

# Thought about SupCP

Jiixin Hu

May 21, 2021

## 1 SupCP performance

In simulation, SupCP performs better as the number of available feature matrices increases, which is counterintuitive. A possible reason to this phenomena is that the valid space for the fitted values of STD becomes smaller due to more “supervision”. Here we only consider the Gaussian data. Let  $\mathcal{Y} \in \mathbb{R}^{d \times d \times d}$ ,  $\mathbf{X}_k \in \mathbb{R}^{d \times p}$ ,  $k \in [3]$ . Consider the tucker rank  $\mathbf{r} = (3, 3, 3)$  and CP rank  $R$ . The dimension of  $\mathbf{M}_k$  can be obtained by the context.

**First**, we start with the unsupervised case. Recall the STD and SupCP model.

$$\begin{aligned} STD &: \mathcal{Y} = \mathcal{C} \times \{\mathbf{M}_1, \mathbf{M}_2, \mathbf{M}_3\} + \mathcal{E} \\ SupCP &: \mathcal{Y} = \llbracket \mathbf{A}_1, \mathbf{A}_2, \mathbf{A}_3 \rrbracket + \mathcal{E}. \end{aligned}$$

The estimation problem of  $\mathcal{Y}$  in STD and SupCP can be formulated in the same form

$$\min_{\text{vec}(\hat{\mathcal{Y}}) \in \mathcal{P}} \left\| \mathcal{Y} - \hat{\mathcal{Y}} \right\|_F^2,$$

where the space  $\mathcal{P}$  depends on the model structure. Particularly,

$$\begin{aligned} \mathcal{P}_{STD} &= \{[\langle \mathcal{C}, \mathbf{M}_{1i} \circ \mathbf{M}_{2j} \circ \mathbf{M}_{3k} \rangle], i, j, k \in [d] \mid \mathcal{C} \in \mathbb{R}^{r \times r \times r}, \mathbf{M}_k \in \mathbb{R}^{d \times r}, \mathbf{M}_k^T \mathbf{M}_k = \mathbf{I}_r\} \\ &= \{\text{span}(\mathbf{M}_1 \otimes \mathbf{M}_2 \otimes \mathbf{M}_3) \mid \mathbf{M}_k \in \mathbb{R}^{d \times r}, \mathbf{M}_k^T \mathbf{M}_k = \mathbf{I}_r\}, \end{aligned} \quad \text{what does it mean?}$$

where the second equation holds since  $\mathcal{C}$  is free, and

CP factor has no orthogonality constraints.

$$\begin{aligned} \mathcal{P}_{SupCP} &= \{[\langle \mathcal{I}_R, \mathbf{A}_{1i} \circ \mathbf{A}_{2j} \circ \mathbf{A}_{3k} \rangle], i, j, k \in [d] \mid \mathbf{A}_k \in \mathbb{R}^{d \times R}, \mathbf{A}_k^T \mathbf{A}_k = \mathbf{I}_R\} \\ &= \{\text{span}([\mathbf{A}_1 \odot \mathbf{A}_2 \odot \mathbf{A}_3] \mathbf{1}_R) \mid \mathbf{A}_k \in \mathbb{R}^{d \times R}, \mathbf{A}_k^T \mathbf{A}_k = \mathbf{I}_R\} \end{aligned}$$

where  $\mathcal{I}_R$  is the identity tensor with entries  $\mathbf{I}_{R,iii} = 1, i \in [R]$  and other entries 0, and  $\odot$  is the Khatri-Rao product (column-wise Kronecker product).

Note that the space  $\mathcal{P}_{SupCP} \subset \mathcal{P}_{STD}$  if  $R < r^3$  since CP decomposition is a special case of Tucker decomposition. Therefore, if the true signal tensor is generated from Tucker decomposition, our Tucker-based STD outperforms CP-based SupCP.

how about  $R \gg r^3$       verify in experiment

**Next**, we use the same idea above for the supervised case, i.e., compare the space  $\mathcal{P}_{SupCP}, \mathcal{P}_{STD}$  with feature matrix. For simplicity, we ignore the extra noise  $\mathcal{E}'$  in SupCP and consider the completely supervised  $\mathbf{A}_1 = \mathbf{X}_1 \mathbf{B}$ .

Now we have the spaces

$$\begin{aligned}\mathcal{P}_{STD} &= \{[\langle \mathcal{C}, (\mathbf{X}_1 \mathbf{M}_1)_i \circ (\mathbf{X}_2 \mathbf{M}_2)_j \circ (\mathbf{X}_3 \mathbf{M}_3)_k \rangle], i, j, k \in [d] \mid \mathcal{C} \in \mathbb{R}^{r \times r \times r}, \mathbf{M}_k \in \mathbb{R}^{p \times r}, \mathbf{M}_k^T \mathbf{M}_k = r\} \\ &= \{\text{span}(\mathbf{X}_1 \mathbf{M}_1 \otimes \mathbf{X}_2 \mathbf{M}_2 \otimes \mathbf{X}_3 \mathbf{M}_3) \mid \mathbf{M}_k \in \mathbb{R}^{p \times r}, \mathbf{M}_k^T \mathbf{M}_k = r\}, \quad \text{subset } \mathbb{R}^{\text{?}}\end{aligned}$$

and

$$\begin{aligned}\mathcal{P}_{SupCP} &= \{[\langle \mathcal{I}_R, (\mathbf{X}_1 \mathbf{B})_i \circ \mathbf{A}_{2j} \circ \mathbf{A}_{3k} \rangle], i, j, k \in [d] \mid \mathbf{B} \in \mathbb{R}^{p \times R}, \mathbf{A}_k \in \mathbb{R}^{d \times R}, \mathbf{A}_k^T \mathbf{A}_k = R\} \\ &= \{\text{span}([\mathbf{X}_1 \mathbf{B} \odot \mathbf{A}_2 \odot \mathbf{A}_3] 1_R) \mid \mathbf{B} \in \mathbb{R}^{p \times R}, \mathbf{A}_k \in \mathbb{R}^{d \times R}, \mathbf{A}_k^T \mathbf{A}_k = R\}, \quad \text{Text}\end{aligned}$$

where  $\mathbf{X}_k$  can be  $\mathbf{I}_d$  if no feature matrix is available on  $k$ -th mode.

What do you mean by “subset”? First space is a vector space of dimension  $d^3$ . Second space is a vectors space of dimension  $pd^2$ . “subset” is only meaningful for comparable spaces. No sense to say  $(1,2)$  subset  $(1,1,1)$ . Consider  $\mathbf{X}_2, \mathbf{X}_3 = \mathbf{I}_d$ . If  $R \leq r$ , then  $\mathcal{P}_{SupCP} \subset \mathcal{P}_{STD}$ ; if  $R > r$ , there exists a vector  $v \in \mathcal{P}_{SupCP} \setminus \mathcal{P}_{STD}$ , since  $\text{span}(\mathbf{X}_1 \mathbf{B}) \not\subset \text{span}(\mathbf{X}_1 \mathbf{M}_1)$ , for all full rank  $\mathbf{M}_1 \in \mathbb{R}^{d \times r}$  and full rank  $\mathbf{B} \in \mathbb{R}^{d \times R}$  under this case.

2. Let  $\mathbf{X}_2, \mathbf{X}_3 \neq \mathbf{I}_d$ . Note that

$$(\mathbf{X}_1 \mathbf{M}_1 \otimes \mathbf{X}_2 \mathbf{M}_2 \otimes \mathbf{X}_3 \mathbf{M}_3) = (\mathbf{X}_1 \otimes \mathbf{X}_2 \otimes \mathbf{X}_3) (\mathbf{M}_1 \otimes \mathbf{M}_2 \otimes \mathbf{M}_3). \quad (*)$$

Then  $\text{span}(\mathbf{X}_1 \mathbf{M}_1 \otimes \mathbf{X}_2 \mathbf{M}_2 \otimes \mathbf{X}_3 \mathbf{M}_3) \subset \text{span}(\mathbf{X}_1 \otimes \mathbf{X}_2 \otimes \mathbf{X}_3)$ , and thus  $\mathcal{P}_{STD}$  becomes smaller when feature matrices are available on more modes and each matrix contains fewer information. I do not under the logic. Both sides (\*) have X. Specify which is for # feature = 3, which is

for # feature =1

Meanwhile, no matter  $R \leq r$  or  $R > r$ , it is possible that  $\mathcal{P}_{SupCP} \not\subset \mathcal{P}_{STD}$ , since there is no supervised constraints on the last two factors  $\mathbf{A}_2, \mathbf{A}_3$ . Have you verified from experiment? Generate from CP then fit two models?

3. Therefore, there always exists model misspecification if the true signal is generated from the STD model. That's the reason why STD outperforms than SupCP.
4. When there is no information or only 1 feature matrix available with  $R \leq r$ ,  $\mathcal{P}_{SupCP}$  is a subset of  $\mathcal{P}_{STD}$  and it is likely that the true signal does not fall in the space  $\mathcal{P}_{SupCP}$ . When there is more information,  $\mathcal{P}_{STD}$  becomes smaller, and it is more likely for  $\mathcal{P}_{SupCP}$  to cover the true signal, particularly when  $R > r$ . An extreme example is  $r = 1, R \geq 2$ . Another example is  $r = 2, R \geq 8$  and  $\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3$  are available. In both examples,  $\mathcal{P}_{STD}$  is a subset of  $\mathcal{P}_{SupCP}$ .

## References