

# Verifications

Jiaxin Hu

June 4, 2021

## 1 No iteration needed

Here we verify that unsupervised algorithm is enough to implement our supervised model numerically under the Gaussian data case. Recall our supervised model

$$\mathcal{Y} = \mathcal{B} \times \{\mathbf{X}_1, \dots, \mathbf{X}_K\} + \mathcal{E}, \quad (1)$$

where  $\mathcal{Y} \in \mathbb{R}^{d_1 \times \dots \times d_K}$ ,  $\mathbf{X}_k \in \mathbb{R}^{d_k \times p_k}$ , and  $\mathcal{E} \in \mathbb{R}^{d_1 \times \dots \times d_K}$  is the noise tensor with i.i.d. standard Gaussian entries. Consider the QR decomposition of the feature matrices, i.e.,  $\mathbf{X}_k = \mathbf{Q}_k \mathbf{R}_k$ , where  $\mathbf{Q}_k \in \mathbb{R}^{d_k \times p_k}$  is a matrix with orthogonal columns, and  $\mathbf{R}_k \in \mathbb{R}^{p_k \times p_k}$  is an upper-triangle matrix. Multiplying  $\mathbf{Q}_k^T$  on both sides of the model (1), we have

$$\mathcal{Y} \times \{\mathbf{Q}_1^T, \dots, \mathbf{Q}_K^T\} = \mathcal{B} \times \{\mathbf{R}_1, \dots, \mathbf{R}_K\} + \mathcal{E} \times \{\mathbf{Q}_1^T, \dots, \mathbf{Q}_K^T\}. \quad (2)$$

Let  $\mathcal{E}' = \mathcal{E} \times \{\mathbf{Q}_1^T, \dots, \mathbf{Q}_K^T\}$  denote the new noise tensor. By (Hoff et al., 2011; Li and Zhang, 2017), we know that  $\mathcal{E}'$  belongs to the class of random Gaussian tensor, whose possess a Kronecker covariance structure. Specifically, the covariance of vectorized  $\mathcal{E}'$  is

$$\text{Cov}(\text{vec}(\mathcal{E}')) = \mathbf{Q}_K^T \mathbf{Q}_K \otimes \dots \otimes \mathbf{Q}_1^T \mathbf{Q}_1 = \mathbf{I}_{\prod_k p_k \times \prod_k p_k}.$$

which

Applying unsupervised Tucker decomposition to the new observation  $\mathcal{Y} \times \{\mathbf{Q}_1^T, \dots, \mathbf{Q}_K^T\}$ , we obtain the low-rank estimate to the new coefficient  $\mathcal{B} \times \{\mathbf{R}_1, \dots, \mathbf{R}_K\}$ . Multiplying  $\hat{\mathbf{R}}_k = (\mathbf{R}_k^T \mathbf{R}_k)^{-1} \mathbf{R}_k^T$  to the coefficient estimate, we obtain the estimate of  $\mathcal{B}$ .

In general, unsupervised decomposition and supervised decomposition tackle different goals. For example, we can not fit an unsupervised decomposition for  $\mathcal{Y}$  in model 1 and then multiply the  $\mathbf{X}_k$  to the unsupervised estimate to obtain the estimate of  $\mathcal{B}$ . Intuitively, such estimated  $\mathcal{B}$  is not desirable since the we do not use the information from  $\mathbf{X}_k$  during the decomposition and the fitted tensor  $\hat{\mathcal{Y}}$  lie in the space jointly determined by  $\{\mathbf{I}_{d_1}, \dots, \mathbf{I}_{d_k}\}$ . However, in supervised decomposition, we restrict the fitted tensor  $\hat{\mathcal{Y}}$  in the space jointly determined by  $\{\mathbf{X}_1, \dots, \mathbf{X}_K\}$ . In model (2), the new observation tensor  $\mathcal{Y} \times \{\mathbf{Q}_1^T, \dots, \mathbf{Q}_K^T\}$  already includes the information from  $\mathbf{X}_k$  via  $\mathbf{Q}_k^T$ , since the column space  $C(\mathbf{X}_k) = C(\mathbf{Q}_k)$ . Additionally, the new noise remain i.i.d.. Then, the unsupervised decomposition works with model (2).

## 2 Robust to overparameterization

consider only

Here we verify that our model is robust to overparameterization. For simplicity, we only consider the case that side information is on the third mode and assume all the feature matrices have orthogonal columns. Recall the model

$$\mathcal{Y} = \mathcal{B} \times_3 \mathbf{X} + \mathcal{E}.$$

The estimate of  $\mathcal{B}$  satisfies

$$\hat{\mathcal{B}} = \arg \min_{\mathcal{B}' \text{ with rank } \mathbf{r}} \|\mathcal{B} + \mathcal{E} \times_3 \mathbf{X}^T - \mathcal{B}'\|_F^2.$$

Suppose we fit the data with feature matrix  $[\mathbf{X}, \tilde{\mathbf{X}}]$ , where  $\tilde{\mathbf{X}}$  has orthogonal columns with  $\mathbf{X}$ . The estimate of  $\mathcal{B}$  under this fit satisfies

$$\begin{aligned} \hat{\mathcal{B}}_{over} &= \arg \min_{\tilde{\mathcal{B}} \text{ with rank } \tilde{\mathbf{r}}} \left\| \mathcal{B} \times_3 \mathbf{X} \times_3 [\mathbf{X}, \tilde{\mathbf{X}}]^T + \mathcal{E} \times_3 [\mathbf{X}, \tilde{\mathbf{X}}]^T - \tilde{\mathcal{B}} \right\|_F^2 \\ &= \arg \min_{\tilde{\mathcal{B}} \text{ with rank } \tilde{\mathbf{r}}} \left\| \mathcal{B} \times_3 [\mathbf{I}, \mathbf{0}]^T + \mathcal{E} \times_3 [\mathbf{X}, \tilde{\mathbf{X}}]^T - \tilde{\mathcal{B}} \right\|_F^2. \end{aligned}$$

Hence, when noise is small enough and set the fitted rank  $\mathbf{r} = \tilde{\mathbf{r}}$ , we will obtain very similar  $\hat{\mathcal{B}}$  and  $\tilde{\mathcal{B}}$  on the first  $p$  slices. Besides, if we know the true rank of  $\mathcal{B}$  is  $\mathbf{r}_{true}$  and set  $\mathbf{r} \geq \mathbf{r}_{true}$ . The overparameterization rank  $\tilde{\mathbf{r}} \geq \mathbf{r}$  will also lead to similar estimates, since higher-rank decomposition can fit a low-rank signal well, and the fitted value does not change very much when noise is small. See Section 4.

### 3 Unsupervised v.s. Supervised decomposition

rephrase as a remark

Here we give a toy example in which unsupervised and supervised decomposition have orthogonal factor estimates. Let the core tensor  $\mathcal{C} \in \mathbb{R}^{4 \times 4 \times 4}$  be a superdiagonal tensor with elements  $(1, 2, 2, 2)$ , and factor matrices  $\mathbf{A} = \mathbf{B} = \mathbf{C} = \begin{bmatrix} \mathbf{I}_4 \\ \mathbf{0} \end{bmatrix} \in \mathbb{R}^{10 \times 4}$ . The observation tensor is generated as

$$\mathcal{Y} = \mathcal{C} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C} + \mathcal{E}.$$

Consider the fit with rank  $(1, 1, 1)$ . The unsupervised estimates of  $c, \mathbf{A}, \mathbf{B}, \mathbf{C}$  satisfies

$$(\hat{c}, \hat{\mathbf{A}}, \hat{\mathbf{B}}, \hat{\mathbf{C}}) = \arg \min_{c' \in \mathbb{R}, \mathbf{A}', \mathbf{B}', \mathbf{C}' \in \mathbb{R}^{10}} \|\mathcal{C} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C} + \mathcal{E} - c' \times_1 \mathbf{A}' \times_2 \mathbf{B}' \times_3 \mathbf{C}'\|_F^2. \quad (3)$$

Consider the supervision on the first mode  $\mathbf{X} = \mathbf{A}[1, 1]$ . By model (2), the estimate  $c, \mathbf{A}, \mathbf{B}, \mathbf{C}$  satisfies

$$\begin{aligned} (\tilde{c}, \tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \tilde{\mathbf{C}}) &= \arg \min_{c' \in \mathbb{R}, \mathbf{A}', \mathbf{B}', \mathbf{C}' \in \mathbb{R}^{10}} \|\mathcal{C} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C} \times_1 \mathbf{X}^T + \mathcal{E} \times_1 \mathbf{X}^T - c' \times_1 \mathbf{A}' \times_2 \mathbf{B}' \times_3 \mathbf{C}'\|_F^2 \\ &= \arg \min_{c' \in \mathbb{R}, \mathbf{A}', \mathbf{B}', \mathbf{C}' \in \mathbb{R}^{10}} \|\mathcal{C} \times_1 [1, 0, 0, 0] \times_2 \mathbf{B} \times_3 \mathbf{C} + \mathcal{E} \times_1 \mathbf{X}^T - c' \times_1 \mathbf{A}' \times_2 \mathbf{B}' \times_3 \mathbf{C}'\|_F^2. \end{aligned}$$

Suppose the noise is 0. In the unsupervised case, the observation in equation (3) is  $\mathcal{C} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C}$ , which is a sum of 4 rank-1 tensors with only one non-zero element on the super-diagonal. Last three rank-1 tensors has non-zero elements 2 and the first one has 1. Thus, the rank-1 estimate should have only one non-zero element with value 2 on the super-diagonal. This indicates that

$\hat{\mathbf{A}}, \hat{\mathbf{B}}, \hat{\mathbf{C}} \in \begin{bmatrix} 0 \\ \mathbf{I}_3 \\ \mathbf{0} \end{bmatrix}$ , where 0 denotes a zero row. In the supervised case, the “new” observation is

$\mathcal{C} \times_1 [1, 0, 0, 0] \times_2 \mathbf{B} \times_3 \mathbf{C}$ , which is a rank-1 tensor with the only non-zero element 1 on the position  $(1, 1, 1)$ . Therefore, the estimate  $\tilde{\mathbf{B}}, \tilde{\mathbf{C}} \in (1, 0, \dots, 0)^T$ .

Therefore, the unsupervised factors  $\hat{\mathbf{B}}, \hat{\mathbf{C}}$  are orthogonal to the supervised factors  $\tilde{\mathbf{B}}, \tilde{\mathbf{C}}$ .

## 4 Fit with higher rank

Here we verify that our model is robust to the fitting with higher rank than the rank of observation tensor. Note that our model can be reduced to an unsupervised decomposition in model 2. Without the loss of generality, we consider the order 3 tensor and only need to verify that a low-rank tensor  $\mathcal{Y}$  with rank  $(r, r, r)$  has infinitely many decompositions with rank  $(r, r, r_3)$  with the same fitted values, where  $r_3 > r$ .

Recall the tucker decomposition.

summarize # 4 in a proposition.

$$\mathcal{Y} = \mathcal{C} \times_1 \mathbf{M}_1 \times_2 \mathbf{M}_2 \times_3 \mathbf{M}_3,$$

where  $\mathcal{C} \in \mathbb{R}^{r \times r \times r}$  and  $\mathbf{M}_k$  has orthogonal columns. If we fit with rank  $(r, r, r_3)$ , we could always find a new core  $\tilde{\mathcal{C}} \in \mathbb{R}^{r \times r \times r_3}$  such that first  $r$  slices  $\tilde{\mathcal{C}}[:, 1 : r] = \mathcal{C}$  and last  $r_3 - r$  slices  $\tilde{\mathcal{C}}[:, (r+1) : r_3] = 0$ , which perfectly fits the data  $\mathcal{Y}$ . Correspondingly, the estimated factor matrix can be  $\tilde{\mathbf{M}}_3 = [\mathbf{M}_3, \mathbf{N}]$  for any  $\mathbf{N}$  has orthogonal columns with  $\mathbf{M}_3$ . This implies we have infinitely many decompositions with  $(r, r, r_3)$ . However, the fitted value would not change since we always have a perfect fit to the data  $\mathcal{Y}$ .

## References

- Hoff, P. D. et al. (2011). Separable covariance arrays via the tucker product, with applications to multivariate relational data. *Bayesian Analysis*, 6(2):179–196.
- Li, L. and Zhang, X. (2017). Parsimonious tensor response regression. *Journal of the American Statistical Association*, 112(519):1131–1146.