

## Supplement to: Bayesian Latent Factor Model for Higher-order Data

**Zerui Tao**

**Xuyang Zhao**

**Toshihisa Tanaka**

*Department of Electrical and Electronic Engineering*

*Tokyo University of Agriculture and Technology, Tokyo, Japan*

*RIKEN Center for Advanced Intelligence Project (AIP), Tokyo, Japan*

ZERUI.TAO@RIKEN.JP

ZHAO@SIP.TUAT.AC.JP

TANAKAT@CC.TUAT.AC.JP

**Qibin Zhao\***

*RIKEN Center for Advanced Intelligence Project (AIP), Tokyo, Japan*

QIBIN.ZHAO@RIKEN.JP

**Editors:** Vineeth N Balasubramanian and Ivor Tsang

### Appendix A. Gibbs Sampler

In this section, we firstly derive the Gibbs sampler of our model. Since the distributions are all conditionally conjugate, the derivation is straightforward. For the parameters, we denote  $\mathcal{Q} = \{\mathcal{Q}^{(d)}\}_{d=1}^D$ ,  $\boldsymbol{\eta} = \{\boldsymbol{\eta}^{(n)}\}_{n=1}^N$ ,  $\boldsymbol{\phi} = \{\boldsymbol{\phi}^{(d)}\}_{d=1}^D$ ,  $\boldsymbol{\delta} = \{\boldsymbol{\delta}^{(d)}\}_{d=1}^{D+1}$  and  $\boldsymbol{\Theta} = \{\mathcal{Q}, \boldsymbol{\eta}, \boldsymbol{\phi}, \boldsymbol{\delta}, \tau\}$ . The joint distribution of the TRLF model is

$$\log p(\mathcal{Y}, \boldsymbol{\Theta}) = \log p(\mathcal{Y} \mid \boldsymbol{\Theta}) + \log p(\boldsymbol{\Theta}), \quad (1)$$

where

$$\log p(\mathcal{Y} \mid \boldsymbol{\Theta}) = \log p(\mathcal{Y} \mid \mathcal{Q}, \boldsymbol{\eta}, \tau) = \sum_{n=1}^N \log p(\mathcal{Y}^{(n)} \mid \mathcal{Q}, \boldsymbol{\eta}^{(n)}, \tau).$$

The prior distribution is

$$\log p(\boldsymbol{\Theta}) = \log p(\mathcal{Q} \mid \boldsymbol{\phi}, \mathbf{u}) + \log p(\boldsymbol{\phi}) + \log p(\boldsymbol{\delta}) + \log p(\tau) + \log p(\boldsymbol{\eta}),$$

where

$$\begin{aligned} \log p(\mathcal{Q} \mid \boldsymbol{\phi}, \mathbf{u}) &= \sum_{d=1}^D \log p(\mathcal{Q}^{(d)} \mid \boldsymbol{\phi}^{(d)}, \mathbf{u}^{(d)}, \mathbf{u}^{(d+1)}), \quad u_i^{(d)} = \prod_{h=1}^l \delta_h^{(d)}, \\ \log p(\boldsymbol{\delta}) &= \sum_{d=1}^{D+1} \log p(\boldsymbol{\delta}^{(d)}), \quad \log p(\boldsymbol{\phi}) = \sum_{d=1}^D \log p(\boldsymbol{\phi}^{(d)}), \quad \log p(\boldsymbol{\eta}) = \sum_{n=1}^N \log p(\boldsymbol{\eta}^{(n)}). \end{aligned}$$

Given the joint distribution Eq. (1), we can establish the Gibbs sampler algorithm by sampling from the conditional distributions sequentially.

---

\* Corresponding author

**Sample  $\mathcal{Q}$**  The conditional distributions of the core tensors  $\mathcal{Q}$  involve subchains of  $\mathcal{Q}$  and  $\boldsymbol{\eta}$ . Hence, to make the notations consistent, we denote  $\mathcal{Q}^{(D+1)} = \boldsymbol{\eta}$ . For  $d = 1, \dots, D$  and  $i = 1, \dots, P_d$ , the conditional distribution is

$$p(\mathcal{Q}^{(d)}[i] \mid -) \sim \mathcal{N}(\text{vec}(\mathcal{Q}^{(d)}[i]) \mid \boldsymbol{\mu}_{\mathcal{Q}^{(d)}}^{(i)}, \boldsymbol{\Lambda}_{\mathcal{Q}^{(d)}}^{(i)}), \quad (2)$$

where

$$\begin{aligned} \boldsymbol{\mu}_{\mathcal{Q}^{(d)}}^{(i)} &= \boldsymbol{\Lambda}_{\mathcal{Q}^{(d)}}^{(i)} \left( \tau \sum_{n=1}^N \sum_{\mathbf{p}-d} \mathcal{Y}_{\mathbf{p}}^{(n)} \text{vec}(\mathcal{Q}^{\neq d, \top}[\overline{\mathbf{p}-d}]) \right), \\ \boldsymbol{\Lambda}_{\mathcal{Q}^{(d)}}^{(i)} &= \left( \tau \sum_{n=1}^N \sum_{\mathbf{p}-d} \text{vec}(\mathcal{Q}^{\neq d, \top}[\overline{\mathbf{p}-d}]) \text{vec}(\mathcal{Q}^{\neq d, \top}[\overline{\mathbf{p}-d}])^\top \right. \\ &\quad \left. + \text{diag}(\text{vec}(\boldsymbol{\phi}_{:,i}^{(d)})) * (\mathbf{U}^{(d+1)} \otimes \mathbf{U}^{(d)}) \right)^{-1}. \end{aligned}$$

where  $\mathbf{U}^{(d)} = \text{diag}(\mathbf{u}^{(d)})$ .

**Sample  $\phi$**  For  $d = 1, \dots, D$ ,  $j = 1, \dots, R_d$  and  $h = 1, \dots, R_{d+1}$ , the conditional distribution of  $\phi$  is

$$p(\phi_{jih}^{(d)} \mid -) \sim \text{Ga}(a_\phi, b_\phi),$$

where

$$a_\phi = \frac{1}{2} + \nu, \quad b_\phi = \nu + (\mathcal{Q}_{jh}^{(d)}[i])^2 u_j^{(d)} u_h^{(d+1)}.$$

**Sample  $\delta$**  We denote  $u_{l,-h}^{(d)} = \prod_{j=1, j \neq h}^l \delta_j^{(d)}$  and  $\mathbf{U}^{(d)} = \text{diag}(\mathbf{u}^{(d)})$ . Hence we have  $u_l^{(d)} = \prod_{j=1}^l \delta_j^{(d)} = u_{l,-h}^{(d)} \delta_h^{(d)}$ . For  $\boldsymbol{\delta}^{(1)}$ , the conditional distribution is

$$\log p(\delta_h^{(1)} \mid -) \propto \log p(\mathcal{Q}^{(1)} \mid \boldsymbol{\phi}^{(1)}, \mathbf{u}^{(1)}, \mathbf{u}^{(2)}) + \log p(\boldsymbol{\delta}^{(1)}), \quad \forall h = 1, \dots, R_1.$$

Hence, we have

$$p(\delta_h^{(1)} \mid -) \sim \text{Ga}(a_\delta^h, b_\delta^h),$$

where

$$a_\delta^h = \alpha_\delta + \frac{P_1 R_2 (R_1 - h + 1)}{2}, \quad b_\delta^h = 1 + \frac{1}{2} \sum_{l=h}^{R_1} \left[ \sum_{i=1}^{P_1} \mathbf{u}^{(2), \top} (\mathcal{Q}_l^{(1)}[i] * \mathcal{Q}_l^{(1)}[i]) \right] u_{l,-h}^{(1)}.$$

For  $d = 2, \dots, D$ , we have

$$\begin{aligned} p(\delta_h^{(d)} \mid -) &\propto \\ &\log p(\mathcal{Q}^{(d-1)} \mid \boldsymbol{\phi}^{(d-1)}, \mathbf{u}^{(d-1)}, \mathbf{u}^{(d)}) + \log p(\mathcal{Q}^{(d)} \mid \boldsymbol{\phi}^{(d)}, \mathbf{u}^{(d)}, \mathbf{u}^{(d+1)}) + \log p(\boldsymbol{\delta}^{(d)}) \end{aligned}$$

Then the conditional distribution is also a Gamma,

$$p(\delta_h^{(d)} \mid -) \sim \text{Ga}(a_\delta^h, b_\delta^h), \quad \forall h = 1, \dots, R_d,$$

where

$$a_\delta^h = \alpha_\delta + \frac{(P_d R_{d+1} + P_{d-1} R_d)(R_d - h + 1)}{2},$$

$$b_\delta^h = 1 + \frac{1}{2} \sum_{l=h}^{R_d} \left[ \sum_{i=1}^{P_d} \mathbf{u}^{(d+1),\top} (\mathbf{Q}_{i,l}^{(d)} [i] * \mathbf{Q}_{i,l}^{(d)} [i]) + \sum_{i=1}^{P_{d-1}} (\mathbf{Q}_{i,l}^{(d-1)} [i] * \mathbf{Q}_{i,l}^{(d-1)} [i]) \mathbf{u}^{(d-1)} \right] u_{l,-h}^{(d)}.$$

Similarly, for  $\delta^{(D+1)}$ , we have

$$p(\delta_h^{(D+1)} | -) \sim Ga(a_\delta^h, b_\delta^h), \quad \forall h = 1, \dots, R_{D+1},$$

where

$$a_\delta^h = \alpha_\delta + \frac{P_D R_{D+1} (R_{D+1} - h + 1)}{2}, \quad b_\delta^h = 1 + \frac{1}{2} \sum_{l=h}^{R_{D+1}} \left[ \sum_{i=1}^{P_D} (\mathbf{Q}_{i,l}^{(D)} [i] * \mathbf{Q}_{i,l}^{(D)} [i]) \mathbf{u}^{(d+1)} \right] u_{l,-h}^{(D+1)}.$$

**Sample  $\boldsymbol{\eta}$**  The conditional distribution is,

$$p(\boldsymbol{\eta}^{(n)} | -) \sim \mathcal{N}(\text{vec}(\boldsymbol{\eta}^{(n)}) | \boldsymbol{\mu}_\boldsymbol{\eta}^{(n)}, \boldsymbol{\Lambda}_\boldsymbol{\eta}^{(n)}), \quad (3)$$

where

$$\boldsymbol{\mu}_\boldsymbol{\eta}^{(n)} = \boldsymbol{\Lambda}_\boldsymbol{\eta}^{(n)} \left( \tau \sum_{\mathbf{p}} \mathcal{Y}_{\mathbf{p}}^{(n)} \text{vec}(\mathbf{Q}^{\leq D, \top} [\bar{\mathbf{p}}]) \right),$$

$$\boldsymbol{\Lambda}_\boldsymbol{\eta}^{(n)} = \left( \tau \sum_{\mathbf{p}} \text{vec}(\mathbf{Q}^{\leq D, \top} [\bar{\mathbf{p}}]) \text{vec}(\mathbf{Q}^{\leq D, \top} [\bar{\mathbf{p}}])^\top + \mathbf{I} \right)^{-1}.$$

**Sample  $\tau$**  The conditional distribution is  $p(\tau | -) \propto Ga(a_\tau, b_\tau)$ , where

$$a_\tau = \frac{N \prod_{d=1}^D P_d}{2} + \alpha_\tau, \quad b_\tau = \beta_\tau + \frac{1}{2} \sum_{n=1}^N \left\| \mathcal{Y}^{(n)} - \hat{\mathcal{Y}}^{(n)} \right\|_F^2.$$

## Appendix B. PX-EM Algorithm

In this section, we firstly derive the original EM algorithm. Then we describe the PX-EM algorithm by adding an additional double-rotation step. In the EM algorithm, we treat  $\boldsymbol{\eta}$  as latent variables and update them in the expectation step. Other parameters are updated in the maximization step. In this section, we denote the expectation w.r.t.  $\boldsymbol{\eta}$  as  $\langle \cdot \rangle$ .

### B.1. Expectation step

Based on Eq. (3), we can compute the expectation of the latent variables  $\boldsymbol{\eta}$ , as

$$\langle \tilde{\boldsymbol{\eta}}^{(n)} \rangle = \boldsymbol{\Lambda}_\boldsymbol{\eta}^{(n)} \left( \tau \sum_{\mathbf{p}} \mathcal{Y}_{\mathbf{p}}^{(n)} \text{vec}(\mathbf{Q}^{\leq D, \top} [\bar{\mathbf{p}}]) \right), \quad (4)$$

$$\langle \tilde{\boldsymbol{\eta}}^{(n)} \cdot (\tilde{\boldsymbol{\eta}}^{(n)})^\top \rangle = \langle \tilde{\boldsymbol{\eta}}^{(n)} \rangle \cdot \langle \tilde{\boldsymbol{\eta}}^{(n)} \rangle^\top + \boldsymbol{\Lambda}_\boldsymbol{\eta}^{(n)}, \quad (5)$$

for  $n = 1, \dots, N$ , where  $\tilde{\boldsymbol{\eta}}^{(n)} = \text{vec}(\boldsymbol{\eta}^{(n)})$  and

$$\boldsymbol{\Lambda}_{\boldsymbol{\eta}}^{(n)} = \left( \tau \sum_{\mathbf{p}} \text{vec}(\mathbf{Q}^{\leq D, \top}[\bar{\mathbf{p}}]) \text{vec}(\mathbf{Q}^{\leq D, \top}[\bar{\mathbf{p}}])^\top + \mathbf{I} \right)^{-1}.$$

The computational complexity is  $\mathcal{O}(K^3 + \sum_{d=1}^D P_d R^4 + NK^2 + NPD R^3)$ , where  $K$  is the factor number,  $D$  is the data order,  $N$  is the sample size,  $P = \prod_{d=1}^D P_d$  is the feature length and  $R_1 = \dots = R_D = R$  is the TR ranks. However, if  $\boldsymbol{\mathcal{Y}}$  admits a TR format of rank  $r$ , the inner product in Eq. (4) can be computed polynomially with  $D$  and the computational complexity reduces to  $\mathcal{O}(K^3 + \sum_{d=1}^D P_d R^4 + NK^2 + N \sum_{d=1}^D P_d R^2 r^2)$ .

## B.2. Maximization step

In the maximization step we give MAP estimate of the rest paramters. These estimates can be derived from the conditional distributions by replacing  $\boldsymbol{\eta}$  by its expectations in Eq. (4) and (5).

**Update  $\mathbf{Q}$**  For the core tensors  $\mathbf{Q}^{(d)}, \forall d = 1, \dots, D$ , according the conditional posterior Eq. (2), we have the update rule

$$\text{vec}(\mathbf{Q}^{(d)}[i]) = \langle \boldsymbol{\Lambda}_Q^{-1} \rangle \left( \tau \sum_{n=1}^N \sum_{\mathbf{p}-d} \boldsymbol{\mathcal{Y}}_{\mathbf{p}-d}^{(n)} \langle \text{vec}(\mathbf{Q}^{\neq d, \top}[\bar{\mathbf{p}}-d]) \rangle \right) \tilde{\mathbf{q}}^{(d)}[i], \quad \forall i = 1, \dots, P_d, \quad (6)$$

where

$$\boldsymbol{\Lambda}_Q = \left( \tau \sum_{n=1}^N \sum_{\mathbf{p}-d} \text{vec}(\mathbf{Q}^{\neq d, \top}[\bar{\mathbf{p}}-d]) \text{vec}(\mathbf{Q}^{\neq d, \top}[\bar{\mathbf{p}}-d])^\top + \text{diag}(\text{vec}(\boldsymbol{\phi}_{:,i}^{(d)})) * (\text{diag}(\mathbf{u}^{(d+1)}) \otimes \text{diag}(\mathbf{u}^{(d)})) \right)^{-1}.$$

The computation of  $\langle \text{vec}(\mathbf{Q}^{\neq d, \top}[\bar{\mathbf{p}}-d]) \rangle$  is straightforward. Then, we have

$$\begin{aligned} & \left\langle \sum_{\mathbf{p}} \text{vec}(\mathbf{Q}^{\neq d, \top}[\bar{\mathbf{p}}]) \text{vec}(\mathbf{Q}^{\neq d, \top}[\bar{\mathbf{p}}])^\top \right\rangle \\ &= \sum_{p_{d-1}} \mathbf{Q}_{r_{d-2}m}^{(d_1)}[p_{d-1}] \mathbf{Q}_{r'_{d-2}n}^{(d-1)}[p_{d-1}] \dots \sum_{n=1}^N \left\langle \text{vec}(\boldsymbol{\eta}^{(n), \top}) \text{vec}(\boldsymbol{\eta}^{(n), \top})^\top \right\rangle \dots \sum_{p_1} \mathbf{Q}_{r_1 i}^{(1)}[p_1] \mathbf{Q}_{r'_1 j}^{(1)}[p_1], \end{aligned}$$

where the expectation term is given by Eq. (5).

The computational complexity for all  $\{\mathbf{Q}^{(d)}\}_{d=1}^D$  is  $\mathcal{O}(\sum_{d=1}^D P_d R^6 + N \sum_{d=1}^D P_d R^4 + NPD^2 R^3)$ . Similarly, if  $\boldsymbol{\mathcal{Y}}$  admits TR format of rank  $r$ , the complexity cab be reduced to  $\mathcal{O}(\sum_{d=1}^D P_d R^6 + N \sum_{d=1}^D P_d R^4 + ND \sum_{d=1}^D P_d R^2 r^2)$ .

**Update  $\boldsymbol{\phi}$**  The update rule for  $\boldsymbol{\phi}$  is

$$\phi_{jih}^{(d)} = \frac{\frac{1}{2} + \nu}{\nu + (\mathbf{Q}_{jh}^{(d)}[i])^2 u_j^{(d)} u_h^{(d+1)}}. \quad (7)$$

**Update  $\delta$**  For  $\delta^{(1)}$ , the update rule is

$$\delta_h^{(1)} = \frac{\alpha_\delta + P_1 R_2 (R_1 - h + 1)/2}{1 + \frac{1}{2} \sum_{l=h}^{R_1} \left[ \sum_{i=1}^{P_1} \mathbf{u}^{(2),\top}(\mathbf{Q}_l^{(1)}[i] * \mathbf{Q}_l^{(1)}[i]) \right] u_{l,-h}^{(1)}}. \quad (8)$$

For  $d = 2, \dots, D$ , the update rule is

$$\delta_h^{(d)} = \frac{\alpha_\delta + (P_d R_{d+1} + P_{d-1} R_d)(R_d - h + 1)/2}{1 + \frac{1}{2} \sum_{l=h}^{R_d} \left[ \sum_{i=1}^{P_d} \mathbf{u}^{(d+1),\top}(\mathbf{Q}_l^{(d)}[i] * \mathbf{Q}_l^{(d)}[i]) + \sum_{i=1}^{P_{d-1}} (\mathbf{Q}_l^{(d-1)}[i] * \mathbf{Q}_l^{(d-1)}[i]) \mathbf{u}^{(d-1)} \right] u_{l,-h}^{(d)}}. \quad (9)$$

For  $\delta^{(D+1)}$ , we have

$$\delta_h^{(D+1)} = \frac{\alpha_\delta + P_D R_{D+1} (R_{D+1} - h + 1)/2}{1 + \frac{1}{2} \sum_{l=h}^{R_{D+1}} \left[ \sum_{i=1}^{P_D} (\mathbf{Q}_l^{(D)}[i] * \mathbf{Q}_l^{(D)}[i]) \mathbf{u}^{(d+1)} \right] u_{l,-h}^{(D+1)}}. \quad (10)$$

**Update  $\tau$**  For the noise precision, the update rule is

$$\tau_{\mathbf{p}} = \frac{NP_1 \cdots P_D / 2 + \alpha_\tau}{\beta_\tau + \sum_{n=1}^N \langle \left\| \mathbf{y}_{\mathbf{p}}^{(n)} - \hat{\mathbf{y}}_{\mathbf{p}}^{(n)} \right\|_F^2 \rangle / 2}. \quad (11)$$

The expectation term can be computed using the similar strategy with updating  $\mathcal{Q}$ .

### B.3. Rotation step

For the parameter-expanded version of the TRLF model, we have

$$\begin{aligned} \mathbf{y} &= \ll \tilde{\mathbf{Q}}^{(1)}, \mathbf{Q}^{(2)}, \dots, \mathbf{Q}^{(D-1)}, \tilde{\mathbf{Q}}^{(D)}, \boldsymbol{\eta} \gg + \boldsymbol{\mathcal{E}}, \\ \boldsymbol{\eta}^{(n)} &\sim \mathcal{MN}(\mathbf{0}, \mathbf{A}, \mathbf{B}), \end{aligned}$$

where

$$\begin{aligned} \tilde{\mathbf{Q}}^{(1)}[i] &= \mathbf{R}_1^{-1} \mathbf{Q}^{(1)}[i], \quad \forall i = 1, \dots, P_1, \\ \tilde{\mathbf{Q}}^{(D)}[i] &= \mathbf{Q}^{(D)}[i] \mathbf{R}_2^{-1}, \quad \forall i = 1, \dots, P_D. \end{aligned}$$

In the PX-EM algorithm, we add MGP priors on the loading core tensors  $\tilde{\mathbf{Q}}^{(1)}$  and  $\tilde{\mathbf{Q}}^{(D)}$ , instead of  $\mathbf{Q}^{(1)}$  and  $\mathbf{Q}^{(D)}$ . Moreover, by setting  $\mathbf{R}_1$  as the upper triangular part of the Cholesky decomposition of  $\mathbf{B}$  and  $\mathbf{R}_2$  as the lower triangular part of the Cholesky decomposition of  $\mathbf{A}$ , it is equivalent to the original TRLF model. The key part is to find the optimal rotation matrix  $\mathbf{R}_1$  and  $\mathbf{R}_2$ , based on  $\mathbf{A}$  and  $\mathbf{B}$ . Like other parameters, the optimal rotation matrices are also obtained by maximizing expectation of the log-likelihood function. In specific, we have

$$\mathbf{A}_*, \mathbf{B}_* = \arg \max_{\mathbf{A}, \mathbf{B}} \sum_{n=1}^N -\frac{R_1}{2} \log |\mathbf{A}| - \frac{R_D}{2} \log |\mathbf{B}| - \frac{1}{2} \text{tr}(\mathbf{A}^{-1} \boldsymbol{\eta}^{(n)} \mathbf{B}^{-1} \boldsymbol{\eta}^{(n),\top}). \quad (12)$$

Though optimization problem Eq. (12) is non-convex for  $\mathbf{A}$  and  $\mathbf{B}$ , we can update them in an alternating manner. Specifically, we have

$$\begin{aligned}
\mathbf{A}_* &= \arg \max_{\mathbf{A}} \sum_{n=1}^N -\frac{R_1}{2} \log |\mathbf{A}| - \frac{1}{2} \langle \text{tr}(\mathbf{A}^{-1} \boldsymbol{\eta}^{(n)} \mathbf{B}^{-1} \boldsymbol{\eta}^{(n),\top}) \rangle \\
&= \langle \sum_{n=1}^N \boldsymbol{\eta}^{(n),\top} \mathbf{B}^{-1} \boldsymbol{\eta}^{(n)} \rangle / (NR_1), \\
\mathbf{B}_* &= \arg \max_{\mathbf{B}} \sum_{n=1}^N -\frac{R_D}{2} \log |\mathbf{B}| - \frac{1}{2} \langle \text{tr}(\mathbf{A}^{-1} \boldsymbol{\eta}^{(n)} \mathbf{B}^{-1} \boldsymbol{\eta}^{(n),\top}) \rangle \\
&= \langle \sum_{n=1}^N \boldsymbol{\eta}^{(n)} \mathbf{A}^{-1} \boldsymbol{\eta}^{(n),\top} \rangle / (NR_D).
\end{aligned} \tag{13}$$

In each iteration of the PX-EM algorithm, we only update  $\mathbf{A}$  and  $\mathbf{B}$  once to get a local solution. It should be noted that if we use the TT format (i.e., TTLF),  $\mathbf{A}$  or  $\mathbf{B}$  is fixed as an identity matrix. Hence, the update rule for TTLF is optimal. The rotation step is used to get sparse latent factors and is optional. In practice, we can add this step for the first several iterations. The PX-EM algorithm is summarized in Algorithm 1.

## Appendix C. Experiments

### C.1. Synthetic Analysis

For the POET model<sup>1</sup>, we use factor number 50. For the InfLF model<sup>2</sup>, we initialize the factor number as 30 and adaptively tune the factor number in the Gibbs sampling process. For the HOLQ model<sup>3</sup>, we use the default settings. For our model, we initialize with TT-rank 30 and prune small factors during training. The TTLF is a special case of TRLF in that the first rank is 1. In the simulation study, the TRLF has similar log-Euclidean distance (LED) with TTLF.

### C.2. Real Data Covariance Estimation

For high dimensional data, if the covariance matrix has no special structure, it is hard to compute the log-likelihood function. Hence, we do not compare with the LW method in this experiment. Moreover, we find that the POET model is hard to fit these data. So we compare our model with InfLF and HOLQ in this subsection. For InfLF and our model, we compute the log-likelihood by using the Woodbury identity. For HOLQ, since it has Kronecker structure, it is also very efficient to compute the matrix determinant and matrix inverse. For all the data, we initialize the InfLF model with factor number 40 and our model with TT-rank 40.

---

1. <https://cran.r-project.org/web/packages/POET/>  
2. <https://cran.r-project.org/web/packages/infinitefactor/>  
3. <https://cran.r-project.org/web/packages/tensr/index.html>

---

**Algorithm 1** PX-EM algorithm for TRLF.

---

**Input:** Data  $\mathcal{Y}$ , hyper-parameters  $\alpha_\delta, \beta_\delta, \alpha_\tau, \beta_\tau, \nu$ .  
**Initialize:** Estimator  $\hat{\mathbf{Y}} = \ll \mathbf{Q}^{(1)}, \dots, \mathbf{Q}^{(D)}, \boldsymbol{\eta} \gg$ .  
**Output:**  $\hat{\mathbf{Y}}$ .

```

repeat
  // E-step
  for  $n = 1 : N$  do
    Update  $\boldsymbol{\eta}^{(n)}$  by Eq. (4) and (5).
  end for
  // M-step
  for  $d = 1 : D$  do
    Update  $\mathbf{Q}^{(d)}$  by Eq. (6).
  end for
  for  $d = 1 : D$  do
    Update  $\boldsymbol{\phi}^{(d)}$  by Eq. (7).
  end for
  for  $d = 1 : D + 1$  do
    Update  $\boldsymbol{\delta}^{(d)}$  by Eq. (8), (9) and (10) respectively.
  end for
  Update  $\boldsymbol{\tau}$  by Eq. (11).
  // Rotation (optional)
  Compute  $\mathbf{A}$  and  $\mathbf{B}$  by Eq. (13) and the corresponding Cholesky decompositions.
  for  $i = 1 : P_1$  do
    Set  $\mathbf{Q}^{(1)}[i] := \mathbf{R}_1 \mathbf{Q}^{(1)}[i]$ .
  end for
  for  $i = 1 : P_D$  do
    Set  $\mathbf{Q}^{(D)}[i] := \mathbf{Q}^{(D)}[i] \mathbf{R}_2$ .
  end for
until Maximum iteration or convergence.

```

---

### C.3. Image Inpainting

The BCPF model also uses sparse Bayesian priors to automatically choose latent factors and requires no hyperparameters. However, for the rest of models, we have to carefully tune the hyperparameters, e.g., the TT/TR-ranks. Here we compute those models under several TT/TR-rank settings and select the best performance. However, it should be noted that this is not realistic in many applications where the true signals are unknown. For our model, we simply initialize with TR format of rank 15, and we truncate the redundant factors while training, in order to reduce computational cost. For the image inpainting task, if we use the tensor train decomposition, it may require much larger TT-ranks to get comparable performance. So we only use tensor ring decomposition in this experiment.

We choose 8 standard pictures as benchmark, as shown in Figure 1. We randomly generate masks of missing rate 0.5, 0.7 and 0.9. Figure 2 shows the visualization results of the Lena image completion. Our model has better recovery results than the competitors, especially when the missing rate is high.

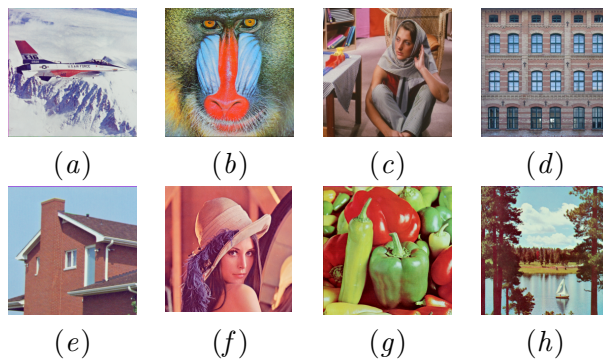


Figure 1: Benchmark images.

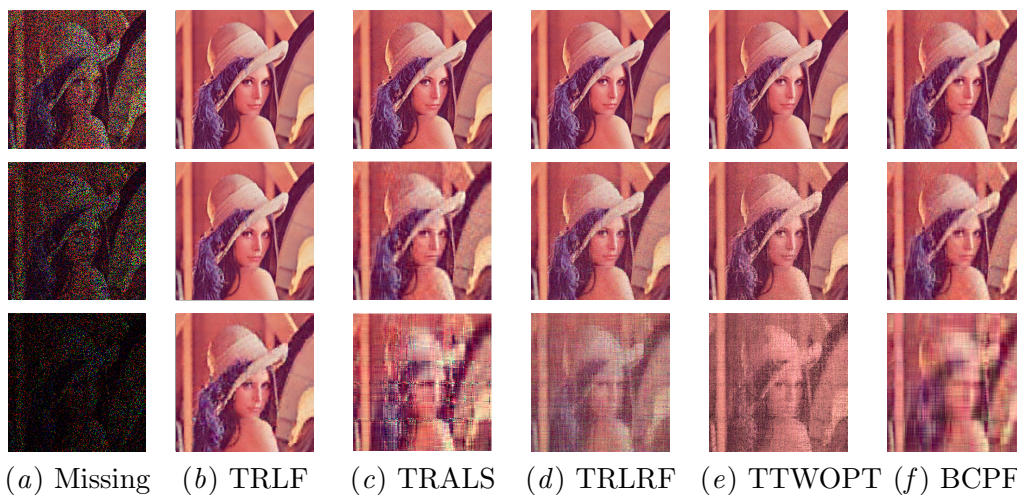


Figure 2: Illustration of the Lena image completion. Each row from the above to the bottom according to missing rate 0.5, 0.7 and 0.9 respectively.