

# Supplemental Material for TKDE-2015-05-0359

Kijung Shin, Lee Sael, U Kang

## 1 PROPOSED METHODS

### 1.1 Proofs of Update Rules

In this section, we present the proofs of the update rules in Section 3.5 of the main paper. Specifically, we prove the CDTF update rule for  $L_1$  regularization (Theorem 7), the CDTF update rule for the non-negativity constraint (Theorem 8), the SALS update rule for coupled tensor factorization (Theorem 9), and the update rule for bias terms commonly used by CDTF and SALS for the bias model (Theorem 10).

*Lemma 1 (Partial Derivative in CDTF):* For a parameter  $a_{i_n k}^{(n)}$ , let  $\hat{r}_{i_1 \dots i_N} = x_{i_1 \dots i_N} - \sum_{s \neq k} \prod_{l=1}^N a_{i_l s}^{(l)}$ ,  $g = -2 \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} \left( \hat{r}_{i_1 \dots i_N} \prod_{l \neq n} a_{i_l k}^{(l)} \right)$ , and  $d = 2 \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} \prod_{l \neq n} (a_{i_l k}^{(l)})^2$ , as in the main paper. Then,

$$\frac{\partial \left( \sum_{(i_1, \dots, i_N) \in \Omega} \left( x_{i_1 \dots i_N} - \sum_{s=1}^K \prod_{l=1}^N a_{i_l s}^{(l)} \right)^2 \right)}{\partial a_{i_n k}^{(n)}} = g + a_{i_n k}^{(n)} d.$$

*Proof:*

$$\begin{aligned} & \frac{\partial \left( \sum_{(i_1, \dots, i_N) \in \Omega} \left( x_{i_1 \dots i_N} - \sum_{s=1}^K \prod_{l=1}^N a_{i_l s}^{(l)} \right)^2 \right)}{\partial a_{i_n k}^{(n)}} \\ &= \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} 2 \left( x_{i_1 \dots i_N} - \sum_{s=1}^K \prod_{l=1}^N a_{i_l s}^{(l)} \right) \\ & \quad \times \frac{\partial \left( x_{i_1 \dots i_N} - \sum_{s=1}^K \prod_{l=1}^N a_{i_l s}^{(l)} \right)}{\partial a_{i_n k}^{(n)}} \\ &= \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} 2 \left( \hat{r}_{i_1 \dots i_N} - \prod_{l=1}^N a_{i_l k}^{(l)} \right) \left( - \prod_{l \neq n} a_{i_l k}^{(l)} \right) \\ &= g + a_{i_n k}^{(n)} d. \quad \square \end{aligned}$$

*Theorem 7: (Correctness of CDTF with  $L_1$ -regularization)* The update rule (12) in the main paper minimizes the loss function (11) with respect to the updated parameter. For an updated parameter  $a_{i_n k}^{(n)}$ , let  $\hat{r}_{i_1 \dots i_N} = x_{i_1 \dots i_N} - \sum_{s \neq k} \prod_{l=1}^N a_{i_l s}^{(l)}$ ,

$g = -2 \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} \left( \hat{r}_{i_1 \dots i_N} \prod_{l \neq n} a_{i_l k}^{(l)} \right)$ , and  $d = 2 \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} \prod_{l \neq n} (a_{i_l k}^{(l)})^2$ , as in the main paper. Then,

$$\arg \min_{a_{i_n k}^{(n)}} L_{Lasso}(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \begin{cases} (\lambda - g)/d & \text{if } g > \lambda \\ -(\lambda + g)/d & \text{if } g < -\lambda \\ 0 & \text{otherwise.} \end{cases}$$

*Proof:* By Lemma 1,

$$\begin{aligned} & \frac{\partial L_{Lasso}(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)})}{\partial a_{i_n k}^{(n)}} \\ &= \frac{\partial \left( \sum_{(i_1, \dots, i_N) \in \Omega} \left( x_{i_1 \dots i_N} - \sum_{s=1}^K \prod_{l=1}^N a_{i_l s}^{(l)} \right)^2 + \lambda \sum_{l=1}^N \|\mathbf{A}^{(l)}\|_1 \right)}{\partial a_{i_n k}^{(n)}} \\ &= g + a_{i_n k}^{(n)} d + \lambda \frac{\partial |a_{i_n k}^{(n)}|}{\partial a_{i_n k}^{(n)}} = \begin{cases} g + a_{i_n k}^{(n)} d + \lambda & \text{if } a_{i_n k}^{(n)} > 0 \\ g + a_{i_n k}^{(n)} d - \lambda & \text{if } a_{i_n k}^{(n)} < 0. \end{cases} \end{aligned}$$

**Case 1:** If  $g > \lambda$  ( $> 0$ ),  $\frac{\partial L_{Lasso}}{\partial a_{i_n k}^{(n)}} > a_{i_n k}^{(n)} d$ . Since  $d \geq 0$ ,  $a_{i_n k}^{(n)}$  should be negative for  $\frac{\partial L_{Lasso}}{\partial a_{i_n k}^{(n)}}$  to be zero. If  $a_{i_n k}^{(n)} < 0$ ,  $\frac{\partial L_{Lasso}}{\partial a_{i_n k}^{(n)}} = g + a_{i_n k}^{(n)} d - \lambda$ , and  $a_{i_n k}^{(n)} = (\lambda - g)/d$  ( $< 0$ ) makes  $\frac{\partial L_{Lasso}}{\partial a_{i_n k}^{(n)}} = 0$ . Since  $\frac{\partial^2 L_{Lasso}}{\partial (a_{i_n k}^{(n)})^2} = d \geq 0$ ,  $a_{i_n k}^{(n)} = (\lambda - g)/d$  minimizes  $L_{Lasso}(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)})$  with respect to  $a_{i_n k}^{(n)}$ .

**Case 2:** Likewise, if  $g < -\lambda$  ( $< 0$ ),  $\frac{\partial L_{Lasso}}{\partial a_{i_n k}^{(n)}} < a_{i_n k}^{(n)} d$ . Since  $d \geq 0$ ,  $a_{i_n k}^{(n)}$  should be positive for  $\frac{\partial L_{Lasso}}{\partial a_{i_n k}^{(n)}}$  to be zero. If  $a_{i_n k}^{(n)} > 0$ ,  $\frac{\partial L_{Lasso}}{\partial a_{i_n k}^{(n)}} = g + a_{i_n k}^{(n)} d + \lambda$ , and  $a_{i_n k}^{(n)} = -(\lambda + g)/d$  ( $> 0$ ) makes  $\frac{\partial L_{Lasso}}{\partial a_{i_n k}^{(n)}} = 0$ . Since  $\frac{\partial^2 L_{Lasso}}{\partial (a_{i_n k}^{(n)})^2} = d \geq 0$ ,  $a_{i_n k}^{(n)} = -(\lambda + g)/d$  minimizes  $L_{Lasso}(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)})$  with respect to  $a_{i_n k}^{(n)}$ .

**Case 3:** On the other hand, if  $-\lambda \leq g \leq \lambda$ ,

$$\frac{\partial L_{Lasso}}{\partial a_{i_n k}^{(n)}} = \begin{cases} g + a_{i_n k}^{(n)} d + \lambda \geq a_{i_n k}^{(n)} d \geq 0 & \text{if } a_{i_n k}^{(n)} > 0 \\ g + a_{i_n k}^{(n)} d - \lambda \leq a_{i_n k}^{(n)} d \leq 0 & \text{if } a_{i_n k}^{(n)} < 0 \end{cases}$$

That is,  $L_{Lasso}(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)})$  decreases if  $a_{i_n k}^{(n)} < 0$  and increases if  $a_{i_n k}^{(n)} > 0$ , given other parameters. Thus,  $L_{Lasso}(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)})$ , which is a continuous

function, is minimized with respect to  $a_{i_n k}^{(n)}$  if  $a_{i_n k}^{(n)} = 0$ . Hence,

$$\arg \min_{a_{i_n k}^{(n)}} L_{Lasso}(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \begin{cases} (\lambda - g)/d & \text{if } g > \lambda \\ -(\lambda + g)/d & \text{if } g < -\lambda \\ 0 & \text{otherwise.} \end{cases}$$

□

*Theorem 8: (Correctness of CDTF with the non-negativity constraint)* The update rule (14) in the main paper minimizes the loss function (1) with respect to the updated parameter under the non-negativity constraint. For an updated parameter  $a_{i_n k}^{(n)}$ , let  $\hat{r}_{i_1 \dots i_N} = x_{i_1 \dots i_N} - \sum_{s \neq k} \prod_{l=1}^N a_{i_l s}^{(l)}$ ,  $g = -2 \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} (\hat{r}_{i_1 \dots i_N} \prod_{l \neq n} a_{i_l k}^{(l)})$ , and  $d = 2 \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} \prod_{l \neq n} (a_{i_l k}^{(l)})^2$ , as in the main paper. Then,

$$\arg \min_{a_{i_n k}^{(n)} \geq 0} L(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \max\left(\frac{-g}{d+2\lambda}, 0\right).$$

*Proof:* By Lemma 1,

$$\begin{aligned} & \frac{\partial L(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)})}{\partial a_{i_n k}^{(n)}} \\ &= \frac{\partial \left( \sum_{(i_1, \dots, i_N) \in \Omega} \left( x_{i_1 \dots i_N} - \sum_{s=1}^K \prod_{l=1}^N a_{i_l s}^{(l)} \right)^2 + \lambda \sum_{l=1}^N \|\mathbf{A}^{(l)}\|_F^2 \right)}{\partial a_{i_n k}^{(n)}} \\ &= g + (d+2\lambda)a_{i_n k}^{(n)}. \end{aligned}$$

Thus,

$$\frac{\partial L}{\partial a_{i_n k}^{(n)}} \begin{cases} > 0 & \text{if } a_{i_n k}^{(n)} > -g/(d+2\lambda) \\ = 0 & \text{if } a_{i_n k}^{(n)} = -g/(d+2\lambda) \\ < 0 & \text{otherwise.} \end{cases}$$

**Case 1:** If  $-g/(d+2\lambda) \geq 0$ , since  $\frac{\partial^2 L}{\partial (a_{i_n k}^{(n)})^2} = d+2\lambda \geq 0$ ,  $L(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)})$  is minimized if  $a_{i_n k}^{(n)} = -g/(d+2\lambda)$  with respect to  $a_{i_n k}^{(n)}$ .

**Case 2:** On the other hand, if  $-g/(d+2\lambda) < 0$ , under the constraint that  $a_{i_n k}^{(n)} \geq 0$ ,  $L(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)})$  is minimized with respect to  $a_{i_n k}^{(n)}$  if  $a_{i_n k}^{(n)} = 0$ . This is because  $L(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)})$ , which is a continuous function, monotonically increases (i.e.,  $\frac{\partial L}{\partial a_{i_n k}^{(n)}} > 0$ ) if  $a_{i_n k}^{(n)} \geq 0$  ( $> -g/(d+2\lambda)$ ).

Hence,

$$\arg \min_{a_{i_n k}^{(n)} \geq 0} L(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \max\left(\frac{-g}{d+2\lambda}, 0\right).$$

□

*Theorem 9: (Correctness of SALS in Coupled Tensor Factorization)* The update rule (16) in the main paper minimizes (15) with respect to the updated parameters. Let  ${}_x \mathcal{R}$  and  ${}_y \mathcal{R}$  be the residual tensors for  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively. That is, for  $C$  updated parameters  $a_{i_1 k_1}^{(1)}, \dots, a_{i_1 k_C}^{(1)}$ ,  $x \hat{r}_{i_1 \dots i_{N_x}} = x_{i_1 \dots i_{N_x}} - \sum_{k=1}^K \prod_{n=1}^{N_x} x a_{i_n k}^{(n)} + \sum_{c=1}^C \prod_{n=1}^{N_x} x a_{i_n k_c}^{(n)}$  and  $y \hat{r}_{i_1 \dots i_{N_y}} =$

$y_{i_1 \dots i_{N_y}} - \sum_{k=1}^K \prod_{n=1}^{N_y} y a_{i_n k}^{(n)} + \sum_{c=1}^C \prod_{n=1}^{N_y} y a_{i_n k_c}^{(n)}$ . Likewise, let  ${}_x \Omega$  and  ${}_y \Omega$  be the sets of indices of the observable entries in  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively. Then,

$$\begin{aligned} & \arg \min_{[a_{i_1 k_1}^{(1)}, \dots, a_{i_1 k_C}^{(1)}]^T} L_{Coupled}(x \mathbf{A}^{(1)}, \dots, x \mathbf{A}^{(N_x)}, y \mathbf{A}^{(1)}, \dots, y \mathbf{A}^{(N_y)}) \\ &= ({}_x \mathbf{B}_{i_1}^{(1)} + {}_y \mathbf{B}_{i_1}^{(1)} + \lambda \mathbf{I}_C)^{-1} ({}_x \mathbf{c}_{i_1}^{(1)} + {}_y \mathbf{c}_{i_1}^{(1)}), \end{aligned}$$

where  ${}_x \mathbf{B}_{i_1}^{(1)}$  and  ${}_y \mathbf{B}_{i_1}^{(1)}$  are  $C$  by  $C$  matrices whose entries are

$$\begin{aligned} ({}_x \mathbf{B}_{i_1}^{(1)})_{c_1 c_2} &= \sum_{(i_1, \dots, i_{N_x}) \in {}_x \Omega_{i_1}^{(1)}} \left( \prod_{n \neq 1} x a_{i_n k_{c_1}}^{(n)} \prod_{n \neq 1} x a_{i_n k_{c_2}}^{(n)} \right), \\ ({}_y \mathbf{B}_{i_1}^{(1)})_{c_1 c_2} &= \sum_{(i_1, \dots, i_{N_y}) \in {}_y \Omega_{i_1}^{(1)}} \left( \prod_{n \neq 1} y a_{i_n k_{c_1}}^{(n)} \prod_{n \neq 1} y a_{i_n k_{c_2}}^{(n)} \right), \forall c_1, c_2, \end{aligned}$$

${}_x \mathbf{c}_{i_1}^{(1)}$  and  ${}_y \mathbf{c}_{i_1}^{(1)}$  are length  $C$  vectors whose entries are

$$\begin{aligned} ({}_x \mathbf{c}_{i_1}^{(1)})_c &= \sum_{(i_1, \dots, i_{N_x}) \in {}_x \Omega_{i_1}^{(1)}} \left( x \hat{r}_{i_1 \dots i_{N_x}} \prod_{n \neq 1} x a_{i_n k_c}^{(n)} \right), \\ ({}_y \mathbf{c}_{i_1}^{(1)})_c &= \sum_{(i_1, \dots, i_{N_y}) \in {}_y \Omega_{i_1}^{(1)}} \left( y \hat{r}_{i_1 \dots i_{N_y}} \prod_{n \neq 1} y a_{i_n k_c}^{(n)} \right), \forall c \end{aligned}$$

and  $\mathbf{I}_C$  is the  $C$  by  $C$  identity matrix.

*Proof:*

$$\begin{aligned} & \frac{\partial L(x \mathbf{A}^{(1)}, \dots, x \mathbf{A}^{(N_x)})}{\partial a_{i_1 k_c}^{(1)}} = \frac{\partial L(x \mathbf{A}^{(1)}, \dots, x \mathbf{A}^{(N_x)})}{\partial (x a_{i_1 k_c}^{(1)})} \\ &= \frac{\partial \left( \sum_{(i_1, \dots, i_{N_x}) \in {}_x \Omega} \left( x_{i_1 \dots i_{N_x}} - \sum_{s=1}^K \prod_{n=1}^{N_x} x a_{i_n s}^{(n)} \right)^2 \right)}{\partial (x a_{i_1 k_c}^{(1)})} \\ &+ \frac{\partial \lambda \sum_{n=1}^{N_x} \|x \mathbf{A}^{(n)}\|_F^2}{\partial (x a_{i_1 k_c}^{(1)})} \\ &= \sum_{(i_1, \dots, i_{N_x}) \in {}_x \Omega_{i_1}^{(1)}} 2 \left( x_{i_1 \dots i_{N_x}} - \sum_{s=1}^K \prod_{n=1}^{N_x} x a_{i_n s}^{(n)} \right) \\ &\times \frac{\partial \left( x_{i_1 \dots i_{N_x}} - \sum_{s=1}^K \prod_{n=1}^{N_x} x a_{i_n s}^{(n)} \right)}{\partial (x a_{i_1 k_c}^{(1)})} + 2\lambda (x a_{i_1 k_c}^{(1)}) \\ &= \sum_{(i_1, \dots, i_{N_x}) \in {}_x \Omega_{i_1}^{(1)}} 2 \left( \sum_{s=1}^C \prod_{n=1}^{N_x} x a_{i_n k_s}^{(n)} - x \hat{r}_{i_1 \dots i_{N_x}} \right) \prod_{n \neq 1} x a_{i_n k_c}^{(n)} \\ &+ 2\lambda a_{i_1 k_c}^{(1)}. \end{aligned}$$

Likewise,

$$\begin{aligned} & \frac{\partial L(y \mathbf{A}^{(1)}, \dots, y \mathbf{A}^{(N_y)})}{\partial a_{i_1 k_c}^{(1)}} \\ &= \sum_{(i_1, \dots, i_{N_y}) \in {}_y \Omega_{i_1}^{(1)}} 2 \left( \sum_{s=1}^C \prod_{n=1}^{N_y} y a_{i_n k_s}^{(n)} - y \hat{r}_{i_1 \dots i_{N_y}} \right) \prod_{n \neq 1} y a_{i_n k_c}^{(n)} \\ &+ 2\lambda a_{i_1 k_c}^{(1)}. \end{aligned}$$

From these,

$$\begin{aligned}
\frac{\partial L_{Coupled}}{\partial a_{i_1 k_c}^{(1)}} &= 0, \forall c, 1 \leq c \leq C \\
\Leftrightarrow \sum_{(i_1, \dots, i_{N_x}) \in x\Omega_{i_1}^{(1)}} &\left( \sum_{s=1}^C \prod_{n=1}^{N_x} x a_{i_n k_s}^{(n)} - x \hat{r}_{i_1 \dots i_{N_x}} \right) \prod_{n \neq 1} x a_{i_n k_c}^{(n)} \\
+ \sum_{(i_1, \dots, i_{N_y}) \in y\Omega_{i_1}^{(1)}} &\left( \sum_{s=1}^C \prod_{n=1}^{N_y} y a_{i_n k_s}^{(n)} - y \hat{r}_{i_1 \dots i_{N_y}} \right) \prod_{n \neq 1} y a_{i_n k_c}^{(n)} \\
+ \lambda a_{i_1 k_c}^{(1)} &= 0, \forall c \\
\Leftrightarrow \sum_{(i_1, \dots, i_{N_x}) \in x\Omega_{i_1}^{(1)}} &\left( \sum_{s=1}^C \left( x a_{i_1 k_s}^{(1)} \prod_{n \neq 1} x a_{i_n k_s}^{(n)} \right) \prod_{n \neq 1} x a_{i_n k_c}^{(n)} \right) \\
+ \sum_{(i_1, \dots, i_{N_y}) \in y\Omega_{i_1}^{(1)}} &\left( \sum_{s=1}^C \left( y a_{i_1 k_s}^{(1)} \prod_{n \neq 1} y a_{i_n k_s}^{(n)} \right) \prod_{n \neq 1} y a_{i_n k_c}^{(n)} \right) \\
+ \lambda a_{i_1 k_c}^{(1)} & \\
= \sum_{(i_1, \dots, i_{N_x}) \in x\Omega_{i_1}^{(1)}} &\left( x \hat{r}_{i_1 \dots i_{N_x}} \prod_{n \neq 1} x a_{i_n k_c}^{(n)} \right) \\
+ \sum_{(i_1, \dots, i_{N_y}) \in y\Omega_{i_1}^{(1)}} &\left( y \hat{r}_{i_1 \dots i_{N_y}} \prod_{n \neq 1} y a_{i_n k_c}^{(n)} \right), \forall c \\
\Leftrightarrow (x \mathbf{B}_{i_1}^{(1)} + y \mathbf{B}_{i_1}^{(1)} + \lambda \mathbf{I}_C) [a_{i_1 k_1}^{(1)}, \dots, a_{i_1 k_C}^{(1)}]^T &= (x \mathbf{c}_{i_1}^{(1)} + y \mathbf{c}_{i_1}^{(1)}).
\end{aligned}$$

Hence,

$$\begin{aligned}
\arg \min_{[a_{i_1 k_1}^{(1)}, \dots, a_{i_1 k_C}^{(1)}]^T} L_{Coupled}(x \mathbf{A}^{(1)}, \dots, x \mathbf{A}^{(N_x)}, y \mathbf{A}^{(1)}, \dots, y \mathbf{A}^{(N_y)}) \\
= (x \mathbf{B}_{i_1}^{(1)} + y \mathbf{B}_{i_1}^{(1)} + \lambda \mathbf{I}_C)^{-1} (x \mathbf{c}_{i_1}^{(1)} + y \mathbf{c}_{i_1}^{(1)}).
\end{aligned}$$

**Theorem 10: (Correctness of the Update Rule for Bias Terms)** The update rule (18) in the main paper minimizes (17) with respect to the updated parameter. For an updated parameter  $b_{i_n}^{(n)}$ , let  $\bar{r}_{i_1 \dots i_N} = x_{i_1 \dots i_N} - \sum_{k=1}^K \prod_{l=1}^N a_{i_l k}^{(l)} - \sum_{l \neq n} b_{i_l}^{(l)} - \mu$ , as in the main paper. Then,

$$\begin{aligned}
\arg \min_{b_{i_n}^{(n)}} L_{Bias}(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(N)}) \\
= \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} \bar{r}_{i_1 \dots i_N} / (\lambda_{\mathbf{b}} + |\Omega_{i_n}^{(n)}|).
\end{aligned}$$

*Proof:*

$$\begin{aligned}
\frac{\partial L_{Bias}(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(N)})}{\partial b_{i_n}^{(n)}} \\
= \frac{\partial \left( \sum_{(i_1, \dots, i_N) \in \Omega} \left( x_{i_1 \dots i_N} - \sum_{s=1}^K \prod_{l=1}^N a_{i_l s}^{(l)} - \sum_{l=1}^N b_{i_l}^{(l)} - \mu \right)^2 \right)}{\partial b_{i_n}^{(n)}} \\
+ \frac{\partial \left( \lambda_{\mathbf{A}} \sum_{l=1}^N \|\mathbf{A}^{(l)}\|_F^2 \right)}{\partial b_{i_n}^{(n)}} + \frac{\partial \left( \lambda_{\mathbf{b}} \sum_{l=1}^N \|\mathbf{b}^{(l)}\|_F^2 \right)}{\partial b_{i_n}^{(n)}}
\end{aligned}$$

$$\begin{aligned}
&= \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} 2 \left( x_{i_1 \dots i_N} - \sum_{s=1}^K \prod_{l=1}^N a_{i_l s}^{(l)} - \sum_{l=1}^N b_{i_l}^{(l)} - \mu \right) \\
&\quad \times \frac{\partial \left( x_{i_1 \dots i_N} - \sum_{s=1}^K \prod_{l=1}^N a_{i_l s}^{(l)} - \sum_{l=1}^N b_{i_l}^{(l)} - \mu \right)}{\partial b_{i_n}^{(n)}} + \frac{\lambda_{\mathbf{b}} \sum_{l=1}^N \|\mathbf{b}^{(l)}\|_F^2}{\partial b_{i_n}^{(n)}} \\
&= \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} -2(\bar{r}_{i_1 \dots i_N} - b_{i_n}^{(n)}) + 2\lambda_{\mathbf{b}} b_{i_n}^{(n)} \\
&= 2(\lambda_{\mathbf{b}} + |\Omega_{i_n}^{(n)}|) b_{i_n}^{(n)} - \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} 2\bar{r}_{i_1 \dots i_N}.
\end{aligned}$$

Since  $\frac{\partial^2 L_{Bias}}{\partial (b_{i_n}^{(n)})^2} = 2(\lambda_{\mathbf{b}} + |\Omega_{i_n}^{(n)}|) \geq 0$ ,  $L_{Bias}(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(N)})$  is minimized with respect to  $b_{i_n}^{(n)}$  if  $b_{i_n}^{(n)} = \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} \bar{r}_{i_1 \dots i_N} / (\lambda_{\mathbf{b}} + |\Omega_{i_n}^{(n)}|)$ , which entails  $\frac{\partial L_{Bias}}{\partial b_{i_n}^{(n)}} = 0$ . Hence,

$$\begin{aligned}
\arg \min_{b_{i_n}^{(n)}} L_{Bias}(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}, \mathbf{b}^{(1)}, \dots, \mathbf{b}^{(N)}) \\
= \sum_{(i_1, \dots, i_N) \in \Omega_{i_n}^{(n)}} \bar{r}_{i_1 \dots i_N} / (\lambda_{\mathbf{b}} + |\Omega_{i_n}^{(n)}|).
\end{aligned}$$

□

## 1.2 Pseudocodes

We present the pseudocodes of the SALS variants described in Section 3.5 of the main paper.

### 1.2.1 SALS for Coupled Tensor Factorization

Algorithm 6 describes SALS for coupled tensor factorization, where two tensors, denoted by  $\mathcal{X}$  and  $\mathcal{Y}$ , share their first mode without loss of generality. We denote the residual tensors for  $\mathcal{X}$  and  $\mathcal{Y}$  by  $x\mathcal{R}$  and  $y\mathcal{R}$ , respectively. The lengths of the  $n$ -th modes of  $\mathcal{X}$  and  $\mathcal{Y}$  are denoted by  $xI_n$  and  $yI_n$ , respectively.

---

#### Algorithm 6: SALS for Coupled Tensor Factorization

---

**Input :**  $\mathcal{X}, \mathcal{Y}, K, \lambda$   
**Output:**  $\mathbf{A}^{(1)}, x\mathbf{A}^{(n)}$  for  $n = 2 \dots N_x, y\mathbf{A}^{(n)}$  for  $n = 2 \dots N_y$

- 1 initialize  $x\mathcal{R}, y\mathcal{R}, \mathbf{A}^{(1)}, x\mathbf{A}^{(n)}$  and  $y\mathbf{A}^{(n)}$  for all  $n \geq 2$
- 2 **for** outer iter = 1.. $T_{out}$  **do**
- 3     **for** split iter = 1.. $\lceil \frac{K}{C} \rceil$  **do**
- 4         choose  $k_1, \dots, k_C$  (from columns not updated yet)
- 5         compute  $x\hat{\mathcal{R}}$  and  $y\hat{\mathcal{R}}$  using (8)
- 6         **for** inner iter = 1.. $T_{in}$  **do**
- 7             **for**  $i_1 = 1 \dots I_1 (= xI_1 = yI_1)$  **do**
- 8                 update  $a_{i_1 k_1}^{(1)}, \dots, a_{i_1 k_C}^{(1)}$  using (16)
- 9             **for**  $n = 2 \dots N_x$  **do**
- 10                 **for**  $i_n = 1 \dots xI_n$  **do**
- 11                     update  $x a_{i_n k_1}^{(n)}, \dots, x a_{i_n k_C}^{(n)}$  using (9)
- 12             **for**  $n = 2 \dots N_y$  **do**
- 13                 **for**  $i_n = 1 \dots yI_n$  **do**
- 14                     update  $y a_{i_n k_1}^{(n)}, \dots, y a_{i_n k_C}^{(n)}$  using (9)
- 15             update  $x\mathcal{R}$  and  $y\mathcal{R}$  using (10)

---

**Algorithm 7: SALS for Bias Model**


---

**Input** :  $\mathcal{X}, K, \lambda_A, \lambda_B$   
**Output**:  $\mathbf{A}^{(n)}$  for all  $n$ ,  $\mathbf{b}^{(n)}$  for all  $n, \mu$

- 1 compute  $\mu$  (the mean of the observable entries of  $\mathcal{X}$ )
- 2 initialize  $\mathcal{R}, \mathbf{A}^{(n)}$  for all  $n$ , and  $\mathbf{b}^{(n)}$  for all  $n$
- 3 **for** *outer iter* = 1.. $T_{out}$  **do**
- 4     **for** *split iter* = 1.. $\lceil \frac{K}{C} \rceil$  **do**
- 5         choose  $k_1, \dots, k_C$  (from columns not updated yet)
- 6         compute  $\mathcal{R}$  using (8)
- 7         **for** *inner iter* = 1.. $T_{in}$  **do**
- 8             **for**  $n = 1..N$  **do**
- 9                 **for**  $i_n = 1..I_n$  **do**
- 10                     update  $a_{i_n k_1}^{(n)}, \dots, a_{i_n k_C}^{(n)}$  using (9)
- 11             update  $\mathcal{R}$  using (10)
- 12     **for**  $n = 1..N$  **do**
- 13         **for**  $i_n = 1..I_n$  **do**
- 14             update  $b_{i_n}^{(n)}$  using (18)
- 15             update  $\mathcal{R}$  using (19)

---

### 1.2.2 SALS for Bias Model

SALS for the bias model is described in Algorithm 7, where each  $(i_1, \dots, i_N)$ th entry of  $\mathcal{R}$  is  $r_{i_1 \dots i_N} = x_{i_1 \dots i_N} - \mu - \sum_{n=1}^N b_{i_n}^{(n)} - \sum_{k=1}^K \prod_{n=1}^N a_{i_n k}^{(n)}$ , as explained in Section 3.5.5 of the main paper.

## 2 OPTIMIZATION ON MAPREDUCE

In this section, we present the details of the optimization techniques described in Section 4 of the main paper.

### 2.1 Local Disk Caching

As explained in Section 4.1 of the main paper, in our MAPREDUCE implementation of CDTF and SALS with local disk caching,  $\mathcal{X}$  entries are distributed across machines and cached in their local disk during the map and reduce stages. Algorithm 8 gives the details of the map and reduce stages. The rest of CDTF and SALS runs in the close stage (cleanup stage in Hadoop) using the cached data.

### 2.2 Direct Communication

In the main paper, we introduce direct communication between reducers using distributed file system to overcome the rigidity of MAPREDUCE model. Algorithm 9 describes the implementation of  $m\mathbf{a}_{*k}^{(n)}$  broadcast in CDTF (line 10 of Algorithm 3 in the main paper) based on this communication method.

### 2.3 Greedy Row Assignment

Our MAPREDUCE implementation of SALS and CDTF uses the greedy row assignment, explained in Section 3.4.3 of the main paper. In this section, we explain our MAPREDUCE implementation of the greedy row assignment. We assume that  $\mathcal{X}$  is stored

**Algorithm 8: Data distribution in CDTF and SALS with local disk caching**


---

**Input** :  $\mathcal{X}, mS_n$  for all  $m$  and  $n$   
**Output**:  $m\Omega^{(n)}$  entries of  $\mathcal{R}(=\mathcal{X})$  for all  $m$  and  $n$

- 1 Map(Key  $k$ , Value  $v$ )
- 2 **begin**
- 3      $((i_1, \dots, i_N), x_{i_1 \dots i_N}) \leftarrow v$
- 4     **for**  $n = 1, \dots, N$  **do**
- 5         find  $m$  where  $i_n \in mS_n$
- 6         emit  $\langle (m, n), ((i_1, \dots, i_N), x_{i_1 \dots i_N}) \rangle$
- 7 **end**
- 8 Partitioner(Key  $k$ , Value  $v$ )
- 9 **begin**
- 10      $(m, n) \leftarrow k$
- 11     assign  $\langle k, v \rangle$  to machine  $m$
- 12 **end**
- 13 Reduce(Key  $k$ , Value  $v[1..|v|]$ )
- 14 **begin**
- 15      $(m, n) \leftarrow k$
- 16     create a file on the local disk to cache  $m\Omega^{(n)}$  entries of  $\mathcal{R}$
- 17     **foreach**  $((i_1, \dots, i_N), x_{i_1 \dots i_N}) \in v$  **do**
- 18         write  $((i_1, \dots, i_N), x_{i_1 \dots i_N})$  to the file
- 19 **end**

---

**Algorithm 9:  $m\mathbf{a}_{*k}^{(n)}$  broadcast in CDTF**


---

**Input** :  $m\mathbf{a}_{*k}^{(n)}$  (parameters to broadcast)  
**Output**:  $\mathbf{a}_{*k}^{(n)}$  (parameters received from others)

- 1 **begin**
- 2     create a data file  $mA$  on the distributed file system (DFS)
- 3     write  $m\mathbf{a}_{*k}^{(n)}$  on the datafile
- 4     create a dummy file  $mD$  on DFS
- 5     **while** *not all data files are read* **do**
- 6         get the list of dummy files from DFS
- 7         **foreach**  $m'D$  in the list **do**
- 8             **if**  $m'A$  are not read **then**
- 9                 read  $m'\mathbf{a}_{*k}^{(n)}$  from  $m'A$
- 10 **end**

---

on the distributed file system. At the first stage,  $|\Omega_{i_n}^{(n)}|$  for all  $n$  and  $i_n$  is computed. Specifically, mappers output  $\langle (n, i_n), 1 \rangle$  for all  $n$  for each entry  $x_{i_1 \dots i_N}$ , and reducers output  $\langle (n, i_n), |\Omega_{i_n}^{(n)}| \rangle$  for all  $n$  and  $i_n$  by counting the number of values for each key. At the second stage, the outputs are aggregated to a single reducer which runs the rest of Algorithm 5 in the main paper.

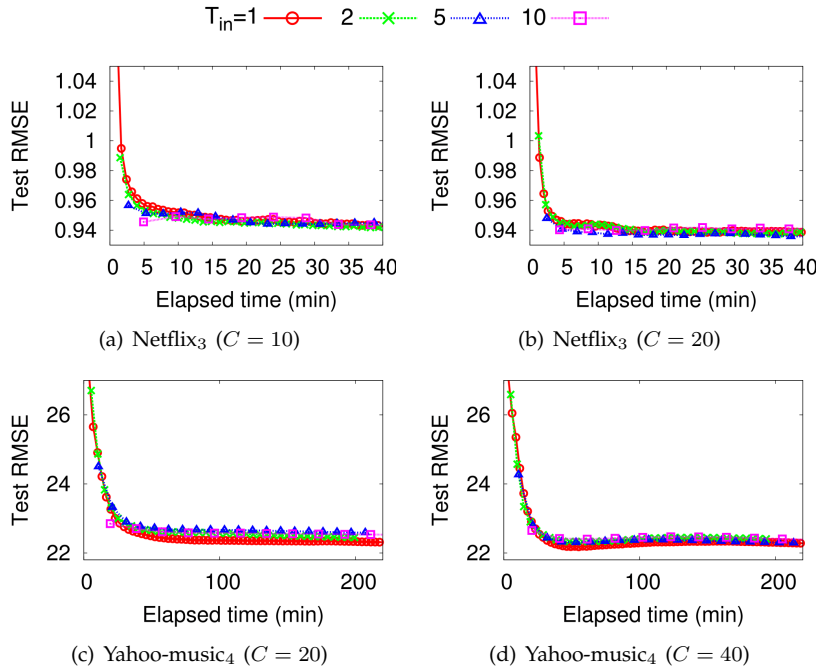
## 3 EXPERIMENTS

In this section, we design and conduct additional experiments to answer the following questions:

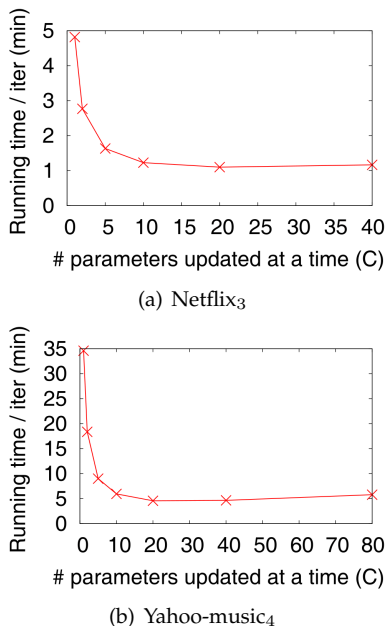
- How do different numbers of inner iterations ( $T_{in}$ ) affect the convergence of SALS?
- How do different numbers of columns updated at a time ( $C$ ) affect the running time of SALS?

### 3.1 Experimental Settings

We ran experiments on a 20-node Hadoop cluster. Each node had an Intel Xeon E3-1230 3.3GHz CPU.



**Fig. 14:** Effects of  $T_{in}$  (i.e., inner iterations) on the convergence of SALS when  $C$  (i.e., the number of columns updated at a time) has large enough values. The effects of  $T_{in}$  on convergence speed and the quality of converged solutions are marginal.



**Fig. 15:** Effects of the number of columns updated at a time ( $C$ ) on the running time of SALS. Running time per iteration decreased until  $C = 20$ , then started to increase.

The maximum heap size per reducer was set to 8GB. Other experimental settings, including datasets and parameter values ( $\lambda$  and  $K$ ), were the same as those in the main paper. The number of reducers was set to 20. We used the root mean square error (RMSE) on a held-out test set, which is commonly used in recommender systems, to measure the accuracy, as in the main paper.

### 3.2 Effects of the Number of Inner Iterations (i.e., $T_{in}$ ) on the Convergence of SALS

We compared the convergence properties of SALS with different  $T_{in}$  values. Especially, we focused on cases where  $C$  (i.e., the number of columns updated at a time) has large enough values. The effect of  $T_{in}$  when  $C$  is set to one and thus SALS is equivalent to CDTF can be found in the main paper. As seen in Figure 14, the effects of  $T_{in}$  on convergence speed and the quality of converged solutions are neither distinct nor consistent. When  $C$  is set to one, however, high  $T_{in}$  values are preferred (see Section 5.7 of the main paper for detailed experimental results).

### 3.3 Effects of the Number of Columns Updated at a Time (i.e., $C$ ) on Running Time of SALS

We measured the running times per iteration in SALS, as we increased  $C$  from 1 to  $K$ . As seen in Figure 15, running time per iteration decreased until  $C = 20$ , then started to increase. As  $C$  increases, the amount of disk I/O declines since it depends on the number of times that the entries of  $\mathcal{R}$  or  $\hat{\mathcal{R}}$  are streamed from disk, which is inversely proportional to  $C$ . Conversely, computational cost increases quadratically with regard to  $C$ . At small  $C$  values, the decrease in the amount of disk I/O was greater and led to a downward trend of running time per iteration. The opposite happened at large  $C$  values.