

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, which possesses a Kronecker covariance structure. Specifically, the covariance of vectorized \mathcal{E}' is

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

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 $\tilde{\mathbf{R}}_k = (\mathbf{R}_k^T \mathbf{R}_k)^{-1} \mathbf{R}_k^T$ 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 generalized inverse of \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

Here we verify that our model is robust to overparameterization. For simplicity, we consider only 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

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

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

Remark 1. Note that unsupervised and supervised decomposition tackles different goals. The unsupervised decomposition aims to describe the data variation in a low-rank structure while supervised decomposition explores the connection between data and the features. Intuitively, in unsupervised decomposition, the data \mathcal{Y} is projected to the space jointly determined by $\{\mathbf{I}_{d_1}, \dots, \mathbf{I}_{d_K}\}$; in supervised decomposition, the data \mathcal{Y} is projected to a specific subspace determined by $\{\mathbf{X}_1, \dots, \mathbf{X}_K\}$. The fitted values \mathcal{Y} with and without supervision may lie in orthogonal spaces.

Example 1. Consider the noiseless order-3 tensor observation $\mathcal{Y} \in \mathbb{R}^{10 \times 10 \times 10}$. We generate the observation as

$$\mathcal{Y} = \mathcal{C} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C},$$

where $\mathcal{C} \in \mathbb{R}^{4 \times 4 \times 4}$ is a superdiagonal tensor with elements (1, 2, 2, 2), and factor matrices $\mathbf{A} = \mathbf{B} = \mathbf{C} = \begin{bmatrix} \mathbf{I}_4 \\ \mathbf{0} \end{bmatrix}$. Fitting with rank (1, 1, 1), the unsupervised estimates of $\mathcal{C}, \mathbf{A}, \mathbf{B}, \mathbf{C}$ satisfies

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

Suppose we have a supervision on the first mode $\mathbf{X} = \mathbf{A}[1]$. By model (2), the estimates $\mathcal{C}, \mathbf{A}, \mathbf{B}, \mathbf{C}$ satisfies

$$\begin{aligned} (\tilde{\mathcal{C}}, \tilde{\mathbf{A}}, \tilde{\mathbf{B}}, \tilde{\mathbf{C}}) &= \arg \min_{\mathcal{C}' \in \mathbb{R}, \mathbf{A}', \mathbf{B}', \mathbf{C}' \in \mathbb{R}^{10}} \|\mathcal{C} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C} \times_1 \mathbf{X}^T - \mathcal{C}' \times_1 \mathbf{A}' \times_2 \mathbf{B}' \times_3 \mathbf{C}'\|_F^2 \\ &= \arg \min_{\mathcal{C}' \in \mathbb{R}, \mathbf{A}', \mathbf{B}', \mathbf{C}' \in \mathbb{R}^{10}} \|\mathcal{C} \times_1 [1, 0, 0, 0] \times_2 \mathbf{B} \times_3 \mathbf{C} - \mathcal{C}' \times_1 \mathbf{A}' \times_2 \mathbf{B}' \times_3 \mathbf{C}'\|_F^2. \end{aligned} \quad (4)$$

In the unsupervised case (3), the observation $\mathcal{C} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C}$ is a sum of 4 rank-1 tensors with only one non-zero element on the super-diagonal. Last three rank-1 tensors have non-zero element 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 (4), the “new” observation $\mathcal{C} \times_1 [1, 0, 0, 0] \times_2 \mathbf{B} \times_3 \mathbf{C}$ is a rank-1 tensor with the only non-zero element 1 on the position (1, 1, 1). Thus, the estimates $\tilde{\mathbf{B}}, \tilde{\mathbf{C}} \in (1, 0, \dots, 0)^T$, which are orthogonal to the unsupervised factors $\hat{\mathbf{B}}, \hat{\mathbf{C}}$.

In practice, we have noisy observations. If the noise is small, this phenomenon still holds. If the noise is big and comparable with the signal, the estimations are unstable and we fail to discuss their properties.

4 Fit with higher rank

Proposition 1. *Let \mathcal{Y} be a low-rank tensor with rank (r, r, r) . There exist infinitely many decompositions of $\mathcal{Y} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$ with core shape (r, r, r_3) , where $r_3 > r$, and the factor matrices on first two modes are unique up to orthogonal transformation. Also, the fitted values are the same with these two kinds of decompositions.*

Proof. The decomposition of \mathcal{Y} is

$$\mathcal{Y} = \mathcal{C} \times_1 \mathbf{M}_1 \times_2 \mathbf{M}_2 \times_3 \mathbf{M}_3, \quad (5)$$

where $\mathcal{C} \in \mathbb{R}^{r \times r \times r}$ is the core tensor and $\mathbf{M}_k \in \mathbb{R}^{d_k \times r_k}, k \in [3]$ has orthogonal columns. Consider a core tensor $\tilde{\mathcal{C}} \in \mathbb{R}^{r \times r \times r_3}$ such that first r slices $\tilde{\mathcal{C}}_{..i} = \mathcal{C}_{..i}, i \in [r]$ and the last $r_3 - r$ slices $\tilde{\mathcal{C}}_{..i} = 0, i = r+1, \dots, r_3$. Correspondingly, let $\tilde{\mathbf{M}}_3 = [\mathbf{M}_3, \mathbf{N}]$, where $\mathbf{N} \in \mathbb{R}^{d_3 \times (r_3 - r)}$ has orthogonal columns with \mathbf{M}_3 . Then, we have the decomposition

$$\mathcal{Y} = \tilde{\mathcal{C}} \times_1 \tilde{\mathbf{M}}_1 \times_2 \tilde{\mathbf{M}}_2 \times_3 \tilde{\mathbf{M}}_3. \quad (6)$$

Note the decomposition holds for infinitely many \mathbf{N} . By setting $\tilde{\mathbf{M}}_1 = \mathbf{M}_1$ and $\tilde{\mathbf{M}}_2 = \mathbf{M}_2$, we have infinitely many decomposition of \mathcal{Y} with core shape (r, r, r_3) .

Notice that the column space $C(\text{Unfold}_1(\mathcal{Y})) = C(\mathbf{M}_1)$ with dimension r . Then, in the decomposition with core shape (r, r, r_3) , we should have $C(\tilde{\mathbf{M}}_1) = C(\mathbf{M}_1)$ with dimension r . Therefore, $\tilde{\mathbf{M}}_1$ is equal to \mathbf{M}_1 up to orthogonal transformations. Similar to the factor $\tilde{\mathbf{M}}_2$ and \mathbf{M}_2 .

Since both of the original decomposition (5) and the degenerate decomposition (6) recover \mathcal{Y} exactly, the fitted values $\mathcal{C} \times_1 \mathbf{M}_1 \times_2 \mathbf{M}_2 \times_3 \mathbf{M}_3$ and $\tilde{\mathcal{C}} \times_1 \tilde{\mathbf{M}}_1 \times_2 \tilde{\mathbf{M}}_2 \times_3 \tilde{\mathbf{M}}_3$ are identical. \square

Remark 2. In practice, we always have noisy observations. Since the noise tensor is usually full rank, our observation \mathcal{Y} is usually high-rank, and we will obtain a full rank core tensor \mathcal{C} . However, if the noise is too small, the estimate of core tensor \mathcal{C} becomes unstable and has an approximate low-rank structure.

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.