
Supplements for “Multiway clustering via tensor block models”

A Proofs

A.1 Stochastic tensor block model

The following property shows that Bernoulli distribution belongs to the sub-Gaussian family with a subgaussianity parameter σ equal to $1/4$.

Property 1. Suppose $x \sim \text{Bernoulli}(p)$, then $x \sim \text{sub-Gaussian}(\frac{1}{4})$.

Proof. For all $\lambda \in \mathbb{R}$, we have

$$\ln(\mathbb{E}(e^{\lambda(x-p)})) = \ln\left(pe^{\lambda(1-p)} + (1-p)e^{-p\lambda}\right) = -p\lambda + \ln(1 + pe^\lambda - p) \leq \frac{\lambda^2}{8}.$$

Therefore $\mathbb{E}(e^{\lambda(x-p)}) \leq e^{\lambda^2(1/4)/2}$. □

A.2 Proof of Proposition 1

Proof. Let \mathbb{P}_Θ denotes the (either Gaussian or Bernoulli) tensor block model, where $\Theta = \mathcal{C} \times_1 \mathbf{M}_1 \times_2 \cdots \times_K \mathbf{M}_K$ parameterizes the mean tensor. Since the mapping $\Theta \mapsto \mathbb{P}_\Theta$ is one-to-one, Θ is identifiable. Now suppose that Θ can be decomposed in two ways, $\Theta = \Theta(\{\mathbf{M}_k\}, \mathcal{C}) = \Theta(\{\tilde{\mathbf{M}}_k\}, \tilde{\mathcal{C}})$. Based on the Assumption 1, we have

$$\Theta = \mathcal{C} \times_1 \mathbf{M}_1 \times_2 \cdots \times_K \mathbf{M}_K = \tilde{\mathcal{C}} \times_1 \tilde{\mathbf{M}}_1 \times_2 \cdots \times_K \tilde{\mathbf{M}}_K, \quad (1)$$

where $\mathcal{C}, \tilde{\mathcal{C}} \in \mathbb{R}^{R_1 \times \cdots \times R_K}$ are two irreducible cores, and $\mathbf{M}_k, \tilde{\mathbf{M}}_k \in \{0, 1\}^{R_k \times d_k}$ are membership matrices for all $k \in [K]$. We will prove by contradiction that \mathbf{M}_k and $\tilde{\mathbf{M}}_k$ induce the same partition of $[d_k]$, for all $k \in [K]$.

Suppose the above claim does not hold. Then there exists a mode $k \in [K]$ such that the $\mathbf{M}_k, \tilde{\mathbf{M}}_k$ induce two different partitions of $[d_k]$. Without loss of generality, we assume $k = 1$. The definition of partition implies that there exists a pair of indices $i \neq j, i, j \in [d_1]$, such that, i, j belong to the same cluster based on \mathbf{M}_1 , but they belong to different clusters based on $\tilde{\mathbf{M}}_1$. Let $\mathcal{A} \neq \mathcal{B}, \mathcal{A}, \mathcal{B} \subset [d_1]$ respectively denote the clusters that i and j belong to, based on $\tilde{\mathbf{M}}_1$. The left-hand side of (1) implies

$$\Theta_{i, i_2, \dots, i_K} = \Theta_{j, i_2, \dots, i_K}, \quad \text{for all } (i_2, \dots, i_K) \in [d_2] \times \cdots \times [d_K]. \quad (2)$$

On the other hand, (1) implies

$$\Theta_{i, i_2, \dots, i_K} = \Theta_{k, i_2, \dots, i_K}, \quad \text{for all } k \in \mathcal{A} \text{ and all } (i_2, \dots, i_K) \in [d_2] \times \cdots \times [d_K], \quad (3)$$

and

$$\Theta_{j, i_2, \dots, i_K} = \Theta_{k, i_2, \dots, i_K}, \quad \text{for all } k \in \mathcal{B} \text{ and all } (i_2, \dots, i_K) \in [d_2] \times \cdots \times [d_K]. \quad (4)$$

Combining (2), (3) and (4), we have

$$\Theta_{i, i_2, \dots, i_K} = \Theta_{k, i_2, \dots, i_K}, \quad \text{for all } k \in \mathcal{A} \cup \mathcal{B} \text{ and all } (i_2, \dots, i_K) \in [d_2] \times \cdots \times [d_K]. \quad (5)$$

Equation (5) implies that \mathcal{A} and \mathcal{B} can be merged into one cluster. This contradicts the irreducibility assumption of the core tensor $\tilde{\mathcal{C}}$. Therefore, \mathbf{M}_1 and $\tilde{\mathbf{M}}_1$ induce a same partition of $[d_1]$, and thus they are equal up to permutation of cluster labels. The proof is now complete. □

A.3 Proof of Theorem 1

The following lemma is useful for the proof of Theorem 1.

Lemma 1. Suppose $\mathcal{Y} = \Theta_{\text{true}} + \mathcal{E}$ with $\Theta_{\text{true}} \in \mathcal{P}$. Let $\hat{\Theta} = \arg \min_{\Theta \in \mathcal{P}} \|\hat{\Theta} - \mathcal{Y}\|_F^2$ be the least-square estimator of Θ_{true} . We have

$$\|\hat{\Theta} - \Theta_{\text{true}}\|_F \leq 2 \sup_{\mu \in \frac{\mathcal{P} - \mathcal{P}'}{|\mathcal{P} - \mathcal{P}'|}} \langle \mu, \mathcal{E} \rangle,$$

where $\mathcal{P} - \mathcal{P}' = \{\Theta - \Theta' : \Theta, \Theta' \in \mathcal{P}\}$ and $\mathcal{S}/|\mathcal{S}| = \{s/\|s\|_2 : s \in \mathcal{S}\}$.

Proof. Based on the definition of least-square estimator, we have

$$\|\hat{\Theta} - \mathcal{Y}\|_F^2 \leq \|\Theta_{\text{true}} - \mathcal{Y}\|_F^2. \quad (6)$$

Combining (6) with the fact

$$\begin{aligned} \|\hat{\Theta} - \mathcal{Y}\|_F^2 &= \|\hat{\Theta} - \Theta_{\text{true}} + \Theta_{\text{true}} - \mathcal{Y}\|_F^2 \\ &= \|\hat{\Theta} - \Theta_{\text{true}}\|_F^2 + \|\Theta_{\text{true}} - \mathcal{Y}\|_F^2 + 2\langle \hat{\Theta} - \Theta_{\text{true}}, \Theta_{\text{true}} - \mathcal{Y} \rangle, \end{aligned}$$

yields

$$\|\hat{\Theta} - \Theta_{\text{true}}\|_F^2 \leq 2\langle \hat{\Theta} - \Theta_{\text{true}}, \mathcal{Y} - \Theta_{\text{true}} \rangle = 2\langle \hat{\Theta} - \Theta_{\text{true}}, \mathcal{E} \rangle.$$

Dividing each side by $\|\hat{\Theta} - \Theta_{\text{true}}\|_F$, we have

$$\|\hat{\Theta} - \Theta_{\text{true}}\|_F \leq 2 \left\langle \frac{\hat{\Theta} - \Theta_{\text{true}}}{\|\hat{\Theta} - \Theta_{\text{true}}\|_F}, \mathcal{E} \right\rangle.$$

The desired inequality follows by noting $\frac{\hat{\Theta} - \Theta_{\text{true}}}{\|\hat{\Theta} - \Theta_{\text{true}}\|_F} \in \frac{\mathcal{P} - \mathcal{P}'}{|\mathcal{P} - \mathcal{P}'|}$. \square

Proof of Theorem 1. To study the performance of the least-square estimator $\hat{\Theta}$, we need to introduce some additional notation. We view the membership matrix M_k as an onto function $M_k: [d_k] \mapsto [R_k]$. With a little abuse of notation, we still use M_k to denote the mapping function and write $M_k \in R_k^{d_k}$ by convention. We use $M = \{M_k\}_{k \in [K]}$ to denote the collection of K membership matrices, and write $\mathcal{M} = \{M : M \text{ is the collection of membership matrices } M_k \text{'s}\}$. For any set J , $|J|$ denotes its cardinality. Note that $|\mathcal{M}| \leq \prod_k R_k^{d_k}$, because each M_k can be identified by a partition of $[d_k]$ into R_k disjoint non-empty sets.

For ease of notation, we define $d = \prod_k d_k$ and $R = \prod_k R_k$. We sometimes identify a tensor in $\mathbb{R}^{d_1 \times \dots \times d_K}$ with a vector in \mathbb{R}^d . By the definition of the parameter space \mathcal{P} , the element $\Theta \in \mathcal{P}$ can be equivalently identified by $\Theta = \Theta(M, C)$, where $M \in \mathcal{M}$ is the collection of K membership matrices and $C = \text{vec}(C) \in \mathbb{R}^R$ is the core tensor. Note that, for a fixed clustering structure M , the space consisting of $\Theta = \Theta(M, \cdot)$ is a linear space of dimension R .

Now consider the least-square estimator

$$\hat{\Theta} = \arg \min_{\Theta \in \mathcal{P}} \{-2\langle \mathcal{Y}, \Theta \rangle + \|\Theta\|_F^2\} = \arg \min_{\Theta \in \mathcal{P}} \{\|\mathcal{Y} - \Theta\|_F^2\}.$$

Based on the Lemma 1,

$$\begin{aligned} \|\hat{\Theta} - \Theta_{\text{true}}\|_F &\leq 2 \sup_{\Theta \in \mathcal{P}} \sup_{\Theta' \in \mathcal{P}} \left\langle \frac{\Theta - \Theta'}{\|\Theta - \Theta'\|_F}, \mathcal{E} \right\rangle \\ &\leq 2 \sup_{M, M' \in \mathcal{M}} \sup_{C, C' \in \mathbb{R}^R} \left\langle \frac{\Theta(M, C) - \Theta(M', C')}{\|\Theta(M, C) - \Theta(M', C')\|_F}, \mathcal{E} \right\rangle. \end{aligned}$$

By union bound, we have, for any $t > 0$,

$$\begin{aligned}
\mathbb{P} \left(\|\hat{\Theta} - \Theta_{\text{true}}\|_F > t \right) &\leq \mathbb{P} \left(\sup_{\mathbf{M}, \mathbf{M}' \in \mathcal{M}} \sup_{\mathbf{C}, \mathbf{C}' \in \mathbb{R}^R} \left| \left\langle \frac{\Theta(\mathbf{M}, \mathbf{C}) - \Theta'(\mathbf{M}', \mathbf{C}')}{\|\Theta(\mathbf{M}, \mathbf{C}) - \Theta'(\mathbf{M}', \mathbf{C}')\|_F}, \mathcal{E} \right\rangle \right| > \frac{t}{2} \right) \\
&\leq \sum_{\mathbf{M}, \mathbf{M}' \in \mathcal{M}} \mathbb{P} \left(\sup_{\mathbf{C}' \in \mathbb{R}^R} \sup_{\mathbf{C} \in \mathbb{R}^R} \left| \left\langle \frac{\Theta(\mathbf{M}, \cdot) - \Theta'(\mathbf{M}', \cdot)}{\|\Theta(\mathbf{M}, \cdot) - \Theta'(\mathbf{M}', \cdot)\|_F}, \mathcal{E} \right\rangle \right| \geq \frac{t}{2} \right) \\
&\leq |\mathcal{M}|^2 C_1^R \exp \left(-\frac{C_2 t^2}{32 \sigma^2} \right) \\
&= \exp \left(2 \sum_k d_k \log R_k + C_1 \prod_k R_k - \frac{C_2 t^2}{32 \sigma^2} \right),
\end{aligned}$$

for two universal constants $C_1, C_2 > 0$. Here the third line follows from [1] (Theorem 1.19) and the fact that $\Theta = \Theta(\mathbf{M}, \cdot)$ lies in a linear space of dimension R . The last line uses $|\mathcal{M}| \leq \prod_k R_k^{d_k}$ and $R = \prod_k R_k$. Choosing $t = C\sigma\sqrt{\prod_k R_k + \sum_k d_k \log R_k}$ yields the desired bound. \square

A.4 Proof of Theorem 2

To prove Theorem 2, first we introduce some notations.

A.4.1 Notations

$\mathbf{c}^{(k)} \in \mathbb{R}^{d_k}$: unknown mode- k cluster membership vector with element $c_{i_k}^{(k)}$ refers to the true label of i_k th fiber in mode k , $\forall k \in [K]$, $i_k \in [d_k]$;

$\hat{\mathbf{c}}^{(k)} \in \mathbb{R}^{d_k}$: mode- k cluster assignment vector with element $\hat{c}_{i_k}^{(k)}$ refers to the assigned label of i_k th fiber in mode k , $\forall k \in [K]$, $i_k \in [d_k]$;

$\mathbf{p}^{(k)} \in \mathbb{R}^{R_k}$: mode- k cluster proportion vector with element $p_{r_k}^{(k)} = \frac{\sum_{i_k=1}^{d_k} \mathbb{I}\{c_{i_k}^{(k)}=r_k\}}{d_k}$, $\forall k \in [K]$, $r_k \in [R_k]$;

$\hat{\mathbf{p}}^{(k)} \in \mathbb{R}^{R_k}$: mode- k label proportion vector with element $\hat{p}_{r_k}^{(k)} = \frac{\sum_{i_k=1}^{d_k} \mathbb{I}\{\hat{c}_{i_k}^{(k)}=r_k\}}{d_k}$, can be seen as a function of $\hat{\mathbf{c}}^{(k)}$, $\forall k \in [K]$, $r_k \in [R_k]$;

$\mathbf{D}^{(k)} = [D_{a_k r_k}^{(k)}] \in \mathbb{R}^{R_k \times R_k}$: mode- k confusion matrix with element $D_{r_k, r'_k}^{(k)} = \frac{1}{d_k} \sum_{i_k=1}^{d_k} \mathbb{I}\{c_{i_k}^{(k)} = r_k, \hat{c}_{i_k}^{(k)} = r'_k\}$, can be seen as a function of $(\hat{\mathbf{c}}^{(1)}, \dots, \hat{\mathbf{c}}^{(K)})$, $\forall k \in [K]$, $r_k \in [R_k]$;

$\mathcal{J}_\tau = \{(\hat{\mathbf{c}}^{(1)}, \dots, \hat{\mathbf{c}}^{(K)}) : \hat{p}_{r_1}^{(1)}(\hat{\mathbf{c}}^{(1)}) > \tau, \dots, \hat{p}_{r_K}^{(K)}(\hat{\mathbf{c}}^{(K)}) > \tau, r_k \in [R_k], k \in [K]\}$;

$\mathcal{I}_d \subset 2^{[d_1]} \times \dots \times 2^{[d_K]}$: is the set of all the blocks that satisfy that $p_{i_k}^{(k)} > \tau$, $\forall i_k \in [d_k]$, $\forall k \in [K]$;

$L_d = \inf\{|I| : I \in \mathcal{I}_d\}$;

$\|\mathbf{A}\|_\infty = \max_{r_1, \dots, r_K} |\mathbf{A}_{r_1, \dots, r_K}|$ for any tensor $\mathbf{A} \in \mathbb{R}^{R_1 \times \dots \times R_K}$.

Remark 1. 1. $\mathbf{D}^{(k)} \mathbf{1} = \mathbf{p}^{(k)}$, $\mathbf{D}^{(k)T} \mathbf{1} = \hat{\mathbf{p}}^{(k)}$. If $\mathbf{D}^{(k)}$ is diagonal, then the assigned labels match the true cluster in mode k , $\forall k \in [K]$.

2. Because our model satisfies the irreducible core assumption, there is always exists a τ such that our estimator $(\hat{\mathbf{c}}^{(1)}, \dots, \hat{\mathbf{c}}^{(K)}) \in \mathcal{J}_\tau$. We denote it as marginal assumption in this proof.

A.4.2 Introduction

To prove the theorem, considering our least-square estimator

$$\begin{aligned}
\hat{\Theta} &= \arg \min_{\Theta \in \mathcal{P}} \{-2 \langle \mathcal{Y}, \Theta \rangle + \|\Theta\|_F^2\} \\
&= \arg \max_{\Theta \in \mathcal{P}} \{\langle \mathcal{Y}, \Theta \rangle - \frac{\|\Theta\|_F^2}{2}\}
\end{aligned}$$

the $\langle \mathcal{Y}, \Theta \rangle = -\frac{\|\Theta\|_F^2}{2}$ is the log-likelihood of the data tensor when our model is a Gaussian tensor block model.

Then the profile log-likelihood $F(\hat{c}^{(1)}, \dots, \hat{c}^{(K)})$ satisfies

$$\begin{aligned} F(\hat{c}^{(1)}, \dots, \hat{c}^{(K)}) &= \sup_{\Theta \in \mathcal{P}} \left\{ \langle \mathcal{Y}, \Theta \rangle - \frac{\|\Theta\|_F^2}{2} \right\} \\ &= \sup_{\Theta \in \mathcal{P}} \left\{ \sum_{i_1, \dots, i_K} y_{i_1, \dots, i_K} c_{r_1(i_1), \dots, r_K(i_K)} - \frac{1}{2} \sum_{i_1, \dots, i_K} c_{r_1(i_1), \dots, r_K(i_K)}^2 \right\} \\ &= \frac{1}{2} \sum_{i_1, \dots, i_K} \overline{y_{r_1(i_1), \dots, r_K(i_K)}}^2 \\ &= \sum_{r_1, \dots, r_K} \prod_{k=1}^K \hat{p}_{r_k}^{(k)} f(\overline{y_{r_1(i_1), \dots, r_K(i_K)}}) \end{aligned}$$

where $f(x) = \frac{x^2}{2}$. Thus our clustering estimator can be represented as

$$(\widehat{\hat{c}}^{(1)}, \dots, \widehat{\hat{c}}^{(K)}) = \arg \max_{(\hat{c}^{(1)}, \dots, \hat{c}^{(K)}) \in \mathcal{J}_\tau} F(\hat{c}^{(1)}, \dots, \hat{c}^{(K)})$$

The error $\|\hat{\Theta} - \Theta\|_F^2$ comes from two aspects: noise and clustering. To measure the error which is from noise, we define a new function $G(\hat{c}^{(1)}, \dots, \hat{c}^{(K)})$:

$$G(\hat{c}^{(1)}, \dots, \hat{c}^{(K)}) = \sum_{r_1, \dots, r_K} [\mathbf{D}^{(1)T} \mathbf{1}]_{r_1} \cdots [\mathbf{D}^{(K)T} \mathbf{1}]_{r_K} f(E_{r_1, \dots, r_K})$$

where $\mathbf{E}(\hat{c}^{(1)}, \dots, \hat{c}^{(K)}) = [E(\hat{c}^{(1)}, \dots, \hat{c}^{(K)})_{r_1, \dots, r_K}] \in \mathbb{R}^{R_1 \times R_2 \times \cdots \times R_K}$,

$$E(\hat{c}^{(1)}, \dots, \hat{c}^{(K)})_{r_1, \dots, r_K} = \frac{\sum_{i_1, \dots, i_K} \sum_{j_1, \dots, j_K} c_{j_1, \dots, j_K} \mathbb{I}\{c_{i_1}^{(1)} = j_1, \hat{c}_{i_1}^{(1)} = r_1\} \cdots \mathbb{I}\{c_{i_K}^{(K)} = j_K, \hat{c}_{i_K}^{(K)} = r_K\}}{\sum_{i_1, \dots, i_K} \mathbb{I}\{\hat{c}_{i_1}^{(1)} = r_1, \dots, \hat{c}_{i_K}^{(K)} = r_K\}}$$

is the average value of $E y_{i_1, \dots, i_K}$ over the block defined by labels r_1, \dots, r_K . Additionally, we define normalized residual matrix $\mathbf{R}(\hat{c}^{(1)}, \dots, \hat{c}^{(K)}) = [\mathbf{R}(\hat{c}^{(1)}, \dots, \hat{c}^{(K)})_{r_1, \dots, r_K}] \in \mathbb{R}^{R_1 \times \cdots \times R_K}$:

$$\mathbf{R}(\hat{c}^{(1)}, \dots, \hat{c}^{(K)})_{r_1, \dots, r_K} = \overline{Y_{r_1, \dots, r_K}} - E(\hat{c}^{(1)}, \dots, \hat{c}^{(K)})_{r_1, \dots, r_K}$$

We use $G(\mathbf{D}^{(1)}, \dots, \mathbf{D}^{(K)}) = \sum_{r_1, \dots, r_K} p_{r_1}^{(1)} \cdots p_{r_K}^{(K)} f(c_{r_1, \dots, r_K})$ to measure the loss. Under the

condition of $\text{MCR}(\hat{\mathbf{M}}_k, \mathbf{P}_k \mathbf{M}_{k, \text{true}}) \geq \varepsilon$ for all $k \in [K]$, we can turn our goal into find the upper bound for the total loss. The following lemma gives the rigorous proof.

Lemma 2. For all $\tau > 0$, for $(\hat{c}^{(1)}, \dots, \hat{c}^{(K)}) \in \mathcal{J}_\tau$ and $\text{MCR}(\hat{\mathbf{M}}_k, \mathbf{P}_k \mathbf{M}_{k, \text{true}}) \geq \varepsilon$, $\exists k \in [K]$,

$$G(\mathbf{D}^{(1)}, \dots, \mathbf{D}^{(K)}) - \sum_{r_1, \dots, r_K} p_{r_1}^{(1)} \cdots p_{r_K}^{(K)} f(c_{r_1, \dots, r_K}) \leq -\frac{\varepsilon \tau^{K-1} \delta_{\min}}{4}$$

Proof of Lemma 2. If $\text{MCR}(\hat{\mathbf{M}}_1, \mathbf{P}_1 \mathbf{M}_{1, \text{true}}) \geq \varepsilon$, then for some r_1 and some $a_1 \neq a'_1$, $\min\{D_{a_1 r_1}^{(1)}, D_{a'_1 r_1}^{(1)}\} \geq \varepsilon$. Since the core tensor is irreducible according to our basic assumption in paper, there exist a_2, \dots, a_K such that $c_{a_1, \dots, a_K} \neq c_{a'_1, \dots, a_K}$. Select the a_2, \dots, a_K such that $(c_{a_1, \dots, a_K} - c_{a'_1, \dots, a_K})^2 = \min_{a_1 \neq a'_1} \max_{a_2, \dots, a_K} (c_{a_1, \dots, a_K} - c_{a'_1, \dots, a_K})^2$. Let $W = [\mathbf{D}^{(1)T} \mathbf{1}]_{r_1} \cdots [\mathbf{D}^{(K)T} \mathbf{1}]_{r_K}$, this is nonzero according to the selection of r_1, \dots, r_K . Now, there exists $c_* \in \mathbb{R}$ such that

$$\begin{aligned}
[\mathcal{N} \times_1 \mathbf{D}^{(1)^T} \times_2 \cdots \times_K \mathbf{D}^{(K)^T}]_{r_1, \dots, r_K} &= D_{a_1 r_1}^{(1)} \cdots D_{a_K r_K}^{(K)} f(c_{a_1, \dots, a_K}) \\
&\quad + D_{a'_1 r_1}^{(1)} \cdots D_{a_K r_K}^{(K)} f(c_{a'_1, \dots, a_K}) \\
&\quad + (W - D_{a_1 r_1}^{(1)} \cdots D_{a_K r_K}^{(K)} - D_{a'_1 r_1}^{(1)} \cdots D_{a_K r_K}^{(K)}) f(c_*)
\end{aligned}$$

Here $\mathcal{N} = [f(c_{a_1, \dots, a_K})] \in \mathbb{R}^{R_1 \times \cdots \times R_K}$ is the loss function evaluated at each block where $[\mathcal{N} \times_1 \mathbf{D}^{(1)^T} \times_2 \cdots \times_K \mathbf{D}^{(K)^T}]_{r_1, \dots, r_K}$ is the weighted value of the loss function. Let $z = \frac{[\mathcal{C} \times_1 \mathbf{D}^{(1)^T} \times_2 \cdots \times_K \mathbf{D}^{(K)^T}]_{r_1, \dots, r_K}}{W}$ where z_{r_1, \dots, r_K} is the (r_1, \dots, r_K) -th weighted entry of the block means. By Taylor expansion and basic inequality $\frac{a+b}{2} \leq \sqrt{\frac{a^2+b^2}{2}}$,

$$\begin{aligned}
&\frac{[\mathcal{N} \times_1 \mathbf{D}^{(1)^T} \times_2 \cdots \times_K \mathbf{D}^{(K)^T}]_{r_1, \dots, r_K}}{W} - f(z) \\
&\geq \frac{\min\{D_{a_1 r_1}^{(1)}, D_{a'_1 r_1}^{(1)}\} D_{a_2 r_2}^{(2)} \cdots D_{a_K r_K}^{(K)}}{4W} (c_{a_1, \dots, a_K} - c_{a'_1, \dots, a_K})^2 \\
&\geq \frac{\varepsilon D_{a_2 r_2}^{(2)} \cdots D_{a_K r_K}^{(K)}}{4W} (c_{a_1, \dots, a_K} - c_{a'_1, \dots, a_K})^2
\end{aligned} \tag{7}$$

Note the inequality ((7)) only holds for a certain $r_1 \in [R_1]$, for any other $r'_1 \in [R_1] \in [R_1]/r_1$, by Jensen's inequality we have

$$\frac{[\mathcal{N} \times_1 \mathbf{D}^{(1)^T} \times_2 \cdots \times_K \mathbf{D}^{(K)^T}]_{r_1, \dots, r_K}}{W} - f(z) \geq 0 \tag{8}$$

With $\sum_{r_k=1}^{R_k} D_{a_k r_k}^{(k)} = \hat{p}_{a_k}^{(k)} \geq \tau$, combining the sum of ((7)) over (r_2, \dots, r_K) and ((8)) gives

$$\begin{aligned}
&G(\mathbf{D}^{(1)}(\hat{\mathbf{c}}^{(1)}), \dots, \mathbf{D}^{(K)}(\hat{\mathbf{c}}^{(K)})) - \sum_{r_1, \dots, r_K} \prod_{k=1}^K p_{r_k}^{(k)} f(c_{r_1, \dots, r_K}) \\
&= \sum_{r_1, \dots, r_K} [\mathbf{D}^{(1)^T} \mathbf{1}]_{r_1} \cdots [\mathbf{D}^{(K)^T} \mathbf{1}]_{r_K} f\left(\frac{[\mathcal{C} \times_1 \mathbf{D}^{(1)^T} \times_2 \cdots \times_K \mathbf{D}^{(K)^T}]_{r_1, \dots, r_K}}{[\mathbf{D}^{(1)^T} \mathbf{1}]_{r_1} \cdots [\mathbf{D}^{(K)^T} \mathbf{1}]_{r_K}}\right) \\
&\leq -\varepsilon \sum_{r_2, \dots, r_K} \frac{D_{a_2 r_2}^{(2)} \cdots D_{a_K r_K}^{(K)}}{4} (c_{a_1, \dots, a_K} - c_{a'_1, \dots, a_K})^2 \\
&\leq -\frac{\varepsilon \tau^{K-1} \delta_{min}}{4}
\end{aligned}$$

Similarly, the proof also goes through if $\text{MCR}(\hat{\mathbf{M}}_k, \mathbf{P}_k \mathbf{M}_{k, true}) \geq \varepsilon$, $k \in [K]$. \square

A.4.3 Proof

Proof of Theorem 2. By lemma 2, we obtained

$$\begin{aligned}
&\mathbb{P}\left(\text{MCR}(\hat{\mathbf{M}}_k, \mathbf{P}_k \mathbf{M}_{k, true}) \geq \varepsilon\right) \\
&\leq \mathbb{P}\left(G(\mathbf{D}^{(1)}, \dots, \mathbf{D}^{(K)}) - \sum_{r_1, \dots, r_K} p_{r_1}^{(1)} \cdots p_{r_K}^{(K)} f(c_{r_1, \dots, r_K}) \leq -\frac{\varepsilon \tau^{K-1} \delta_{min}}{4}\right) \\
&= \mathbb{P}\left(G(\mathbf{D}^{(1)}(\widehat{\mathbf{c}}^{(1)}), \dots, \mathbf{D}^{(K)}(\widehat{\mathbf{c}}^{(K)})) - F(\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(K)}) \leq -\frac{\varepsilon \tau^{K-1} \delta_{min}}{4}\right)
\end{aligned} \tag{9}$$

Additionally, letting $r_d = \sup_{\mathcal{J}_\tau} |F(\hat{\mathbf{c}}^{(1)}, \dots, \hat{\mathbf{c}}^{(K)}) - G(\mathbf{D}^{(1)}(\hat{\mathbf{c}}^{(1)}), \dots, \mathbf{D}^{(K)}(\hat{\mathbf{c}}^{(K)}))|$ which refers to the loss caused only by noise, when $G(\mathbf{D}^{(1)}(\widehat{\mathbf{c}}^{(1)}), \dots, \mathbf{D}^{(K)}(\widehat{\mathbf{c}}^{(K)})) - F(\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(K)}) \leq -\frac{\varepsilon \tau^{K-1} \delta_{min}}{4}$, we have

$$F(\hat{\mathbf{c}}^{(1)}, \dots, \hat{\mathbf{c}}^{(K)}) - F(\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(K)}) \leq 2r_d - \frac{\varepsilon \tau^{K-1} \delta_{min}}{4} \quad (10)$$

Plug the inequality ((10)) back into inequality ((9)), we obtain

$$\begin{aligned} & \mathbb{P} \left(\text{MCR}(\hat{\mathbf{M}}_k, \mathbf{P}_k \mathbf{M}_{k,true}) \geq \varepsilon \right) \\ & \leq \mathbb{P} \left(F(\widehat{\mathbf{c}}^{(1)}, \dots, \widehat{\mathbf{c}}^{(K)}) - F(\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(K)}) \leq 2r_d - \frac{\varepsilon \tau^{K-1} \delta_{min}}{4} \right) \\ & \leq \mathbb{P} \left(r_d \geq \frac{\varepsilon \tau^{K-1} \delta_{min}}{8} \right) \end{aligned} \quad (11)$$

Now we convert our problem into find the upper bound of $\mathbb{P} \left(r_d \geq \frac{\varepsilon \tau^{K-1} \delta_{min}}{8} \right)$. Consider $\mathbb{P}(r_d \leq t)$, because f is locally lipschitz continuous with lipschitz constant $c = \sup |f'(\mu)|$ for μ in a neighborhood of the convex hull of the entries of \mathcal{C} ,

$$\begin{aligned} & |F(\hat{\mathbf{c}}^{(1)}, \dots, \hat{\mathbf{c}}^{(K)}) - G(\mathbf{D}^{(1)}(\hat{\mathbf{c}}^{(1)}), \dots, \mathbf{D}^{(K)}(\hat{\mathbf{c}}^{(K)}))| \\ & \leq \sum_{r_1, \dots, r_K} \hat{p}_{r_1}^{(1)} \hat{p}_{r_2}^{(2)} \dots \hat{p}_{r_K}^{(K)} |f(\overline{Y_{r_1, \dots, r_K}}) - f(E_{r_1, \dots, r_K})| \\ & \leq c \|\mathbf{R}(\hat{\mathbf{c}}^{(1)}, \dots, \hat{\mathbf{c}}^{(K)})\|_\infty \end{aligned} \quad (12)$$

Combining ((11)), ((12)), Hoeffding's inequality, $L_d \geq \tau^K \prod_{k=1}^K d_k$ and $C_2 = \frac{\tau^{3K-2}}{128c^2}$ yields the desired conclusion. \square

A.5 Sparse estimator

Lemma 3. Consider the regularized least-square estimation,

$$\hat{\Theta}^{sparse} = \arg \min_{\Theta \in \mathcal{P}} \{ \|\mathcal{Y} - \Theta\|_F^2 + \lambda \|\mathcal{C}\|_\rho \}, \quad (13)$$

where $\mathcal{C} = \llbracket c_{r_1, \dots, r_K} \rrbracket \in \mathbb{R}^{R_1 \times \dots \times R_K}$ is the block-mean tensor, $\|\mathcal{C}\|_\rho$ is the penalty function with ρ being an index for the tensor norm, and λ is the penalty tuning parameter. We have

$$\hat{\mathcal{C}}_{r_1, \dots, r_K}^{sparse} = \begin{cases} \hat{c}_{r_1, \dots, r_K}^{ols} \mathbf{1} \left\{ |\hat{c}_{r_1, \dots, r_K}^{ols}| \geq \sqrt{\frac{\lambda}{n_{r_1, \dots, r_K}}} \right\} & \text{if } \rho = 0, \\ \text{sign}(\hat{c}_{r_1, \dots, r_K}^{ols}) \left(|\hat{c}_{r_1, \dots, r_K}^{ols}| - \frac{\lambda}{2n_{r_1, \dots, r_K}} \right)_+ & \text{if } \rho = 1, \end{cases} \quad (14)$$

where $a_+ = \max(a, 0)$ and $\hat{c}_{r_1, \dots, r_K}^{ols}$ denotes the ordinary least-square estimate as in Algorithm 1.

Proof. We formulate the estimation of \mathcal{C} as a regularized least-square regression. Note that $\Theta \in \mathcal{P}$ implies that

$$\Theta = \mathcal{C} \times_1 \mathbf{M}_1 \times \dots \times_K \mathbf{M}_K.$$

Define $\mathbf{X} = \mathbf{M}_1 \otimes \dots \otimes \mathbf{M}_K \in \mathbb{R}^{d \times R}$, where $d = \prod_k d_k$ and $R = \prod_k R_k$, and $\beta = \text{vec}(\mathcal{C}) \in \mathbb{R}^R$. Here \mathbf{X} is a membership matrix that indicates the block allocation among tensor entries. Specifically, \mathbf{X} consists of orthogonal columns with $\mathbf{X}^T \mathbf{X} = \text{diag}(n_1, \dots, n_R)$, where n_r is the number of entries in the tensor block that corresponds to the r -th column of \mathbf{X} .

For a given set of M'_k s, the optimization (15) with respect to \mathcal{C} is equivalent to a regularized linear regression with $\mathbf{Y} = \text{vec}(\mathcal{Y})$ as the response and \mathbf{X} as the design matrix:

$$L(\beta) = \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_\rho. \quad (15)$$

When $\lambda = 0$ (no penalty), the minimizer is $\hat{\beta}^{\text{ols}} = (\hat{\beta}_1^{\text{ols}}, \dots, \hat{\beta}_R^{\text{ols}}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$, where $\hat{\beta}_r^{\text{ols}} = \frac{1}{n_r} \mathbf{y}_r \mathbf{1}_{n_r}^T$ for all $r \in [R]$.

Case 1: $\rho = 0$.

Note that \mathbf{X} induces a partition of indices $[d]$ into R blocks. With a little abuse of notation, we use $\mathbf{r} = \{i \in [d] : \mathbf{X}(i) = r\}$ to denote the tensor indices that belong to the r th block, and use $\mathbf{y}_r \in \mathbb{R}^{n_r}$ to denote the corresponding tensor entries. By the orthogonality of \mathbf{X} , we have

$$\begin{aligned} L(\beta) &= \sum_{r=1}^R \|\mathbf{y}_r - \beta_r \mathbf{1}_{n_r}\|_2^2 + \lambda \sum_{r=1}^R \mathbb{1}\{\beta_r \neq 0\} \\ &= \sum_{r=1}^R \underbrace{(\|\mathbf{y}_r - \beta_r \mathbf{1}_{n_r}\|_2^2 + \lambda \mathbb{1}\{\beta_r \neq 0\})}_{:=L_r(\beta_r)} \end{aligned}$$

The optimization can be separated into each of β_r 's. For any $r \in [R]$, the sub-optimization $\min_{\beta_r} L_r(\beta_r)$ has a closed-form solution

$$\min_{\beta_r} L_r(\beta_r) = \begin{cases} \mathbf{y}_r^T \mathbf{y}_r - n_r \left(\hat{\beta}_r^{\text{ols}} \right)^2 + \lambda & \text{if } \hat{\beta}_r^{\text{ols}} \neq 0, \\ \mathbf{y}_r^T \mathbf{y}_r & \text{if } \hat{\beta}_r^{\text{ols}} = 0, \end{cases}$$

with

$$\arg \min_{\beta_r} L_r(\beta_r) = \begin{cases} 0 & \text{if } n_r \left(\hat{\beta}_r^{\text{ols}} \right)^2 \leq \lambda, \\ \hat{\beta}_r^{\text{ols}} & \text{otherwise.} \end{cases} \quad (16)$$

Solution (16) can be simplified as $\hat{\beta}_r^{\text{sparse}} = \hat{\beta}_r^{\text{ols}} \mathbb{1}\{|\hat{\beta}_r^{\text{ols}}| \leq \sqrt{\frac{\lambda}{n_r}}\}$. The proof is complete by noting that $\hat{c}_{r_1, \dots, r_R}^{\text{sparse}} = \hat{\beta}_r^{\text{sparse}}$ and $n_{r_1, \dots, r_R} = n_r$ for all $(r_1, \dots, r_R) \in [R_1] \times \dots \times [R_K]$.

Case 2: $\rho = 1$.

Similar as in Case 1, we write the optimization (15) as

$$L(\beta) = \sum_{r=1}^R \underbrace{(\|\mathbf{y}_r - \beta_r \mathbf{1}_{n_r}\|_2^2 + \lambda |\beta_r|)}_{:=L_r(\beta_r)},$$

where, with a little abuse of notation, we still use $L_r(\beta_r)$ to denote the sub-optimization. To solve $\arg \min_{\beta_r} L_r(\beta_r)$, we use the properties of subderivative. Taking the subderivative with respect to β_r , we obtain

$$\frac{\partial L_r(\beta_r)}{\partial \beta_r} = \begin{cases} 2n_r \beta_r - 2n_r \hat{\beta}_r^{\text{ols}} + \lambda & \text{if } \beta_r > 0, \\ [2n_r \beta_r - 2\hat{\beta}_r^{\text{ols}} - \lambda, 2n_r \beta_r - \hat{\beta}_r^{\text{ols}} + \lambda] & \text{if } \beta_r = 0, \\ 2n_r \beta_r - 2n_r \hat{\beta}_r^{\text{ols}} + \lambda & \text{if } \beta_r < 0. \end{cases}$$

Because $\hat{\beta}_r^{\text{sparse}}$ minimizes $L_r(\beta_r)$ if and only if $0 \in \frac{\partial L_r(\beta_r)}{\partial \beta_r}$, we have:

$$\hat{\beta}_r^{\text{sparse}} = \begin{cases} \hat{\beta}_r^{\text{ols}} + \frac{\lambda}{2n_r} & \text{if } \hat{\beta}_r^{\text{ols}} < -\frac{\lambda}{2n_r}, \\ 0 & \text{if } \hat{\beta}_r^{\text{ols}} \in [-\frac{\lambda}{2n_r}, \frac{\lambda}{2n_r}], \\ \hat{\beta}_r^{\text{ols}} - \frac{\lambda}{2n_r} & \text{if } \hat{\beta}_r^{\text{ols}} > \frac{\lambda}{2n_r}. \end{cases} \quad (17)$$

The solution (17) can be simplified as

$$\hat{\beta}_r^{\text{sparse}} = \text{sign}(\hat{\beta}_r^{\text{ols}}) \left(|\hat{\beta}_r^{\text{ols}}| - \frac{\lambda}{2n_r} \right)_+, \quad \text{for all } r \in [R].$$

□

B Supplementary Figures and Tables

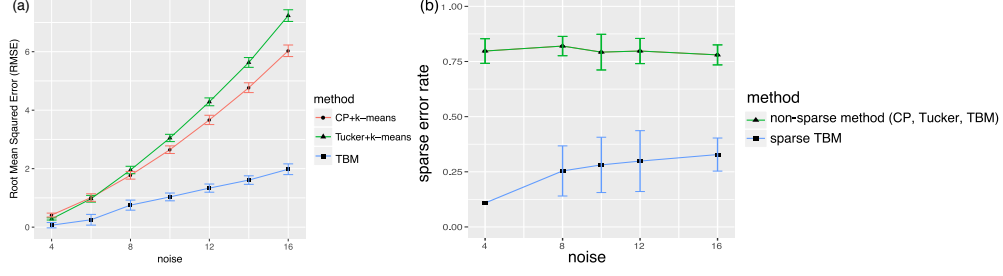


Figure S1: (a) estimation error and (b) sparse error rate against noise for sparse tensors of dimension $(40, 40, 40)$ when $p = 0.8$.

Dimensions (d_1, d_2, d_3)	True clustering sizes (R_1, R_2, R_3)	Noise (σ)	Estimated clustering sizes ($\hat{R}_1, \hat{R}_2, \hat{R}_3$)
(40, 40, 40)	(4, 4, 4)	4	(4, 4, 4) \pm (0, 0, 0)
(40, 40, 40)	(4, 4, 4)	8	(3.94, 3.96, 3.96) \pm (0.03, 0.03, 0.03)
(40, 40, 40)	(4, 4, 4)	12	(3.08, 3.12, 3.12) \pm (0.10, 0.10, 0.10)
(40, 40, 80)	(4, 4, 4)	4	(4, 4, 4) \pm (0, 0, 0)
(40, 40, 80)	(4, 4, 4)	8	(4, 4, 4) \pm (0, 0, 0)
(40, 40, 80)	(4, 4, 4)	12	(3.96, 3.96, 3.92) \pm (0.04, 0.04, 0.04)
(40, 40, 40)	(2, 3, 4)	4	(2, 3, 4) \pm (0, 0, 0)
(40, 40, 40)	(2, 3, 4)	8	(2, 3, 3.96) \pm (0, 0, 0.03)
(40, 40, 40)	(2, 3, 4)	12	(2, 2.96, 3.60) \pm (0, 0.05, 0.09)

Table S1: The simulation results for estimating $\mathbf{R} = (R_1, R_2, R_3)$. Bold number indicates no significant difference between the estimate and the ground truth, based on a z -test with a level 0.05.

Tissues	Over-expressed genes	Block-means	Under-expressed genes	Block-means
Cluster 1	GFAP, MBP	10.88	GPR6, DLX5, DLX6, NKX2-1	-8.40
Cluster 2	GFAP, MBP	5.98	CDH9, RXFP1, CRH, ARX, CARTPT, DLX1, FEZF2	-9.49
Cluster 3	GFAP, MBP	8.34	AVPR1A, CCKAR, CHRN4, CYP19A1, HOXA4, LBX1, SLC6A3	-8.45
			TBR1, SLC17A6, SLC30A3	-8.17
Cluster 4	GFAP, MBP	8.83	AVPR1A, CCKAR, CHRN4, CYP19A1, HOXA4, LBX1, SLC6A3	-8.40
			DAO, EN2, EOMES	-6.57

Table S2: Top expression blocks from the multi-tissue gene expression analysis. The tissue clusters are described in Supplementary Section D.

Countries	Countries	Relation types
Cluster 1	Clusters 4 and 5	reltreaties, booktranslations, relbooktranslations, relexports, exports3
Clusters 1 and 4	Cluster 5	relintergovorgs, relngo, intergovorgs3, ngoorgs3
Cluster 3	Clusters 1, 4, and 5	commonbloc0, blockpositionindex
Clusters 1 and 3	Clusters 4 and 5	
Cluster 1	Cluster 3	timesinceally, independence
Cluster 4	Cluster 5	
Cluster 4	Cluster 5	treaties, conferences, weightedunvote, unweightedunvote, intergovorgs, ngo, officialvisits, exportbooks, relexportbooks, tourism, reltourism, tourism3, exports, militaryalliance, commonbloc2

Table S3: Top blocks from the *Nations* data analysis. The countries clusters are described in Supplementary Section D.

C Time complexity

The total cost of our Algorithm 1 is $\mathcal{O}(d)$ per iteration, where $d = \prod_k d_k$ denotes the total number of tensor entries. The per-iteration computational cost scales linearly with the sample size, and this

complexity is comparable to the classical tensor methods such as CP and Tucker decomposition. More specifically, each iteration of Algorithm 1 consists of updating the core tensor \mathcal{C} and K membership matrices M_k 's. The update of \mathcal{C} requires $\mathcal{O}(d)$ operations and the update of M_k requires $\mathcal{O}(R_k \frac{d}{d_k})$ operations. Therefore the total cost is $\mathcal{O}(d + d \sum_k \frac{R_k}{d_k})$.

D Additional information for real data analysis

Multi-tissue gene expression. The gene expression data we analyzed is part of the GTEx v6 datasets (<https://www.gtexportal.org/home/datasets>). We cleaned and preprocessed the data following the steps in [2]. We focused on the 13 brain tissues, 193 individuals, and 362 annotated genes provided by Atlax of the Developing Human Brain (<http://www.brainspan.org/ish>). After applying the ℓ_0 penalized TBM to the mean-centered data tensor, we identified the following four clusters of tissues:

- Cluster 1: Substantia nigra, Spinal cord (cervical c-1)
- Cluster 2: Cerebellum, Cerebellar Hemisphere
- Cluster 3: Caudate (basal ganglia), Nucleus accumbens (basal ganglia), Putamen (basal ganglia)
- Cluster 4: Cortex, Hippocampus, Anterior cingulate cortex (BA24), Frontal Cortex (BA9), Hypothalamus, Amygdala

We found that most tissue clusters are spatially restricted to specific brain regions, such as the two cerebellum tissues (cluster 2), three basal ganglia tissues (cluster 3), and the cortex tissues (cluster 4). Supplementary Table S2 reports the associated gene cluster for each tissue cluster. Because our method attaches importance to blocks by the absolute mean estimates, our method is able to detect both over- and under-expression patterns. Blocks with highly positive means correspond to over-expressed genes, whereas blocks with highly negative means correspond to under-expressed genes.

Nations dataset. This is a $14 \times 14 \times 56$ binary tensor consisting of 56 political relations of 14 countries between 1950 and 1965 [3]. The tensor entry indicates the presence or absence of a political action, such as “treaties”, “sends tourists to”, between the nations. We applied the ℓ_0 penalized TBM to the binary-valued data tensor, and we identified the following five clusters of countries:

- Cluster 1: Brazil, Egypt, India, Israel, Netherlands
- Cluster 2: Burma, Indonesia, Jordan
- Cluster 3: China, Cuba, Poland, USSA
- Cluster 4: USA
- Cluster 5: UK

Supplementary Table S3 reports the cluster constitutions for top blocks. Because the tensor entries take value on either 0 or 1, the top blocks mostly have mean one.

References

- [1] Phillippe Rigollet and Jan-Christian Hütter. High dimensional statistics. *Lecture notes for course 18S997*, 2015.
- [2] Miaoyan Wang, Jonathan Fischer, and Yun S Song. Three-way clustering of multi-tissue multi-individual gene expression data using constrained tensor decomposition. *Annals of Applied Statistics, in press*, 2019.
- [3] Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. A three-way model for collective learning on multi-relational data. In *International Conference on Machine Learning*, volume 11, pages 809–816, 2011.