

# Matrix Co-Factorization on Compressed Sensing

Jiho Yoo<sup>1</sup> and Seungjin Choi<sup>1,2</sup>

<sup>1</sup> Department of Computer Science

<sup>2</sup> Division of IT Convergence Engineering

Pohang University of Science and Technology

San 31 Hyoja-dong, Nam-gu, Pohang 790-784, Korea

{zentasis, seungjin}@postech.ac.kr

## Abstract

In this paper we address the problem of matrix factorization on compressively-sampled measurements which are obtained by random projections. While this approach improves the scalability of matrix factorization, its performance is not satisfactory. We present a matrix co-factorization method where compressed measurements and a small number of uncompressed measurements are jointly decomposed, sharing a factor matrix. We evaluate the performance of three matrix factorization methods in terms of Cramér-Rao bounds, including: (1) matrix factorization on uncompressed data (MF); (2) matrix factorization on compressed data (CS-MF); (3) matrix co-factorization on compressed and uncompressed data (CS-MCF). Numerical experiments demonstrate that CS-MCF improves the performance of CS-MF, emphasizing the useful behavior of exploiting side information (a small number of uncompressed measurements).

## 1 Introduction

Matrix factorization is a decomposition method where a data matrix  $\mathbf{X} \in \mathbb{R}^{M \times N}$  is approximated as a product of two or more factor matrices which uncover the latent structure of the data. The 2-factor decomposition seeks latent factor matrices  $\mathbf{U} \in \mathbb{R}^{M \times K}$  and  $\mathbf{V} \in \mathbb{R}^{N \times K}$  which minimize the following objective function

$$\mathcal{J}_{\text{MF}} = \frac{1}{2} \|\mathbf{X} - \mathbf{UV}^\top\|_F^2, \quad (1)$$

where  $\|\cdot\|_F$  denotes the Frobenious norm. When columns in  $\mathbf{X}$  are treated as data points in  $M$ -dimensional space,  $\mathbf{U}$  and  $\mathbf{V}$  are referred to be as *basis matrix* and *encoding matrix*, respectively. Exemplary matrix factorization methods include principal component analysis (PCA) [Jolliffe, 2002], non-negative matrix factorization (NMF) [Lee and Seung, 1999], probabilistic matrix factorization (PMF) [Srebro *et al.*, 2005], to name a few. Matrix factorization has established itself as a powerful technique in various applications such as collaborative prediction [Koren *et al.*, 2009], document clustering [Xu *et al.*, 2003; Yoo and Choi, 2010], music transcription [Smaragdis and Brown, 2003], and brain wave analysis [Lee

and Choi, 2009]. In some of these applications such as collaborative prediction, the data matrix is sparse since most of entries are missing or unobserved. Thus, matrix factorization methods can be applied to a large scale problem without much extra care, since only observed elements are used for learning. On the other hand, data matrix is dense in some applications where spectrograms of music or brain wave data are analyzed. In such applications, it is not easy to handle large-scale data.

Compressed sensing is a new sensing/sampling paradigm which ensures near-optimal recovery of sparse signals from a small number of linear measurements [Donoho, 2006; Candes and Tao, 2006]. Compressed sensing is a promising approach to handling a large-scale dense data matrix since it is allowed to perform a processing (learning or data analysis) in the compressed domain (in which a compressively-sampled matrix is much smaller than the original input data matrix in size) rather than in the ambient space. Let  $\mathbf{Y} \in \mathbb{R}^{P \times N}$  be the compressively-sampled data matrix which is obtained by  $\mathbf{Y} = \Phi\mathbf{X}$ , in which  $\Phi \in \mathbb{R}^{P \times M}$  is a linear sensing matrix which is often given by a random projection. A direct application of the compressed sensing framework to matrix factorization, referred to as CS-MF, yields the following objective function

$$\mathcal{J}_{\text{CS-MF}} = \frac{1}{2} \|\mathbf{Y} - \Phi\mathbf{UV}^\top\|_F^2. \quad (2)$$

CS-MF determines  $\mathbf{U}$  and  $\mathbf{V}$  which minimize (2), improving the scalability since  $\mathbf{Y}$  and  $\Phi$  requires less space in storage, compared to  $\mathbf{X}$ . Recently, nonnegative matrix factorization (NMF) in the compressed domain, referred to as CS-NMF, was presented in [O'Grady and Rickard, 2008], where (2) is minimized with nonnegativity constraints imposed on factor matrices  $\mathbf{U}$  and  $\mathbf{V}$  as well as the data matrix  $\mathbf{X}$ . In contrast to the standard compressed sensing framework where  $\mathbf{U}$  is known in advance, CS-MF or CS-NMF estimates both  $\mathbf{U}$  and  $\mathbf{V}$  in the compressed-domain, which degrades the reconstruction performance.

Matrix co-factorization (MCF) is a technique that decompose two or more data matrices jointly, sharing some factor matrices in factorizations. In addition to the data matrix, side information matrices are also considered for the joint decomposition, including label information [Yu *et al.*, 2005; Lee *et al.*, 2010], link information [Zhu *et al.*, 2007], inter-subject variations [Lee and Choi, 2009], and spectrograms

of music [Yoo *et al.*, 2010]. MCF was also applied to the general entity-relationship models [Singh and Gordon, 2008; Yoo and Choi, 2009] to learn the characteristics of each entity from the relationships. In this paper we present a method of MCF on compressively-sampled data, referred to as CS-MCF, where we decompose the compressed data matrix and a small portion of uncompressed original data jointly in order to improve the performance of matrix factorization in the compressed domain. We compute the Cramér-Rao bound (CRB) for matrix factorization on uncompressed data (MF), CS-MF, and CS-MCF in the case of Gaussian likelihood, showing that CS-MCF indeed improves the performance over CS-MF in terms of the reconstruction quality. We consider alternating least squares (ALS) algorithms for MF, CS-MF, and CS-MCF. In addition, we also consider multiplicative updates, as in NMF, comparing NMF, CS-NMF (NMF in the compressed domain), and CS-NMCF (nonnegative matrix co-factorization in the compressed domain with small uncompressed information).

## 2 Matrix Co-Factorization on Compressively Sampled Data

Compressed sensing [Donoho, 2006; Candes and Tao, 2006] is a new framework which performs the measurement and compression simultaneously to reduce computation required in the measurement process. The fundamental theoretical results of compressed sensing is based on the sparsity of the representation of given data, which makes possible to reduce the required number of measurement less than the amount computed from the classical Shannon-Nyquist sampling theorem. Compressed sensing assumes that the basis which sparsely represents the given data is known, however, the appropriate representation basis for given data is usually not known in advance. In that cases we have to learn a representation matrix from a given data, however, the problem of learning representation basis from the compressed measurement has not been studied much.

Matrix factorization in the compressed domain considers a decomposition of the compressively-sampled data matrix  $\mathbf{Y} = \Phi \mathbf{X}$  where  $\Phi \in \mathbb{R}^{P \times M}$  is a random projection. Thus given  $\mathbf{X}$  and  $\Phi$ , CS-MF seeks  $\mathbf{Y} \approx \Phi \mathbf{U} \mathbf{V}^\top$ , determining factor matrices  $\mathbf{U}$  and  $\mathbf{V}$  which minimizes (2). In contrast to the compressed sensing framework where the basis matrix  $\mathbf{U}$  is known in advance, CS-MF should estimate both  $\mathbf{U}$  and  $\mathbf{V}$  given compressed data  $\mathbf{Y}$ . We believe that this explains the poor reconstruction performance of CS-MF in our experiments.

Now we explain matrix co-factorization in the compressed domain with exploiting partial uncompressed data. We partition the data matrix  $\mathbf{X} \in \mathbb{R}^{M \times N}$  into two sub-matrices  $\mathbf{X}^c \in \mathbb{R}^{M \times C}$  and  $\mathbf{X}^u \in \mathbb{R}^{M \times (N-C)}$ . Assuming that  $C > (N - C)$ , the sub-matrix  $\mathbf{X}^c$  is compressed by a random projection  $\Phi$ , leading to  $\mathbf{Y} = \Phi \mathbf{X}^c$  with the abuse of notation. CS-MCF considers a joint decomposition of the compressed data matrix  $\mathbf{Y}$  and partial uncompressed data  $\mathbf{X}^u$ , minimizing the following objective function:

$$\mathcal{J}_{\text{CS-MCF}} = \frac{1}{2} \|\mathbf{Y} - \Phi \mathbf{U} \mathbf{V}^\top\|_F^2 + \frac{\lambda}{2} \|\mathbf{X}^u - \mathbf{U} \mathbf{W}^\top\|_F^2, \quad (3)$$

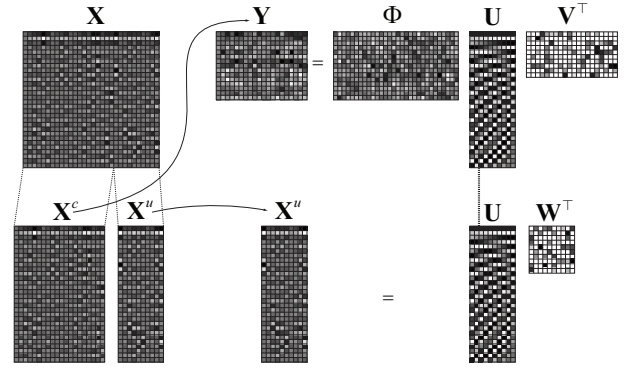


Figure 1: Pictorial illustration of CS-MCF. The input data matrix  $\mathbf{X}$  is partitioned into  $\mathbf{X}^c$  (to be compressed) and  $\mathbf{X}^u$  (partial uncompressed data). Then  $\mathbf{Y} = \Phi \mathbf{X}^c$  and  $\mathbf{X}^u$  are jointly decomposed, sharing the factor matrix  $\mathbf{U}$ , leading to  $\mathbf{Y} \approx \Phi \mathbf{U} \mathbf{V}^\top$  and  $\mathbf{X}^u \approx \mathbf{U} \mathbf{W}^\top$ .

where the factor matrix  $\mathbf{U} \in \mathbb{R}^{M \times K}$  is shared, capturing the common characteristics in representing both  $\mathbf{Y}$  and  $\mathbf{X}^u$ , and  $\mathbf{V} \in \mathbb{R}^{C \times K}$  and  $\mathbf{W} \in \mathbb{R}^{(N-C) \times K}$  capture individual characteristics. The parameter  $\lambda$  controls the balance between two decomposition. Pictorial illustration of CS-MCF is given in Fig. 1.

### 2.1 ALS Updates

ALS updates for  $\mathbf{U}$ ,  $\mathbf{V}$ ,  $\mathbf{W}$  are easily derived by solving

$$\begin{aligned} \frac{\partial \mathcal{J}_{\text{CS-MCF}}}{\partial \mathbf{U}} &= 0, \\ \frac{\partial \mathcal{J}_{\text{CS-MCF}}}{\partial \mathbf{V}} &= 0, \\ \frac{\partial \mathcal{J}_{\text{CS-MCF}}}{\partial \mathbf{W}} &= 0, \end{aligned}$$

for  $\mathbf{U}$ ,  $\mathbf{V}$ ,  $\mathbf{W}$ , respectively, and the corresponding gradients of (3) are calculated as

$$\begin{aligned} \frac{\partial \mathcal{J}_{\text{CS-MCF}}}{\partial \mathbf{U}} &= -\Phi^\top \mathbf{Y} \mathbf{V} + \Phi^\top \Phi \mathbf{U} \mathbf{V}^\top \mathbf{V} \\ &\quad - \lambda \mathbf{X}^u \mathbf{W} + \lambda \mathbf{U} \mathbf{W}^\top \mathbf{W}, \\ \frac{\partial \mathcal{J}_{\text{CS-MCF}}}{\partial \mathbf{V}} &= -\mathbf{Y}^\top \Phi \mathbf{U} + \mathbf{V} \mathbf{U}^\top \Phi^\top \Phi \mathbf{U}, \\ \frac{\partial \mathcal{J}_{\text{CS-MCF}}}{\partial \mathbf{W}} &= -\mathbf{X}^u \mathbf{U} + \mathbf{W} \mathbf{U}^\top \mathbf{U}, \end{aligned}$$

leading to ALS updates

$$\begin{aligned} \text{vec}(\mathbf{U}) &\leftarrow ((\mathbf{V}^\top \mathbf{V}) \otimes (\Phi^\top \Phi) + \lambda (\mathbf{W}^\top \mathbf{W}) \otimes \mathbf{I}_M)^{-1} \\ &\quad \text{vec}(\Phi^\top \mathbf{Y} \mathbf{V} + \lambda \mathbf{X}^u \mathbf{W}), \\ \mathbf{V} &\leftarrow \mathbf{Y}^\top \Phi \mathbf{U} (\mathbf{U}^\top \Phi^\top \Phi \mathbf{U})^{-1}, \\ \mathbf{W} &\leftarrow \mathbf{X}^u \mathbf{U} (\mathbf{U}^\top \mathbf{U})^{-1}, \end{aligned}$$

where  $\mathbf{I}_M$  is the identity matrix with size  $M \times M$ ,  $\text{vec}(\mathbf{U}) = [\mathbf{u}_1^\top, \dots, \mathbf{u}_K^\top]^\top$  is the 'vec' operator which creates a column vector from the matrix  $\mathbf{U}$  by stacking columns of  $\mathbf{U}$ , and  $\otimes$  denotes the Kronecker product (tensor product).

Table 1: Algorithm description for ALS updates and multiplicative updates.

Model	Updates	Model	Updates
MF	$U \leftarrow XV(V^\top V)^{-1}$ $V \leftarrow X^\top U(U^\top U)^{-1}$	NMF	$U \leftarrow U \odot \frac{XV}{UV^\top V}$ $V \leftarrow V \odot \frac{X^\top U}{VU^\top U}$
CS-MF	$\text{vec}(U) \leftarrow ((V^\top V) \otimes (\Phi^\top \Phi))^{-1} \text{vec}(\Phi^\top YV)$ $V \leftarrow Y^\top \Phi U(U^\top \Phi^\top \Phi U)^{-1}$	CS-NMF	$U \leftarrow U \odot \frac{\Phi^\top YV}{\Phi^\top \Phi U V^\top V}$ $V \leftarrow V \odot \frac{Y^\top \Phi U}{VU^\top \Phi^\top \Phi U}$
CS-MCF	$\text{vec}(U) \leftarrow ((V^\top V) \otimes (\Phi^\top \Phi) + \lambda(W^\top W) \otimes I_M)^{-1}$ $\text{vec}(\Phi^\top YV + \lambda X^u W)$ $V \leftarrow Y^\top \Phi U(U^\top \Phi^\top \Phi U)^{-1}$ $W \leftarrow X^{u^\top} U(U^\top U)^{-1}$	CS-NMCF	$U \leftarrow U \odot \frac{\Phi^\top YV + \lambda X^{u^\top} W}{\Phi^\top \Phi U V^\top V + \lambda U W^\top W}$ $V \leftarrow V \odot \frac{Y^\top \Phi U}{VU^\top \Phi^\top \Phi U}$ $W \leftarrow W \odot \frac{X^{u^\top} U}{WU^\top U}$

## 2.2 Multiplicative Updates

As in NMF, we derive multiplicative updates for  $U, V, W$  which iteratively minimize (3) with nonnegativity constraints imposed on factor matrices, given the nonnegative data matrix. Our derivation follows the technique used in [Yoo and Choi, 2008] where the gradient of an objective function  $\mathcal{J}$  is decomposed as

$$\frac{\partial \mathcal{J}}{\partial \Theta} = \left[ \frac{\partial \mathcal{J}}{\partial \Theta} \right]^+ - \left[ \frac{\partial \mathcal{J}}{\partial \Theta} \right]^-,$$

where  $\left[ \frac{\partial \mathcal{J}}{\partial \Theta} \right]^+ > 0$  and  $\left[ \frac{\partial \mathcal{J}}{\partial \Theta} \right]^- > 0$ , and then multiplicative updates for parameters  $\Theta$  are given by

$$\Theta \leftarrow \Theta \odot \frac{\left[ \frac{\partial \mathcal{J}}{\partial \Theta} \right]^+}{\left[ \frac{\partial \mathcal{J}}{\partial \Theta} \right]^+},$$

where  $\odot$  denotes Hadamard product (element-wise product) and the division is done in an element-wise manner. Applying this technique to the minimization of (3) yields

$$U \leftarrow U \odot \frac{\Phi^\top YV + \lambda X^u W}{\Phi^\top \Phi U V^\top V + \lambda U W^\top W},$$

$$V \leftarrow V \odot \frac{Y^\top \Phi U}{VU^\top \Phi^\top \Phi U},$$

$$W \leftarrow W \odot \frac{X^u U}{WU^\top U}.$$

ALS or multiplicative updating algorithms for MF and CS-MF are also derived in the same way, which are summarized in Table 1.

## 3 Cramér-Rao Bounds

We calculate the Cramér-Rao bound (CRB) to evaluate the performance of matrix factorization, CS-MF, and CS-MCF. CRB [Kay, 1993] provide the lower bound of the variance of unbiased estimators of a deterministic parameter in the following form,

$$\mathbb{E} \left\{ (\theta - \tilde{\theta})(\theta - \tilde{\theta})^\top \right\} \geq \mathcal{I}^{-1}, \quad (4)$$

where  $\theta$  is the estimated parameter of the model,  $\tilde{\theta}$  is the true value of the parameter,  $A \geq B$  denotes that  $A - B$  is the positive semi-definite matrix, and  $\mathcal{I}$  is the Fisher information matrix, each element of which is computed by

$$\mathcal{I}_{ij} = \mathbb{E} \left\{ -\frac{\partial^2 \log p(\mathbf{x}|\theta)}{\partial \theta_i \partial \theta_j} \right\}, \quad (5)$$

where  $p(\mathbf{x}|\theta)$  is the likelihood. In the case of matrix factorization, we have factor matrices  $U \in \mathbb{R}^{M \times K}$  and  $V \in \mathbb{R}^{N \times K}$  as the parameters for the model, so  $\theta$  becomes

$$\begin{aligned} \theta &= [\text{vec}(U)^\top, \text{vec}(V)^\top]^\top \\ &= [\mathbf{u}_1^\top, \mathbf{u}_2^\top, \dots, \mathbf{u}_M^\top, \mathbf{v}_1^\top, \mathbf{v}_2^\top, \dots, \mathbf{v}_N^\top]^\top, \end{aligned}$$

where  $\mathbf{u}_i$  and  $\mathbf{v}_j$  represent the column vector consist of the elements in the  $i$ -th and  $j$ -th row of  $U$  and  $V$ , respectively. Then, the Fisher information matrix has the following representation,

$$\begin{aligned} \mathcal{I} &= \begin{bmatrix} \mathcal{I}_U & \mathcal{I}_{UV} \\ \mathcal{I}_{UV}^\top & \mathcal{I}_V \end{bmatrix} \\ &= \begin{bmatrix} \mathcal{I}_{\mathbf{u}_1} & \dots & \mathcal{I}_{\mathbf{u}_1 \mathbf{u}_M} & \mathcal{I}_{\mathbf{u}_1 \mathbf{v}_1} & \dots & \mathcal{I}_{\mathbf{u}_1 \mathbf{v}_N} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \mathcal{I}_{\mathbf{u}_M \mathbf{u}_1} & \dots & \mathcal{I}_{\mathbf{u}_M} & \mathcal{I}_{\mathbf{u}_M \mathbf{v}_1} & \dots & \mathcal{I}_{\mathbf{u}_M \mathbf{v}_N} \\ \mathcal{I}_{\mathbf{v}_1 \mathbf{u}_1} & \dots & \mathcal{I}_{\mathbf{v}_1 \mathbf{u}_M} & \mathcal{I}_{\mathbf{v}_1} & \dots & \mathcal{I}_{\mathbf{v}_1 \mathbf{v}_N} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \mathcal{I}_{\mathbf{v}_N \mathbf{u}_1} & \dots & \mathcal{I}_{\mathbf{v}_N \mathbf{u}_M} & \mathcal{I}_{\mathbf{v}_N \mathbf{v}_1} & \dots & \mathcal{I}_{\mathbf{v}_N} \end{bmatrix}. \end{aligned}$$

where each part is computed for the corresponding parameter, for example,

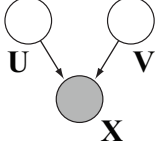
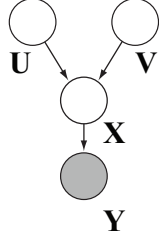
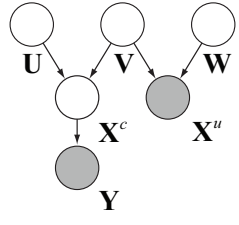
$$[\mathcal{I}_{\mathbf{u}_i}]_{k,k'} = \mathbb{E} \left\{ -\frac{\partial^2 \log p(\mathbf{X}|U, V)}{\partial u_{ik} \partial u_{ik'}} \right\}, \quad (6)$$

$$[\mathcal{I}_{\mathbf{u}_i \mathbf{u}_{i'}}]_{k,k'} = \mathbb{E} \left\{ -\frac{\partial^2 \log p(\mathbf{X}|U, V)}{\partial u_{ik} \partial u_{i'k'}} \right\}, \quad (7)$$

where  $[\mathbf{A}]_{i,j}$  represents the  $(i, j)$ -th element in the matrix  $\mathbf{A}$ . The Fisher information matrix of CS-MCF has a different parametrization with the matrix  $U, V$  and  $W$ , which is

$$\mathcal{I} = \begin{bmatrix} \mathcal{I}_U & \mathcal{I}_{UV} & \mathcal{I}_{UW} \\ \mathcal{I}_{UV}^\top & \mathcal{I}_V & \mathcal{I}_{VW} \\ \mathcal{I}_{UW}^\top & \mathcal{I}_{VW}^\top & \mathcal{I}_W \end{bmatrix},$$

Table 2: Graphical model, likelihood, Fisher information matrix (FIM) and Cramér-Rao lower bound for matrix factorization (MF), matrix factorization for compressed measurements (CS-MF), and matrix co-factorization for compressed measurements (CS-MCF).

	MF	CS-MF	CS-MCF
Graphical Model			
Likelihood	$p(\mathbf{X} \mathbf{U}, \mathbf{V}) = \prod_j \mathcal{N}(\mathbf{x}_j \mathbf{U}\mathbf{v}_j, \Sigma)$	$p(\mathbf{Y} \mathbf{U}, \mathbf{V}) = \prod_j \mathcal{N}(\mathbf{y}_j \Phi\mathbf{U}\mathbf{v}_j, \Sigma)$	$p(\mathbf{Y}, \mathbf{X}^u \mathbf{U}, \mathbf{V}, \mathbf{W}) = \prod_j \mathcal{N}(\mathbf{y}_j \Phi\mathbf{U}\mathbf{v}_j, \Sigma_b) \prod_k \mathcal{N}(\mathbf{x}_k^u \mathbf{U}\mathbf{w}_k, \Sigma_a)$
FIM			
$\mathcal{I}_U$	$(\mathbf{V}^\top \Sigma^{-1} \mathbf{V}) \otimes \mathbf{I}_M$	$(\mathbf{V}^\top \mathbf{V}) \otimes (\Phi^\top \Sigma^{-1} \Phi)$	$(\mathbf{V}^\top \mathbf{V}) \otimes (\Phi^\top \Sigma_a^{-1} \Phi) + (\mathbf{W}^\top \Sigma_b^{-1} \mathbf{W}) \otimes \mathbf{I}_M$
$\mathcal{I}_V$	$(\mathbf{U}^\top \Sigma^{-1} \mathbf{U}) \otimes \mathbf{I}_N$	$(\mathbf{U}^\top \Phi^\top \Sigma^{-1} \Phi \mathbf{U}) \otimes \mathbf{I}_N$	$(\mathbf{U}^\top \Phi^\top \Sigma_a^{-1} \Phi \mathbf{U}) \otimes \mathbf{I}_C$
$\mathcal{I}_{UV}$	$\text{vec}(\mathbf{U}\Sigma^{-1})\text{vec}(\mathbf{V})^\top$	$\text{vec}(\mathbf{U}\Phi^\top \Sigma^{-1} \Phi)\text{vec}(\mathbf{V})^\top$	$\text{vec}(\mathbf{U}^\top \Phi^\top \Sigma_a^{-1} \Phi)\text{vec}(\mathbf{V})^\top$
$\mathcal{I}_W$			$(\mathbf{U}^\top \Sigma_b^{-1} \mathbf{U}) \otimes \mathbf{I}_{N-C}$
$\mathcal{I}_{UW}$			$\text{vec}(\mathbf{U}\Sigma_b^{-1})\text{vec}(\mathbf{W})^\top$
$\mathcal{I}_{VW}$			$\mathbf{0}$
CRB ( $\mathbf{U}$ known)	$[\mathcal{I}_V]^{-1}$	$[\mathcal{I}_V]^{-1}$	$\begin{bmatrix} \mathcal{I}_V & \mathcal{I}_{VW} \\ \mathcal{I}_{VW}^\top & \mathcal{I}_W \end{bmatrix}^{-1}$
CRB ( $\mathbf{U}$ unknown)	$\begin{bmatrix} \mathcal{I}_U & \mathcal{I}_{UV} \\ \mathcal{I}_{UV}^\top & \mathcal{I}_V \end{bmatrix}^{-1}$	$\begin{bmatrix} \mathcal{I}_U & \mathcal{I}_{UV} \\ \mathcal{I}_{UV}^\top & \mathcal{I}_V \end{bmatrix}^{-1}$	$\begin{bmatrix} \mathcal{I}_U & \mathcal{I}_{UV} & \mathcal{I}_{UW} \\ \mathcal{I}_{UV}^\top & \mathcal{I}_V & \mathcal{I}_{VW} \\ \mathcal{I}_{UW}^\top & \mathcal{I}_{VW}^\top & \mathcal{I}_W \end{bmatrix}^{-1}$

and each element of the parts is computed in the similar way to (6) and (7).

The CRB for matrix factorization, CS-MF and CS-MCF are computed in the following subsections. We use the maximum likelihood estimator (MLE), which is known to be asymptotically unbiased and efficient under mild regularity conditions [Kay, 1993]. The analysis is based on these asymptotic properties, assuming that there are sufficient number of samples. The models, computed Fisher information matrices and the form of CRB are summarized in Table 2.

Since the solution of the matrix factorization is not uniquely determined, the direct comparison between the estimated parameters and the ground-truth values is meaningless. Instead of directly comparing the parameters themselves, we use the expected reconstruction error

$$\begin{aligned} \mathcal{E}_{ij} &= \mathbb{E} \left\{ (\mathbf{X}_{ij} - \widetilde{\mathbf{X}}_{ij})^2 \right\} \\ &= \mathbb{E} \left\{ (\mathbf{u}_i^\top \mathbf{v}_j - \widetilde{\mathbf{u}}_i^\top \widetilde{\mathbf{v}}_j)^2 \right\}, \end{aligned}$$

where  $\widetilde{\mathbf{X}}_{ij}$ ,  $\widetilde{\mathbf{u}}_i$ , and  $\widetilde{\mathbf{v}}_j$  are the ground-truth values of the parameters, and  $\mathbf{X}_{ij}$  is the element of the reconstructed input

matrix from the estimated parameters  $\mathbf{u}_i$  and  $\mathbf{v}_j$ . Although the factor matrices is not uniquely determined, the reconstruction error  $\mathcal{E}_{ij}$  is the same for all the possible decompositions. The lower bound for the reconstruction error is represented by using the variances of the estimators, such as

$$\begin{aligned} \mathcal{E}_{ij} &= \mathbb{E} \left\{ (\mathbf{u}_i^\top \mathbf{v}_j - \widetilde{\mathbf{u}}_i^\top \mathbf{v}_j + \widetilde{\mathbf{u}}_i^\top \mathbf{v}_j - \widetilde{\mathbf{u}}_i^\top \widetilde{\mathbf{v}}_j)^2 \right\} \\ &= \mathbb{E} \left\{ \mathbf{v}_j^\top (\mathbf{u}_i - \widetilde{\mathbf{u}}_i) (\mathbf{u}_i - \widetilde{\mathbf{u}}_i)^\top \mathbf{v}_j \right\} \\ &\quad + \mathbb{E} \left\{ \widetilde{\mathbf{u}}_i^\top (\mathbf{v}_j - \widetilde{\mathbf{v}}_j) (\mathbf{v}_j - \widetilde{\mathbf{v}}_j)^\top \widetilde{\mathbf{u}}_i \right\} \\ &\quad + 2\mathbb{E} \left\{ \mathbf{v}_j^\top (\mathbf{u}_i - \widetilde{\mathbf{u}}_i) (\mathbf{v}_j - \widetilde{\mathbf{v}}_j)^\top \widetilde{\mathbf{u}}_i \right\}. \end{aligned}$$

If we assume that the estimator is unbiased, the last term of the above equation becomes zero, and the remaining terms are lower-bounded by the computed CRB,

$$\mathcal{E}_{ij} \geq \mathbb{E} \left\{ \mathbf{v}_j^\top [\mathcal{I}^{-1}]_{\mathbf{u}_i} \mathbf{v}_j \right\} + \widetilde{\mathbf{u}}_i^\top [\mathcal{I}^{-1}]_{\mathbf{v}_j} \widetilde{\mathbf{u}}_i \quad (8)$$

where  $[\mathcal{I}^{-1}]_{\mathbf{u}_i}$  represents the part of the inverse Fisher information matrix corresponding to the position of  $\mathcal{I}_{\mathbf{u}_i}$ . Then,

the error bound becomes

$$\begin{aligned} \mathcal{E}_{ij} \geq & \tilde{\mathbf{v}}_j^\top [\mathcal{I}^{-1}]_{\mathbf{u}_i} \tilde{\mathbf{v}}_j + \text{Tr} \left\{ [\mathcal{I}^{-1}]_{\mathbf{u}_i} [\mathcal{I}^{-1}]_{\mathbf{v}_j} \right\} \\ & + \tilde{\mathbf{u}}_i^\top [\mathcal{I}^{-1}]_{\mathbf{v}_j} \tilde{\mathbf{u}}_i. \end{aligned} \quad (9)$$

Note that if we assume that the basis matrix  $\mathbf{U}$  is known in advance, we only have to deal with the parameter  $\mathbf{V}$ . With the computation similar to the above, the lower bound of the reconstruction error is computed in the following form,

$$\mathcal{E}_{ij} \geq \tilde{\mathbf{u}}_i^\top [\mathcal{I}^{-1}]_{\mathbf{v}_j} \tilde{\mathbf{u}}_i, \quad (10)$$

which is smaller than the case when the basis is unknown in the amount of the terms corresponding to  $\mathbf{U}$ , which is  $\tilde{\mathbf{v}}_j^\top [\mathcal{I}^{-1}]_{\mathbf{u}_i} \tilde{\mathbf{v}}_j + \text{Tr} \left\{ [\mathcal{I}^{-1}]_{\mathbf{u}_i} [\mathcal{I}^{-1}]_{\mathbf{v}_j} \right\}$ .

### 3.1 Cramér-Rao Bounds for Matrix Factorization

The Fisher information matrix is computed from the likelihood of the model. If we assume that the noise is additive and distributed as the Gaussian distribution, the model is written by

$$\mathbf{X} = \mathbf{U}\mathbf{V}^\top + \mathbf{N},$$

where  $\mathbf{N}_{ij} \sim \mathcal{N}(0, \sigma^2)$ , and  $\mathcal{N}(\mu, \sigma^2)$  is the Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ . Then, the likelihood of the model becomes

$$p(\mathbf{X}|\mathbf{U}, \mathbf{V}) = \prod_j \mathcal{N}(\mathbf{x}_j | \mathbf{U}\mathbf{v}_j, \sigma^2 \mathbf{I}_K).$$

#### Calculation of $\mathcal{I}_V$

The derivative of the log-likelihood with respect to  $\mathbf{v}_j$  is computed by

$$\frac{\partial \log p(\mathbf{X}|\mathbf{U}, \mathbf{V})}{\partial \mathbf{v}_j} = \mathbf{U}^\top \Sigma^{-1} \mathbf{x}_j - \mathbf{U}^\top \Sigma^{-1} \mathbf{U} \mathbf{v}_j,$$

where  $\Sigma = \sigma^2 \mathbf{I}_M$ . If we differentiate above again with respect to  $\mathbf{v}_{j'}$ , where  $j' \neq j$ , the result becomes zero. If  $j' = j$ , then,

$$\frac{\partial^2 \log p(\mathbf{X}|\mathbf{U}, \mathbf{V})}{\partial \mathbf{v}_j \partial \mathbf{v}_j} = -\mathbf{U}^\top \Sigma^{-1} \mathbf{U}.$$

Therefore, the part of Fisher information matrix for  $\mathbf{V}$  is computed by

$$\mathcal{I}_V = \mathbb{E} \left\{ (\mathbf{U}^\top \Sigma^{-1} \mathbf{U}) \otimes \mathbf{I}_N \right\} = (\mathbf{U}^\top \Sigma^{-1} \mathbf{U}) \otimes \mathbf{I}_N.$$

#### Calculation of $\mathcal{I}_U$

In this case, we use the likelihood of the model having the following form,

$$p(\mathbf{X}|\mathbf{U}, \mathbf{V}) = \prod_i \mathcal{N}(\mathbf{x}_i | \mathbf{V} \mathbf{u}_i, \sigma^2 \mathbf{I}_K),$$

where  $\mathbf{x}_i$  represents the  $i$ -th row of the input matrix  $\mathbf{X}$ . The first derivative with respect to  $\mathbf{u}_i$  is computed as

$$\frac{\partial \log p(\mathbf{X}|\mathbf{U}, \mathbf{V})}{\partial \mathbf{u}_i} = \mathbf{V}^\top \Sigma^{-1} \mathbf{x}_i - \mathbf{V}^\top \Sigma^{-1} \mathbf{V} \mathbf{u}_i. \quad (11)$$

Again, the second derivative becomes zero if  $i' \neq i$ , and if  $i' = i$ ,

$$\frac{\partial^2 \log p(\mathbf{X}|\mathbf{U}, \mathbf{V})}{\partial \mathbf{u}_i \partial \mathbf{u}_i} = -\mathbf{V}^\top \Sigma^{-1} \mathbf{V}.$$

Therefore, Fisher information matrix for  $\mathbf{U}$  is written by

$$\mathcal{I}_U = \mathbb{E} \left\{ (\mathbf{V}^\top \Sigma^{-1} \mathbf{V}) \otimes \mathbf{I}_M \right\} = (\mathbf{V}^\top \Sigma^{-1} \mathbf{V}) \otimes \mathbf{I}_M.$$

#### Calculation of $\mathcal{I}_{UV}$

To compute  $\mathcal{I}_{UV}$ , we differentiate (11) with respect to  $\mathbf{v}_j$ , which becomes

$$\frac{\partial^2 \log p(\mathbf{X}|\mathbf{U}, \mathbf{V})}{\partial \mathbf{u}_i \partial \mathbf{v}_j} = \frac{1}{\sigma^2} (\mathbf{x}_{ij} \mathbf{I}_K - \mathbf{u}_i^\top \mathbf{v}_j \mathbf{I}_K - \mathbf{v}_j \mathbf{u}_i^\top).$$

Since  $\mathbb{E} \{ \mathbf{x}_{ij} \mathbf{I}_K \} = \mathbf{u}_i^\top \mathbf{v}_j \mathbf{I}_K$ , the first two terms of the above equation vanishes when we take the expectation. Therefore, the Fisher information matrix for  $\mathbf{U}$  and  $\mathbf{V}$  is written by

$$\mathcal{I}_{UV} = \text{vec}(\mathbf{U} \Sigma^{-1}) \text{vec}(\mathbf{V})^\top.$$

### 3.2 Cramér-Rao Bound for CS-MF

We model the CS-MF with the following two-step procedure,

$$\begin{aligned} \mathbf{X} &= \mathbf{U}\mathbf{V}^\top + \mathbf{N}_a, \\ \mathbf{Y} &= \Phi \mathbf{X} + \mathbf{N}_b, \end{aligned}$$

where  $[\mathbf{N}_a]_{i,j} \sim \mathcal{N}(0, \sigma_a^2)$  is the representation error and  $[\mathbf{N}_b]_{i,j} \sim \mathcal{N}(0, \sigma_b^2)$  is the measurement error. The likelihood is written by

$$\begin{aligned} p(\mathbf{Y}|\mathbf{U}, \mathbf{V}) &= \int p(\mathbf{Y}|\mathbf{X}) p(\mathbf{X}|\mathbf{U}, \mathbf{V}) d\mathbf{X} \\ &= \prod_j \mathcal{N}(\mathbf{y}_j | \Phi \mathbf{U} \mathbf{v}_j, \Sigma), \end{aligned}$$

where  $\Sigma = \sigma_a^2 \Phi \Phi^\top + \sigma_b^2 \mathbf{I}_P$ , where  $\Phi \in \mathbb{R}^{P \times M}$ .

#### Calculation of $\mathcal{I}_V$

The derivative of the log-likelihood with respect to  $\mathbf{v}_j$  is computed by

$$\frac{\partial \log p(\mathbf{Y}|\mathbf{U}, \mathbf{V})}{\partial \mathbf{v}_j} = \mathbf{U}^\top \Phi^\top \Sigma^{-1} \mathbf{y}_j - \mathbf{U}^\top \Phi^\top \Sigma^{-1} \Phi \mathbf{U} \mathbf{v}_j.$$

Differentiate again with  $\mathbf{v}_{j'}$  becomes zero when  $j' \neq j$ , and if  $j' = j$ , then

$$\frac{\partial^2 \log p(\mathbf{Y}|\mathbf{U}, \mathbf{V})}{\partial \mathbf{v}_j \partial \mathbf{v}_j} = -\mathbf{U}^\top \Phi^\top \Sigma^{-1} \Phi \mathbf{U}.$$

Therefore, the Fisher information matrix for  $\mathbf{V}$  is computed by

$$\begin{aligned} \mathcal{I}_V &= \mathbb{E} \left\{ \mathbf{U}^\top \Phi^\top \Sigma^{-1} \Phi \mathbf{U} \otimes \mathbf{I}_N \right\} \\ &= \mathbf{U}^\top \Phi^\top \Sigma^{-1} \Phi \mathbf{U} \otimes \mathbf{I}_N. \end{aligned}$$

### Calculation of $\mathcal{I}_U$

The derivative with respect to  $\mathbf{u}_i$  is written by

$$\begin{aligned} & \frac{\partial \log p(\mathbf{Y}|\mathbf{U}, \mathbf{V})}{\partial \mathbf{u}_i} \\ &= \sum_j \left( \mathbf{v}_j \mathbf{y}_j^\top \Sigma^{-1} \phi_i - \mathbf{v}_j \phi_i^\top \Sigma^{-1} \Phi \mathbf{U} \mathbf{v}_j \right). \end{aligned}$$

The derivative of above equation with respect to  $\mathbf{u}_{i'}$  is simplified by

$$\frac{\partial^2 \log p(\mathbf{Y}|\mathbf{U}, \mathbf{V})}{\partial \mathbf{u}_i \partial \mathbf{u}_{i'}} = -(\Phi^\top \Sigma^{-1} \Phi)_{ii'} \mathbf{V}^\top \mathbf{V},$$

which leads the Fisher information matrix for  $\mathbf{U}$  to be

$$\mathcal{I}_U = (\mathbf{V}^\top \mathbf{V}) \otimes (\Phi^\top \Sigma^{-1} \Phi).$$

### Calculation of $\mathcal{I}_{UV}$

If we differentiate (12) with respect to  $\mathbf{v}_j$ , we obtain

$$\begin{aligned} & \frac{\partial^2 \log p(\mathbf{Y}|\mathbf{U}, \mathbf{V})}{\partial \mathbf{u}_{i^*} \partial \mathbf{v}_{j^*}} = \phi_{i^*}^\top \Sigma^{-1} \mathbf{y}_{j^*} \mathbf{I}_K \\ & - \phi_{i^*}^\top \Sigma^{-1} \Phi \mathbf{U} \mathbf{v}_{j^*} \mathbf{I}_K - \mathbf{U}^\top \Phi^\top \Sigma^{-1} \phi_{i^*} \mathbf{v}_{j^*}^\top. \end{aligned}$$

The first two terms disappear when we take the expectation, so the Fisher information matrix for  $\mathbf{U}$  and  $\mathbf{V}$  becomes

$$\mathcal{I}_{UV} = \text{vec}(\mathbf{U}^\top \Phi^\top \Sigma^{-1} \Phi) \text{vec}(\mathbf{V})^\top.$$

The theoretical bound of CS-MF is larger than that of matrix factorization. The main difference of the computed Fisher information matrix is from the change of the term  $\Sigma^{-1}$  of matrix factorization to the term  $\Phi^\top \Sigma^{-1} \Phi$ . Note that the  $\Sigma^{-1}$  in the matrix factorization is

$$\Sigma^{-1} = \sigma_a^{-2} \mathbf{I}_M. \quad (12)$$

However,  $\Phi^\top \Sigma^{-1} \Phi$  in the CS-MF is written by

$$\Phi^\top \Sigma^{-1} \Phi = \Phi^\top (\sigma_a^2 \Phi \Phi^\top + \sigma_b^2 \mathbf{I}_P)^{-1} \Phi.$$

This is usually smaller than the corresponding term in the matrix factorization. For example, if we assume that  $\sigma_b^2 = 0$  and each element of sensing matrix is drawn independently from  $\mathcal{N}(0, 1)$ , we approximate the term using  $\Phi^\top \Phi \approx P \mathbf{I}_M$  and  $\Phi \Phi^\top \approx M \mathbf{I}_P$ , such that

$$\begin{aligned} \Phi^\top \Sigma^{-1} \Phi &\approx \Phi^\top (\sigma_a^2 M \mathbf{I}_P)^{-1} \Phi \\ &\approx \frac{1}{M} \sigma_a^{-2} \Phi^\top \Phi \\ &\approx \frac{P}{M} \sigma_a^{-2} \mathbf{I}_M, \end{aligned}$$

which is  $P/M$  times smaller than that part of matrix factorization (12). Since the CRB is calculated as the inverse of the Fisher information matrix, CS-MF has larger CRB than that of the matrix factorization.

### 3.3 Cramér-Rao Bound for CS-MCF

The relations in CS-MCF are modeled as in the previous section, such that

$$\begin{aligned} \mathbf{X}^c &= \mathbf{U} \mathbf{V}^\top + \mathbf{N}_a, \\ \mathbf{Y} &= \Phi \mathbf{X}^c + \mathbf{N}_b, \\ \mathbf{X}^u &= \mathbf{U} \mathbf{W}^\top + \mathbf{N}_c, \end{aligned}$$

where  $[\mathbf{N}_a]_{i,j} \sim \mathcal{N}(0, \sigma_a^2)$ ,  $[\mathbf{N}_c]_{i,j} \sim \mathcal{N}(0, \sigma_a^2)$  and  $[\mathbf{N}_b]_{i,j} \sim \mathcal{N}(0, \sigma_b^2)$ . In this case we have three parameters  $\mathbf{U}$ ,  $\mathbf{V}$ , and  $\mathbf{W}$ , so the Fisher information matrix is computed as

$$\mathcal{I} = \begin{bmatrix} \mathcal{I}_U & \mathcal{I}_{UV} & \mathcal{I}_{UW} \\ \mathcal{I}_{UV}^\top & \mathcal{I}_V & \mathcal{I}_{VW} \\ \mathcal{I}_{UW}^\top & \mathcal{I}_{VW}^\top & \mathcal{I}_W \end{bmatrix}.$$

The likelihood is

$$\begin{aligned} & p(\mathbf{Y}, \mathbf{X}^u | \mathbf{U}, \mathbf{V}, \mathbf{W}) \\ &= p(\mathbf{Y} | \mathbf{U}, \mathbf{V}) p(\mathbf{X}^u | \mathbf{U}, \mathbf{W}) \\ &= \prod_j \mathcal{N}(\mathbf{y}_j | \Phi \mathbf{U} \mathbf{v}_j, \Sigma_b) \prod_l \mathcal{N}(\mathbf{x}_l^u | \mathbf{U} \mathbf{w}_l, \Sigma_a), \end{aligned}$$

where  $\Sigma_b = \sigma_a^2 \Phi \Phi^\top + \sigma_b^2 \mathbf{I}_P$  and  $\Sigma_a = \sigma_a^2 \mathbf{I}_M$ .

### Calculation of $\mathcal{I}_U$

Differentiating the log-likelihood with respect to  $\mathbf{u}_i$  leads

$$\begin{aligned} & \frac{\partial \log p(\mathbf{Y}, \mathbf{X}^u | \mathbf{U}, \mathbf{V}, \mathbf{W})}{\partial \mathbf{u}_i} \\ &= \sum_j \left( \mathbf{v}_j \mathbf{y}_j^\top \Sigma_b^{-1} \phi_i - \mathbf{v}_j \phi_i^\top \Sigma_b^{-1} \Phi \mathbf{U} \mathbf{v}_j \right) \\ & \quad + \mathbf{W}^\top \Sigma_a^{-1} \mathbf{x}_i - \mathbf{W}^\top \Sigma_a^{-1} \mathbf{W} \mathbf{u}_i. \end{aligned} \quad (13)$$

Differentiating again with respect to  $\mathbf{u}_{i'}$  yields

$$\begin{aligned} & \frac{\partial^2 \log p(\mathbf{Y}, \mathbf{X}^u | \mathbf{U}, \mathbf{V}, \mathbf{W})}{\partial \mathbf{u}_i \partial \mathbf{u}_{i'}} \\ &= -(\Phi^\top \Sigma_b^{-1} \Phi)_{ii'} \mathbf{V}^\top \mathbf{V} - \mathbf{W}^\top \Sigma_a^{-1} \mathbf{W}, \end{aligned}$$

where  $i' = i$ , and yields

$$\frac{\partial^2 \log p(\mathbf{Y}, \mathbf{X}^u | \mathbf{U}, \mathbf{V}, \mathbf{W})}{\partial \mathbf{u}_i \partial \mathbf{u}_{i'}} = -(\Phi^\top \Sigma_b^{-1} \Phi)_{ii'} \mathbf{V}^\top \mathbf{V},$$

where  $i' \neq i$ . Therefore, the Fisher information matrix for  $\mathbf{U}$  becomes

$$\mathcal{I}_U = (\mathbf{V}^\top \mathbf{V}) \otimes (\Phi^\top \Sigma_b^{-1} \Phi) + (\mathbf{W}^\top \Sigma_a^{-1} \mathbf{W}) \otimes \mathbf{I}_M.$$

### Calculation of $\mathcal{I}_V$ , $\mathcal{I}_W$ , $\mathcal{I}_{UV}$ , and $\mathcal{I}_{VW}$

$\mathcal{I}_V$ ,  $\mathcal{I}_W$ ,  $\mathcal{I}_{UV}$ , and  $\mathcal{I}_{VW}$  are computed in the similar way as in matrix factorization and CS-MF, in the following forms,

$$\begin{aligned} \mathcal{I}_V &= (\mathbf{U}^\top \Phi^\top \Sigma_b^{-1} \Phi \mathbf{U}) \otimes \mathbf{I}_C, \\ \mathcal{I}_W &= (\mathbf{U}^\top \Sigma_a^{-1} \mathbf{U}) \otimes \mathbf{I}_{N-C}, \\ \mathcal{I}_{UV} &= \text{vec}(\mathbf{U}^\top \Phi^\top \Sigma_b^{-1} \Phi) \text{vec}(\mathbf{V})^\top, \\ \mathcal{I}_{UW} &= \text{vec}(\mathbf{U} \Sigma_a^{-1}) \text{vec}(\mathbf{W})^\top. \end{aligned}$$

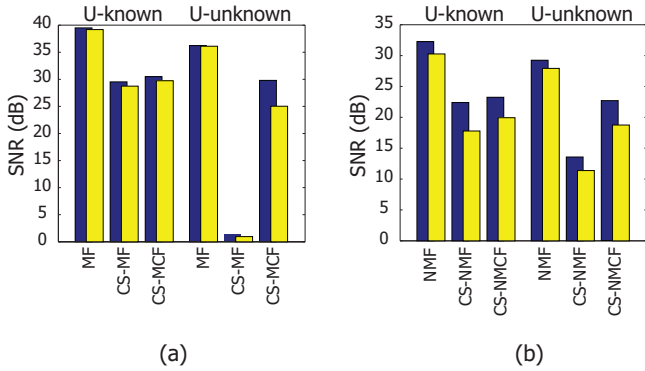


Figure 2: SNR computed from the CRB (dark) and SNR achieved from the algorithms (light), for the (a) ALS algorithms with Gaussian data (b) Multiplicative algorithms with nonnegative data. In each graph, left three are for the case of known representation  $U$ , and right three are for the case of unknown  $U$ .

#### Calculation of $\mathcal{I}_{VW}$

Fisher information matrix for  $V$  and  $W$ ,  $\mathcal{I}_{VW}$ , becomes zero matrix of dimension  $CK \times (N - C)K$  because the term  $W$  banishes in the derivative with respect to  $v_j$ .

The Fisher information matrix of CS-MCF is the combination of the Fisher information matrix of matrix factorization and Fisher information matrix of CS-MF. The parameter  $V$  in matrix factorization and CS-MF is divided into  $V$  and  $W$  in CS-MCF. The  $\mathcal{I}_V$  and  $\mathcal{I}_{UV}$  in CS-MCF follows the form of the corresponding parts of CS-MF, but the  $\mathcal{I}_W$  and  $\mathcal{I}_{UW}$  follows the form of the corresponding parts of matrix factorization. Moreover,  $\mathcal{I}_U$  is the direct combination of the  $\mathcal{I}_U$  of matrix factorization and  $\mathcal{I}_U$  of CS-MF. In the previous section, we showed that in general the CRB of matrix factorization is smaller than the CRB of CS-MF, so this kind of combination leads smaller CRB than that of CS-MF.

## 4 Numerical Experiments

For the numerical experiments, we computed the bounds of the reconstruction performances using the equation (9), and compared them with the actual reconstruction performances of the algorithms listed in the Table 1. We also computed the bounds for the cases where we assume that the parameter  $U$  is known in advance (10). In that cases, the actual reconstruction performances were measured based on the modified algorithms, where the parameter  $U$  is fixed with the ground-truth values.

To compute the CRB, we have to know the ground-truth value of the parameters  $U$ ,  $V$ , and  $W$ , which makes it difficult to use the real-world data where the ground-truth factor matrices are unknown. Therefore, we used synthetic data for the experiments. The factor matrices  $U \in \mathbb{R}^{100 \times 10}$ ,  $V \in \mathbb{R}^{100 \times 10}$  and the sensing matrix  $\Phi \in \mathbb{R}^{20 \times 100}$  are randomly generated, and the input matrices were drawn from them with additive Gaussian noise. For the ALS-based algorithms, each element of the factor matrices and the sensing matrix is sampled from the Gaussian distribution with mean

0 and variance 1. For the multiplicative update algorithms which have the nonnegativity constraints, each element is sampled from a uniform distribution between 0 and 1. For the co-factorizations, 20 out of 100 columns were chosen for the uncompressed prior.

As a measure of the reconstruction performances, we used signal-to-noise ratio (SNR)

$$\text{SNR} = 10 \log_{10} \frac{\|\widetilde{X}\|_F^2}{\|\widetilde{X} - UV^T\|_F^2},$$

where  $\widetilde{X}$  is the true value of the input matrix,  $U$  and  $V$  are the factor matrices obtained by running the algorithms. The SNRs were averaged over the 100 different trials (Figure 2). The bounds were better for the cases where the matrix  $U$  is known in advance, than the cases where the matrix  $U$  is unknown. This indicates that the estimation of both  $U$  and  $V$  is more difficult than the estimation of  $V$  only. The bounds of CS-MF and CS-NMF were seriously degraded compared to the bounds of MF and NMF, as well as the actual performances, especially for the cases that the matrix  $U$  is unknown. However, the co-factorization based methods showed much improved bounds, as well as the performances, compared to CS-MF and CS-NMF. The use of uncompressed prior as the side information actually helps to improve the reconstruction performance in the learning from the compressed data.

As an illustrative example, we ran the algorithms for the compressively sampled data on the cameraman image (Figure 3). We divided the image using 8-by-8 patches to build an input matrix, and compress each column into the size of 49. 25 % of the data is randomly selected to be used for the uncompressed prior. CS-MCF brought much better reconstruction result than the CS-MF.

## 5 Conclusions

We have presented a method for matrix co-factorization in the compressed domain where compressively-sampled data and partial uncompressed data were jointly decomposed, sharing a factor matrix. We calculated CRBs for three matrix factorization methods (MF, CS-MF, CS-MCF) in the case of Gaussian likelihood, showing that CS-MCF improves the reconstruction performance over CS-MF. Numerical experiments on image confirmed the better performance of CS-MCF over CS-MF.

**Acknowledgments:** This work was supported by NIPAMSRSA Creative IT/SW Research Project, the Converging Research Center Program funded by Korea MEST (No. 2010K001171), NIPA ITRC Support Program (NIPA-2011-C1090-1131-0009), and NRF WCU Program (R31-2010-000-10100-0).

## References

[Candes and Tao, 2006] E. Candes and T. Tao. Near optimal signal recovery from random projections: Universal encoding strategies? *IEEE Transactions on Information Theory*, 52(12):5406–5425, 2006.

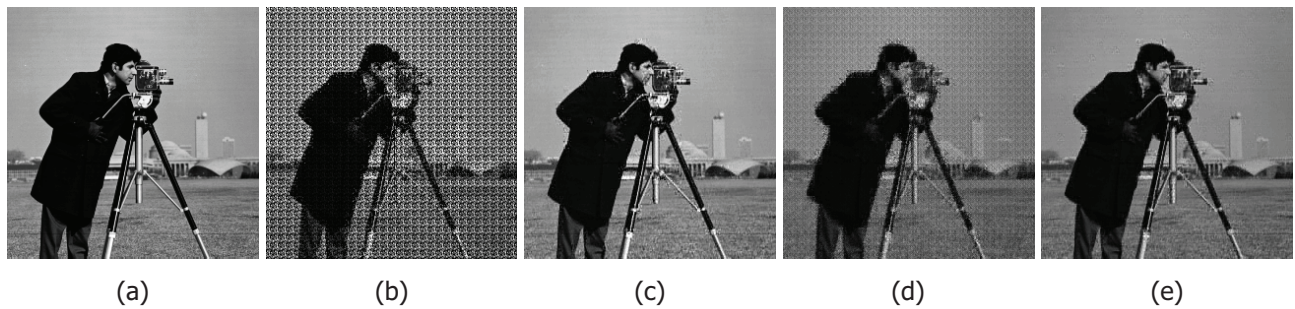


Figure 3: Reconstruction example of an image: (a) original image (b) the reconstruction result of CS-MF (c) the reconstruction result of CS-MCF (d) the reconstruction result of CS-NMF (e) the reconstruction result of CS-NMCF

- [Donoho, 2006] D. L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306, 2006.
- [Jolliffe, 2002] I. T. Jolliffe. *Principal Component Analysis*. Springer-Verlag, 2 edition, 2002.
- [Kay, 1993] S. M. Kay. *Fundamentals of Statistical Signal Processing: Estimation Theory*. Prentice Hall, 1993.
- [Koren et al., 2009] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *IEEE Computer*, 42(8):30–37, 2009.
- [Lee and Choi, 2009] H. Lee and S. Choi. Group nonnegative matrix factorization for EEG classification. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, Clearwater Beach, Florida, 2009.
- [Lee and Seung, 1999] D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401:788–791, 1999.
- [Lee et al., 2010] H. Lee, J. Yoo, and S. Choi. Semi-supervised nonnegative matrix factorization. *IEEE Signal Processing Letters*, 17(1):4–7, 2010.
- [O’Grady and Rickard, 2008] P. D. O’Grady and S. T. Rickard. Non-negative matrix factorisation of compressibly sampled non-negative signals. In *Proceedings of the 8th IMA International Conference on Mathematics in Signal Processing*, Cirencester, UK, 2008.
- [Singh and Gordon, 2008] A. P. Singh and G. J. Gordon. Relational learning via collective matrix factorization. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, Las Vegas, Nevada, 2008.
- [Smaragdakis and Brown, 2003] P. Smaragdakis and J. C. Brown. Non-negative matrix factorization for polyphonic music transcription. In *Proceedings of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, pages 177–180, New Paltz, NY, 2003.
- [Srebro et al., 2005] N. Srebro, J. D. M. Rennie, and T. Jaakkola. Maximum-margin matrix factorization. In *Advances in Neural Information Processing Systems (NIPS)*, volume 17. MIT Press, 2005.
- [Xu et al., 2003] W. Xu, X. Liu, and Y. Gong. Document clustering based on non-negative matrix factorization. In *Proceedings of the ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, Toronto, Canada, 2003.
- [Yoo and Choi, 2008] J. Yoo and S. Choi. Orthogonal non-negative matrix factorization: Multiplicative updates on Stiefel manifolds. In *Proceedings of the 9th International Conference on Intelligent Data Engineering and Automated Learning (IDEAL)*, Daejeon, Korea, 2008.
- [Yoo and Choi, 2009] J. Yoo and S. Choi. Weighted nonnegative matrix co-tri-factorization for collaborative prediction. In *Proceedings of the 1st Asian Conference on Machine Learning (ACML)*, Nanjing, China, 2009.
- [Yoo and Choi, 2010] J. Yoo and S. Choi. Orthogonal non-negative matrix-tri-factorization for co-clustering: Multiplicative updates on stiefel manifolds. *Information Processing and Management*, 46(5):559–570, 2010.
- [Yoo et al., 2010] J. Yoo, M. Kim, K. Kang, and S. Choi. Nonnegative matrix partial co-factorization for drum source separation. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Dallas, Texas, 2010.
- [Yu et al., 2005] K. Yu, S. Yu, and V. Tresp. Multi-label informed latent semantic indexing. In *Proceedings of the ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, Salvador, Brazil, 2005.
- [Zhu et al., 2007] S. Zhu, K. Yu, Y. Chi, and Y. Gong. Combining content and link for classification using matrix factorization. In *Proceedings of the ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, Amsterdam, The Netherlands, 2007.