster C_j as β_j , and the latent vector of a non-anchor item z'_k as β'_k .

By the SVD on $\bar{\mathbf{X}}$, β'_k can be expressed by the combination

$$\beta'_k = \sum_{j \in I(z'_k)} \sqrt{\frac{\deg(z_j)}{\deg(z'_k)}} \beta_j - \mathbf{S}_l^{-1} \mathbf{U}^\top \mathbf{r}. \quad (6)$$

If the entries in \mathbf{r} are small enough to be considered as residuals, or in other words the anchor items are associated with disjoint sets of users, we could still expect approximately orthogonal structure between the latent vectors of the anchor items.

PROPOSITION 2. If the anchor items are associated with disjoint sets of users, i.e. $\bar{\mathbf{x}}_j^\top \bar{\mathbf{x}}_{j'} = 0$ for $\forall j, j' \in \{1, 2, \dots, s\}$, $j \neq j'$, the latent vectors of the anchor items are approximately orthogonal $\beta_j^\top \beta_{j'} \approx 0$ for $\forall j \neq j'$.

The conclusion of Proposition 2 follows from the fact that $\bar{\mathbf{x}}_j^\top \bar{\mathbf{x}}_{j'} = \beta_j^\top \mathbf{S}_l^2 \beta_{j'} = 0$, and the first l singular values in \mathbf{S}_l are approximately the same when the cluster sizes are similar.

In practice, we can consider the anchor items as the most representative items of each cluster. For example, the pure ‘‘SciFi’’ books are the anchor items of the ‘‘SciFi’’ cluster, because they have the most unique features that define what is ‘‘SciFi’’. The anchor items are not necessarily the most popular items or the most informative items in each cluster, but the items that are most distant from all the other clusters and attract a very unique group of users. Usually the anchor items are in the long-tail part, i.e. less popular.

Although the anchor items exhibit the approximately orthogonal structure, we do not know which items are the anchor items. The following result suggests that it is possible to distinguish the anchor items and the non-anchor items by examining the length of the latent vectors.

PROPOSITION 3. For a non-anchor item z'_k , denote $I(z'_k)$ the set of anchor items from which z'_k is constructed. Suppose that the anchor items in $I(z'_k)$ are non-overlapping, i.e. $\bar{\mathbf{x}}_j^\top \bar{\mathbf{x}}_{j'} = 0$ for $\forall j, j' \in I(z'_k)$, $j \neq j'$, and the cluster sizes of the anchor items are similar, i.e. the lengths of the anchor latent vectors are similar (denoted as $\|\beta_j\|_2 \approx b_0$ for $\forall j \in I(z'_k)$, where $\|\cdot\|_2$ is the l_2 norm). We have $\|\beta'_k\|_2 < b_0$. (See appendix for proof.)

This result states that the length of a non-anchor latent vector is smaller than the lengths of the anchor latent vectors from which the former is constructed. Therefore, we can find the anchor items by sorting the lengths of the latent vectors in descending order.

In order to verify the conclusions in Propositions 2 and 3, we compute the latent vectors by the rank- l truncated SVD of $\bar{\mathbf{X}}$, and select $10 \cdot l$ items whose latent vectors have the longest lengths. Then the selected latent vectors are grouped into clusters by the k-means algorithm. By the conclusions of the Propositions, the cluster centers of the selected latent vectors should be approximately orthogonal. In Figure 4, we show the results from a real dataset with NHC and SVD. Because there is no hyperedge-degree and vertex-degree normalizations in SVD, the longest latent vectors are not necessarily the anchor items. We can observe that there exists an approximately orthogonal structure in the NHC results, while such structure is less clear with SVD.

There are existing works about the categorization of the items or users based on their impact on the recommender

system, e.g. [12]. Our classification of anchor and non-anchor items provides a new perspective in the context of latent space. In summary, (1) if the latent vectors are all orthogonal, each corresponds to a different anchor item (the full concentrated case). (2) if they are not, there is an orthogonal subset such that each latent vector in this subset corresponds to a different anchor item, and all the others (non-anchor items) are approximately in the subspace spanned by the orthogonal subset. (3) generally, the anchor items are rated by different groups of users.

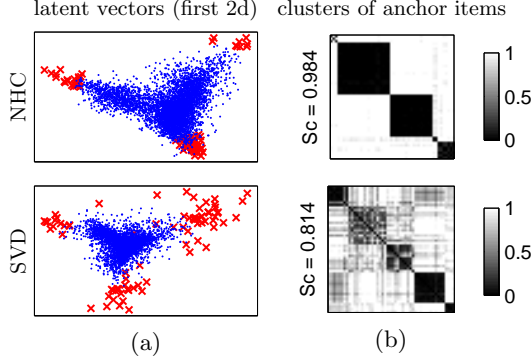


Figure 4: The selected anchor items with longer latent vector lengths in the dataset BookCrossing (see Section 6 for the description of the dataset). (a) The anchor items are shown in red x-markers. (b) The cosine distance matrix of the clustered anchor items. Sc is the sum of average distances between clusters minus the sum of average distances within clusters, normalized by the largest possible value with a perfect orthogonal structure. A Sc value close to 1 suggests a better orthogonal structure.

4. THE NORMALIZED RELATIONAL SVD

We now present a new algorithm based on the NHC Laplacian and the above analysis. Once the latent vectors are computed for the items, we model the user y_i 's interests (or latent profile) as another vector θ_i in the same latent space. The predicted score for user y_i on item z_k is $P(y_i, z_k) = \theta_i^\top \alpha_k$, and the vector θ_i should take the values such that the predictor P coincides with the existing ratings (or as close as possible).

In order to simulate the scenario where the recommendations are made for new users and evaluate the algorithm, we split the rating matrix \mathbf{R} into three parts

$$\mathbf{R} = \begin{bmatrix} \mathbf{R}_m \\ \mathbf{R}_{tr}, \mathbf{R}_{te} \end{bmatrix} \begin{matrix} \text{subset } Y_m \\ \text{subset } Y_t \end{matrix}. \quad (7)$$

Firstly the users are separated into two subsets. Y_m is the subset of old users, and Y_t is the subset of new users. The sub-matrix \mathbf{R}_m contains all the ratings from the users in Y_m , and the ratings in \mathbf{R}_m are used for constructing the hypergraph and computing the latent vectors. The ratings from the users in Y_t are further split into two parts: \mathbf{R}_{tr} and \mathbf{R}_{te} . \mathbf{R}_{tr} contains the “known” ratings of the new users. Based on \mathbf{R}_{tr} and the latent vectors, we can predict or recommend more items for the users in Y_t . Then the predictions are evaluated by the ratings in \mathbf{R}_{te} . Note that only \mathbf{R}_m and \mathbf{R}_{tr} are taken as the inputs of the algorithm. \mathbf{R}_{te} is used for evaluation.

Since our algorithm works with the hypergraph transformation and SVD, we call it HSVD, which contains the following steps:

Step 1: Compute the l -dimensional NHC transformation \mathbf{F} by the truncated SVD $\bar{\mathbf{X}}_m \approx \mathbf{U}_l \mathbf{S}_l \mathbf{F}^\top$, where $\bar{\mathbf{X}}_m$ is computed from the binary matrix $\mathbf{X}_m = (\mathbf{R}_m > 0)$ (see Eq. (2)). Normalize each of the l columns of \mathbf{F} to 1 (l_2 norm).

Step 2: For each user $y_i \in Y_t$, we choose the θ_i that minimizes the error $\|\mathbf{R}_{tr}(i, \cdot) - \theta_i^\top \mathbf{F}^\top\|_2^2$. The predicted score vector is $\hat{r}_i = \theta_i^\top \mathbf{F}^\top$. In matrix form, the full prediction matrix is $\hat{\mathbf{R}}_{te} = \Theta^\top \mathbf{F}^\top$, where Θ is obtained by solving the linear system $\mathbf{F}\Theta = \mathbf{R}_{tr}^\top$ in a least-squares sense.

Step 3: For each user $y_i \in Y_t$, recommend the top- N unrated items with the highest scores in $\hat{\mathbf{R}}_{te}$.

There are several normalizations in the HSVD algorithm which did not exist in previous matrix factorization methods. When computing $\bar{\mathbf{X}}_m$ in step 1, the hyperedge-degree normalization (from the induced graph transformation itself) and the vertex-degree normalization (from the normalized Laplacian of the induced graph) ensure that the anchor latent vectors would have longer lengths. By Proposition 1, the normalization of the columns of \mathbf{F} in step 1 is also indispensable for producing an (approximately) orthogonal structure in the latent space.

When computing the θ_i in step 2, two parts are actually taken into consideration, which is illustrated in Figure 5. Firstly we consider the latent vectors of the items which are rated by user y_i (the triangle points in Figure 5). We would like to choose a θ_i that is close to these rated α_k 's. Since the row vector $\mathbf{R}(i, \cdot)$ is weighted by the ratings, θ_i should be even closer to the α_k 's with higher ratings, which can be expressed as $\theta_i = \arg \min_{\theta} \sum_{k, \mathbf{R}(i, k) > 0} (\mathbf{R}(i, k) - \theta_i^\top \alpha_k)^2$. Because every non-anchor item is approximately a linear combination of some anchor items (e.g. α_3 is a combination of α_1 and α_5 in Figure 5), θ_i should be in the subspace spanned by the latent vectors of the underlying anchor items. Secondly, for those α_k 's that are not rated by user y_i (the circle points in Figure 5), θ_i should stay as far from them as possible, ideally orthogonal to these unrated latent vectors. The orthogonality between the cluster centers of the anchor items suggests that it is possible to find a θ_i which is orthogonal to the anchor items that the user y_i is not interested in. This can be expressed as $\theta_i = \arg \min_{\theta} \sum_{k, \mathbf{R}(i, k) = 0} (\theta_i^\top \alpha_k)^2$. Putting the two parts together, it suggests that in matrix form the θ_i 's can be obtained by solving the linear system $\mathbf{F}\Theta = \mathbf{R}_{tr}^\top$ in a least-squares sense, which is in the step 2 of the HSVD algorithm.

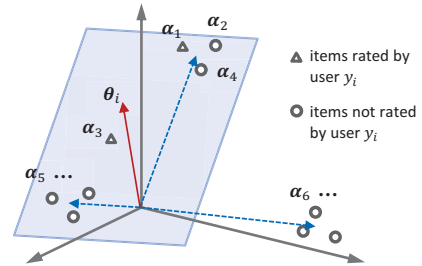


Figure 5: Illustration of the HSVD algorithm in 3d latent space. The dashed lines are cluster centers of the anchor items.

Proposition 3 shows that the lengths of the anchor latent vectors are longer than the non-anchor latent vectors, even if the anchor items are less popular. For a fixed θ_i , because the score $\hat{R}_{te}(i, k)$ is simply the inner product of θ_i and α_k , the anchor items of longer length would have bigger chance to be selected in the top- N recommendations. This would allow us to diversify the recommendation list and recommend more less-popular items.

4.1 Scalability

The main computational cost of our algorithm comes from the truncated SVD routine. Usually the numerical method for truncated SVD on matrix $\bar{\mathbf{X}}$ is in an iterative manner, and operated as an eigen-decomposition of the Hermitian matrix $\bar{\mathbf{X}}^\top \bar{\mathbf{X}}$. The running time of the latter problem depends on the complexity of each iteration and the total number of iterations. In each iteration the basic operation contains two matrix vector multiplications $\bar{\mathbf{X}}^\top \bar{\mathbf{X}} \mathbf{v}$ where \mathbf{v} is a dense vector. If the number of non-zero entries in $\bar{\mathbf{X}}$ (or the number of ratings in \mathbf{R}) is M , the complexity of each iteration would be $O(M)$. On the other hand, the total number of iterations to converge depends on the gap between the l largest singular value and the $l + 1$ largest singular value. A bigger gap leads to a smaller number of iterations [6, 11]. If the number of iterations is n_I , the overall complexity of the HSVD algorithm is $O(n_I M)$, which is linear to the number of ratings. The parameter l (dimensionality of the truncated SVD) is considered as a constant, which can be determined by cross-validation in practice.

5. RELATED WORK

In the introduction, we have mentioned the SVD [17], the SVD on non-empty entries [5], and the asymmetric SVD [14, 9] for recommender systems. Another matrix factorization technique that is widely used in recommender systems is the *non-negative matrix factorization* (NMF). In NMF, the rating matrix is approximately decomposed into two low-rank, sparse and non-negative matrices $\mathbf{R} \approx \mathbf{G}\mathbf{H}^\top$ such that $\|\mathbf{R} - \mathbf{G}\mathbf{H}^\top\|^2$ is minimized [7, 16]. The idea behind NMF is that the set of orthogonal latent factors controls the values in \mathbf{R} , and each latent factor takes one axis in the latent space. This idea is essentially similar to the HSVD since we have shown that the anchor items also exhibit an orthogonal structure in the latent space, although they are not aligned to the axes.

In order to compare our algorithm with the existing approaches, we adapt the SVD and the NMF into our setting, i.e. split the rating matrix into three parts as in Eq. (7) and use the sub-matrix $\mathbf{R}' = [\mathbf{R}_m^\top \mathbf{R}_{tr}^\top]^\top$ to compute the SVD or NMF, then evaluate the results on \mathbf{R}_{te} . The SVD and NMF algorithms in this setting are denoted as SVDR and NMFR.

To make a better baseline, we also apply the SVDR and NMFR on \mathbf{R}_m and predict with \mathbf{R}_{tr} in an asymmetric manner, i.e. follow the same procedures in Section 4 but replace $\bar{\mathbf{X}}_m$ with \mathbf{R}_m (or replace \mathbf{F} with the matrix \mathbf{H} in NMF). These methods are denoted as ASVDR and ANMFR, where the leading ‘‘A’’ means ‘‘asymmetric’’.

We use the routine implemented by TimelyDevelopment for computing the SVD on non-empty entries [18], which is denoted as SVDN. Note that SVDN is exactly the same as SVDR except for the SVD computation routine.

For the asymmetric SVD, we use the approach proposed in [14]. In this approach, a decomposition is first computed

by SVDN $\mathbf{R} \approx \mathbf{U}_i \mathbf{S}_i \mathbf{V}_i^\top$, where $\mathbf{V}_i = [\alpha_1 \cdots \alpha_n]^\top$. Then a user y_i is represented by the currently rated items: $\theta_i = \sum_{k, \mathbf{R}(i, k) > 0} \alpha_k$. Finally a score is computed by the rule $\hat{R}(i, k) = a_k + \theta_i^\top \alpha_k$ for each unrated item, where a_k is the average rating of item z_k . We call this method ASVDN.

6. EXPERIMENTAL RESULTS

In a top- N recommender system, the primary goal is to discover the unrated items that might be liked by a user [3]. We first test the algorithm’s ability to discover the unrated items, which is denoted as the scenario PREDICT-ALL. The rating matrix is split into three parts as in Eq. (7), and we randomly select \mathbf{R}_{tr} to have 5 ratings in each row (for each user). The remaining ratings in the Y_t part are assigned to \mathbf{R}_{te} . This setting simulates the situation where new users just come into the system with only a few ratings, which is denoted as PREDICT-ON-5. Similarly, we create another setting PREDICT-ON-20 where each row of \mathbf{R}_{tr} has 20 ratings to simulate the users with more available ratings. Thus in the scenario PREDICT-ALL, there are two settings, namely PREDICT-ON-5 and PREDICT-ON-20. Note that all the items in \mathbf{R}_{te} (rated or unrated) could be recommended in the top- N recommendation list. In all the experiments, \mathbf{R}_{tr} and \mathbf{R}_{te} are randomly selected in 5 different runs.

To test the results, we check if the set of recommended items coincides with the liked items (ratings higher than the median) for each user in the corresponding row of \mathbf{R}_{te} . Let \hat{R}_i^N and R_i denote these two sets for user y_i . We measure the average precision

$$\text{precision} = (1/|Y_t|) \sum_{y_i \in Y_t} \left(|\hat{R}_i^N \cap R_i| / |\hat{R}_i^N| \right). \quad (8)$$

	Lastfm		YahooMusic		BookCrossing	
SVDR	0.113	0.083	0.237	0.207	0.078	0.063
ASVDR	0.114	0.087	0.237	0.210	0.078	0.062
NMFR	0.114	0.090	0.218	0.189	0.073	0.054
ANMFR	0.110	0.089	0.211	0.190	0.071	0.056
SVDN	0.002	0.003	0.001	0.001	0.005	0.003
ASVDN	0.003	0.002	0.007	0.005	0.005	0.003
HSVD	0.180	0.158	0.258	0.227	0.075	0.065

Table 1: Average precisions in PREDICT-ALL. The left column of each dataset: PREDICT-ON-5, the right column: PREDICT-ON-20. The bold number indicates the method that performs significantly better than the others (p -value < 0.05 in paired t-test).

	Lastfm		YahooMusic		BookCrossing	
SVDR	0.198	0.140	0.490	0.418	0.743	0.744
ASVDR	0.198	0.139	0.496	0.435	0.743	0.745
NMFR	0.541	0.545	0.763	0.774	0.731	0.740
ANMFR	0.279	0.222	0.579	0.496	0.743	0.745
SVDN	0.177	0.244	0.249	0.224	0.389	0.165
ASVDN	0.537	0.531	0.755	0.783	0.696	0.714
HSVD	0.515	0.524	0.765	0.766	0.745	0.747

Table 2: Average precisions in RANK-CANDIDATES. (see the caption of Table 1 for more details)

In the scenario PREDICT-ALL, any item which is not in \mathbf{R}_{tr} can be considered in the recommendation list. But sometimes we already know a set of candidate items that a user

has suggested implicitly. For example, a user has viewed a list of pages of artists, but did not make any further actions like purchase a CD or download a track. In this case, we can limit our recommendations on the candidate items and only pick up the items that are truly liked by the user, which could be simulated by only recommending the rated items in R_{te} . We try to rank these items such that the liked ones are ranked higher. This scenario is denoted as RANK-CANDIDATES. The same measure (precision defined in Eq. (8)) is evaluated in RANK-CANDIDATES.

We have tested all the methods listed in Section 5 with the above settings on three datasets. Lastfm¹ contains the relations between users and artists collected from last.fm, and the ratings in Lastfm are the actual counts of plays with which a user has listened to an artist. YahooMusic² is also in the music domain, but includes relations between users and tracks. The ratings in YahooMusic are in the range of 1 to 100. The last dataset BookCrossing contains ratings from users to books [23]. The range of ratings in BookCrossing is from 1 to 10. In each dataset, we take a subset such that the induced graph of the hypergraph representation is connected. When the dataset contains a hypergraph with multiple connected components, we could simply apply the algorithm to each connected component respectively.

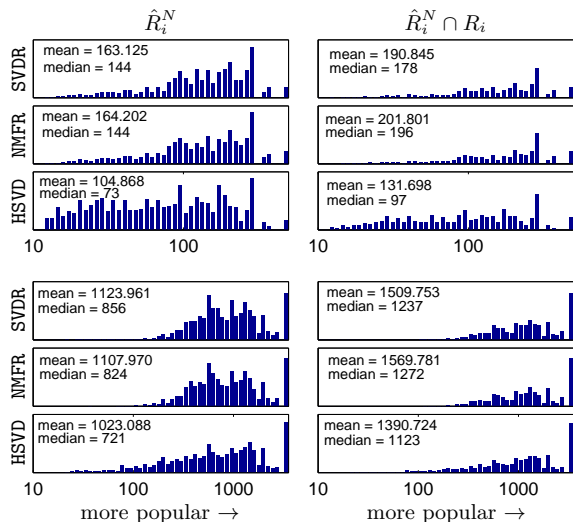


Figure 6: The distribution of popularity of the items (in the PREDICT-ALL scenario with $N = 20$). The x-axis (popularity) is in log-scale. The mean and median of the popularity of the items are also shown in each sub-figure. Top: BookCrossing. Bottom: YahooMusic.

Table 1 and Table 2 show the results of PREDICT-ALL and RANK-CANDIDATES respectively (with $N = 20$). In the first case, our proposed method performs best within the music domain, and shows results as good as the best methods with BookCrossing. The approaches ignoring the empty values (SVDN and ASVDN) does not perform very well. These methods are originally designed to minimize the root mean squared error on a testing set of known ratings, which is not fully compatible with the scenario PREDICT-ALL. In the

¹ Available at <http://mtg.upf.edu/node/1671>

² Yahoo! Webscope dataset ydata-ymusic-kddcup-2011-track2 http://labs.yahoo.com/Academic_Relations

RANK-CANDIDATES, no method performs significantly better than the others over all datasets. But the NMF-based method NMFR and two asymmetric methods ASVDN and HSVD are always as good as the best one.

In recommender systems, we are mainly interested in recommending non-popular (long-tail) items, because the users are usually already aware of the popular items from other sources [19]. For the purpose of diversifying recommendations, we measure the popularity distribution of the items in \hat{R}_i^N and $\hat{R}_i^N \cap R_i$, i.e. the set of all recommended items and the set of successfully recommended items (Figure 6). The popularity of an item z_k is the number of ratings in the k -th column of the rating matrix, i.e. $\text{deg}(z_k)$. We observe that HSVD produces more diverse recommendations with the same (or better) level of precision. Especially with BookCrossing, there are much more less-popular recommendations from HSVD compared to SVDR and NMFR.

Figure 7 explains why our approach produces more diverse recommendations. By using SVD without normalization, the lengths of the anchor items are not necessarily the longest ones. In fact, we can observe that there is a strong correlation between the popularity and the length of the latent vector with ASVDR, which would promote the popular items in the recommendations. In our approach HSVD, the longest latent vectors correspond to the anchor items of smaller popularity.

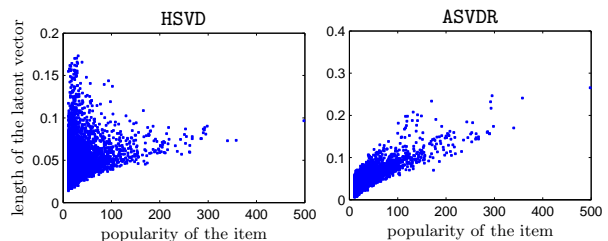


Figure 7: The length of the latent vectors and the popularity of all the items in BookCrossing.

7. CONCLUSION

In this work, we present a new perspective of SVD for constructing a latent space from the training data, which is justified by the transformation of normalized hypergraph cut. We show that the latent vectors representing the items in the latent space can be grouped into (approximately) orthogonal clusters, and the lengths of the vectors can be used to distinguish the anchor items and the non-anchor items. These properties are then utilized for making top- N recommendations in the proposed algorithm HSVD. Instead of explicitly constructing the clusters, the recommendations are generated by solving a linear system. We provide new explanations for the significantly better performance of the asymmetric SVD approaches and better diversity in top- N recommendations. Experiments conducted with three dataset on two domains also confirm our analysis.

Future works could include further studies of the latent space structures, and engineering methods that exploit such structures. For example, we can explicitly extract the anchor items to construct some clusters, and analyze the features of a non-anchor item by studying the underlying anchor clusters. It is also possible to include additional information,

e.g. contextual information [2] and social information [21], or rating bias corrections [10] into our scheme for further improvements.

8. REFERENCES

- [1] S. Agarwal, K. Branson, and S. Belongie. Higher order learning with graphs. In *ICML*, 2006.
- [2] L. Baltrunas, B. Ludwig, and F. Ricci. Matrix factorization techniques for context aware recommendation. In *Recsys*, 2011.
- [3] P. G. Campos, F. Díez, and M. Sánchez-Montañés. Towards a more realistic evaluation: testing the ability to predict future tastes of matrix factorization-based recommenders. In *Recsys*, 2011.
- [4] F. Chung. *Spectral Graph Theory*. American Mathematical Society, 1997.
- [5] S. Funk. Netflix update: Try this at home, <http://sifter.org/~simon/journal/20061211.html>. 2006.
- [6] G. H. Golub and C. F. Van Loan. *Matrix computations*. Johns Hopkins University Press, 1996.
- [7] P. Hoyer. Non-negative matrix factorization with sparseness constraints. *Journal of Machine Learning Research*, 5:1457–1469, 2004.
- [8] G. Karypis, R. Aggarwal, V. Kumar, and S. Shekhar. Multilevel hypergraph partitioning: Application in vlsi domain. In *Proc. of the 34th annual Design Automation Conference*, 1997.
- [9] Y. Koren. Factorization meets the neighborhood: a multifaceted collaborative filtering model. In *SIGKDD*, 2008.
- [10] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.
- [11] D. Mavroudis. Mind the eigen-gap, or how to accelerate semi-supervised spectral learning algorithms. In *IJCAI*, 2011.
- [12] B. K. Mohan, B. J. Keller, and N. Ramakrishnan. Scouts, promoters, and connectors: The roles of ratings in nearest-neighbor collaborative filtering. *ACM Transactions on the Web*, 1(2):8, 2007.
- [13] A. Y. Ng, M. I. Jordan, Y. Weiss, et al. On spectral clustering: Analysis and an algorithm. In *NIPS*, 2002.
- [14] A. Paterek. Improving regularized singular value decomposition for collaborative filtering. In *Proc. of KDD Cup and Workshop*, 2007.
- [15] L. Pu and B. Faltings. Hypergraph learning with hyperedge expansion. In *ECML-PKDD*, 2012.
- [16] S. Rendle and L. Schmidt-Thieme. Online-updating regularized kernel matrix factorization models for large-scale recommender systems. In *Recsys*, 2008.
- [17] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl. Application of dimensionality reduction in recommender system—a case study. In *WebKDD*, 2000.
- [18] TimelyDevelopment. <http://www.timelydevelopment.com/demos/NetflixPrize.aspx>. 2006.
- [19] S. Vargas and P. Castells. Rank and relevance in novelty and diversity metrics for recommender systems. In *Recsys*, 2011.
- [20] L. Wu, X. Ying, X. Wu, and Z. Zhou. Line orthogonality in adjacency eigenspace with application to community partition. In *IJCAI*, 2011.
- [21] Q. Yuan, L. Chen, and S. Zhao. Factorization vs. regularization: fusing heterogeneous social relationships in top-n recommendation. In *Recsys*, 2011.
- [22] D. Zhou, J. Huang, and B. Scholkopf. Learning with hypergraphs: Clustering, classification, and embedding. In *NIPS*, 2007.
- [23] C. Ziegler, S. McNeel, J. Konstan, and G. Lausen. Improving recommendation lists through topic diversification. In *WWW*, 2005.

APPENDIX

Proof for Proposition 1:

Firstly, it can be shown that $\alpha_k = \alpha_{k'}$ if $z_k, z_{k'} \in C_j$, because the k -th column and the k' -th column of $S_s \mathbf{F}^\top$ must be the same to obtain the same \bar{x}_k and $\bar{x}_{k'}$, which implies that $\alpha_k = \alpha_{k'}$. Secondly, we consider the full SVD of $\bar{\mathbf{X}}$ and list the rows of \mathbf{V} corresponding to items in C_j :

$$\mathbf{V}_{(j)}^\top = \begin{bmatrix} \beta_j & \beta_j & \cdots & \beta_j \\ \gamma_1 & \gamma_2 & \cdots & \gamma_{|C_j|} \end{bmatrix}. \quad (9)$$

It is always possible to find a linear combination of the first s columns of \mathbf{V} (denoted as $[\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_s] = \mathbf{F}$) such that $\mathbf{F} \mathbf{t}_j = [\mathbf{0}^\top \cdots \mathbf{1}^\top \cdots \mathbf{0}^\top]^\top$, where \mathbf{t}_j are the coefficients, and the 1's on the right hand side correspond to the items in C_j . Because all the columns in \mathbf{V} are orthogonal to each other, a column is also orthogonal to a linear combination of some other columns (e.g. $\mathbf{F} \mathbf{t}_j$). Thus the entries in each dimension of the vectors $\{\gamma_1, \dots, \gamma_{|C_j|}\}$ sum up to 0. In other words, $\sum_{k=\{1,2,\dots,|C_j|\}} \gamma_k = 0$, and $\mathbf{V}_{(j)}^\top \mathbf{1} = |C_j| [\beta_j^\top \mathbf{0}^\top]^\top$. Since all the rows in \mathbf{V} are also orthogonal to each other, we have $(\mathbf{V}_{(j)}^\top \mathbf{1})^\top (\mathbf{V}_{(j')}^\top \mathbf{1}) = 0$, which implies that $\beta_j^\top \beta_{j'} = 0$.

Proof for Proposition 3:

By Eq. (6) we know that $\beta'_k = \sum_{j \in I(z'_k)} \sqrt{\frac{\deg(z_j)}{\deg(z'_k)}} \beta_j$ (the residual \mathbf{r} is zero because the involved anchor items are non-overlapping). Thus $\|\beta'_k\|_2^2 = \sum_{j \in I(z'_k)} \frac{\deg(z_j)}{\deg(z'_k)} \|\beta_j\|_2^2 + \sum_{j, j' \in I(z'_k), j \neq j'} 2 \frac{\sqrt{\deg(z_j) \deg(z_{j'})}}{\deg(z'_k)} \beta_j^\top \beta_{j'}$. Since the anchor items are non-overlapping, we have $\sum_{j \in I(z'_k)} \deg(z_j) = \deg(z'_k)$. The above equation can be written as

$$\|\beta'_k\|_2^2 = b_0^2 + \sum_{j, j' \in I(z'_k), j \neq j'} 2 \frac{\sqrt{\deg(z_j) \deg(z_{j'})}}{\deg(z'_k)} \beta_j^\top \beta_{j'}. \quad (10)$$

Following the proof for Proposition 1, it can be shown that

$$\beta_j^\top \beta_{j'} = \frac{-1}{|C_j| |C_{j'}|} \left(\sum_{p \in I^{-1}(j)} \sqrt{\frac{\deg(z_j)}{\deg(z'_p)}} \beta'_p \right)^\top \left(\sum_{q \in I^{-1}(j')} \sqrt{\frac{\deg(z_{j'})}{\deg(z'_q)}} \beta'_q \right), \quad (11)$$

where $I^{-1}(j)$ is the set of non-anchor items that contain a fraction of an anchor item in cluster C_j . Since each β'_p is a linear combination of the β_j 's of the anchor items, and we know from Proposition 2 that $\beta_j^\top \beta_{j'} \approx 0$ for $\forall j \neq j'$, it is easy to show that $\beta_j^\top \beta_{j'} < 0$. Applying this result to Eq. (10), it concludes to $\|\beta'_k\|_2 < b_0$.