
Multi-way block localization via tensor-based clustering

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We consider the task of simultaneously clustering each mode of a large noisy tensor.
2 We assume that the tensor elements are distributed with a block-specific mean
3 and propose a least-square estimation for multi-way clustering. An ℓ_1 penalty is
4 applied to the block-means in order to select and identify important blocks. We
5 show that our method is applicable to large tensors with a wide range of multi-way
6 cluster structure, including a single block, multiple blocks, checkerboard clusters,
7 1-way or lower-way blocks. Our proposal amounts to a sparse, multi-way version
8 of k -mean clustering, and a relaxation of our proposal yields the tensor Tucker
9 decomposition. The performance of our proposals are demonstrated in simulations
10 and on...

11 1 Introduction

12 In recent years, much interest has centered around the unsupervised analysis of high-dimensional
13 high-order tensor data.

Here is an example of tensor clustering by using our proposed method.

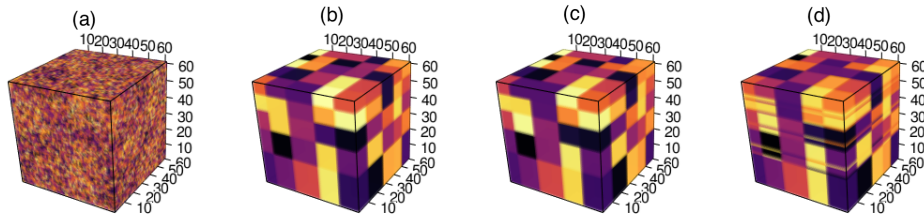


Figure 1: (a): a $60 \times 60 \times 60$ tensor with 5 clusters in each mode; (b): true underlying mean signal within each cluster; (c): mean signal estimated by our proposed approach with true number of clusters $(5, 5, 5)$; (d): mean signal estimated by k -means clustering on each mode with true number of clusters $(5, 5, 5)$.

14

15 2 Preliminaries

16 We say that an event A occurs “with high probability” if $\mathbb{P}(A)$ tends to 1 as the dimension $d_{\min} =$
17 $\min\{d_1, \dots, d_k\}$ tends to infinity. We say that A occurs “with very high probability” if $\mathbb{P}(A)$ tends
18 to 1 faster than any polynomial of d .

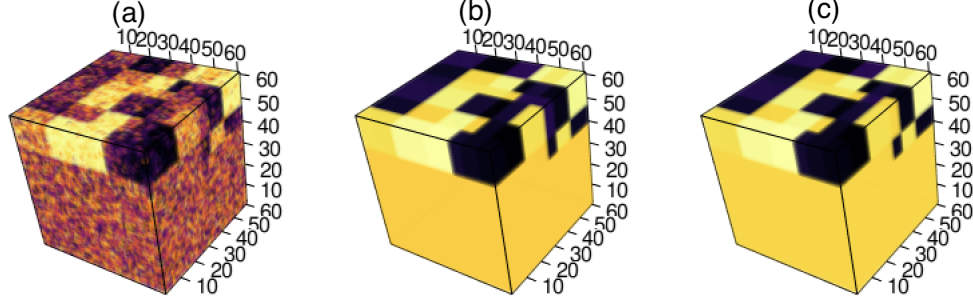


Figure 2: (a): $60 \times 60 \times 60$ sparse tensor; (b) true underlying means; (c) mean signal estimated by our approach with estimated number of clusters and estimated λ .

We use lower-case letters (a, b, u, v, \dots) for scalars and vectors. We use upper-case boldface letters ($\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$) for matrices, and calligraphy letter ($\mathcal{A}, \mathcal{B}, \mathcal{C}, \dots$) for tensors of order $K \geq 3$. $\mathbf{x} \otimes \mathbf{y}$ is the Kronecker product of two vectors. For any set J , $|J|$ denotes its cardinality. $[d]$ represents the set $\{1, 2, \dots, d\}$.

A clustering of d objects can be represented by a partition of the index set $[d] = \{1, \dots, d\}$ into R disjoint non-empty subsets. We refer to R the clustering size. It is often convenient to represent the clustering (or partition) using the “membership matrix”. A membership matrix \mathbf{M} is an d -by- R matrix whose (i, j) -entry is 1 if and only if the element i belongs to the cluster j , and 0 otherwise. Based on the definition, \mathbf{M} is binary matrix with orthogonal columns, and the matrix elements sum up to d . The membership matrix \mathbf{M} can also be viewed as a mapping $\mathbf{M} : [d] \mapsto [R]$. With a little abuse of notation, we still use the \mathbf{M} to denote the mapping, and use $\mathbf{M}(i)$ to denote the cluster label that entry i belongs to.

Throughout the paper, we will use the terms “partition”, “clustering”, and “membership matrix” exchangeably.

For a higher-order tensor, the above concepts can be applied to each of the modes. We use the term “mode- k clustering” to refer to the partition along the k -th mode of the tensor, and reserve “block” to mean a multi-way block in the Cartesian product of mode- k clusters.

3 Tensor block model

Let $\mathcal{Y} = \llbracket y_{i_1, \dots, i_K} \rrbracket \in \mathbb{R}^{d_1 \times \dots \times d_K}$ denote an order- K , (d_1, \dots, d_K) -dimensional data tensor. The main assumption on tensor block model is that the observed data tensor \mathcal{Y} is a noisy realization of an underlying tensor that exhibits a checkbox structure (see Figure 1). Specifically, suppose that there are R_k clusters along the k -th mode of the tensor for $k \in [K]$. If the tensor entry y_{i_1, \dots, i_K} belongs to the block jointly determined by the r_k -th mode- k cluster with $r_k \in [R_k]$, then we assume that

$$y_{i_1, \dots, i_K} = c_{r_1, \dots, r_K} + \varepsilon_{i_1, \dots, i_K}, \quad \text{for } (i_1, \dots, i_K) \in [d_1] \times \dots \times [d_K], \quad (1)$$

where μ_{r_1, \dots, r_K} is the mean of the tensor block indexed by (r_1, \dots, r_K) , and $\varepsilon_{i_1, \dots, i_K}$ ’s are independent, mean-zero noise terms to be specified later. Our goal is to (i) find partitions along each of the modes, and (ii) estimate the block means $\{c_{r_1, \dots, r_K}\}$, such that a corresponding blockwise-constant checkbox structure emerges in the data tensor.

The above tensor block model (1) falls into a larger class of non-overlapping, constant-mean clustering models [1], in that each tensor entry belongs to exactly one block with a common mean. The model (1) can be equivalently expressed as a special tensor Tucker model,

$$\mathcal{Y} = \mathcal{C} \times_1 \mathbf{M}_1 \times_2 \dots \times_K \mathbf{M}_K + \mathcal{E}, \quad (2)$$

where $\mathcal{C} \in \mathbb{R}^{R_1 \times \dots \times R_K}$ is a core tensor consisting of block means, $\mathbf{M}_k \in \{0, 1\}^{d_k \times R_k}$ are membership matrices indicating the block allocations along mode k for $k \in [K]$, and $\mathcal{E} = \llbracket \varepsilon_{i_1, \dots, i_K} \rrbracket$ is the noise tensor. The distinction between our model (2) and a classical Tucker model is that we require the factors \mathbf{M}_k to be membership matrices. Our model (2) can be viewed as a super-sparse Tucker model, in the sense that the each column of \mathbf{M}_k consists of one copy of 1’s and massive 0’s.

We now introduce the assumptions on the noise tensor \mathcal{E} . We assume that $\varepsilon_{i_1, \dots, i_K}$'s are independent, mean-zero, σ -subgaussian noises, where $\sigma > 0$ is the subgaussianity parameter. More precisely,

$$\mathbb{E} e^{\lambda \varepsilon_{i_1, \dots, i_K}} \leq e^{\lambda^2 \sigma^2 / 2}, \quad \text{for all } (i_1, \dots, i_K) \in [d_1] \times \dots \times [d_K] \text{ and } \lambda \in \mathbb{R}. \quad (3)$$

Th assumption (3) is fairly general, which includes many common noises, such as Gaussian errors, Bernoulli errors, bounded errors, or even combinations of them. In particular, we consider two examples of the tensor block model that commonly appear in the literature:

Example 1 (Gaussian Multi-Cluster Model) *Let \mathcal{Y} be a continuous-valued tensor. The Gaussian Multi-cluster model $y_{i_1, \dots, i_K} \sim_{i.i.d.} N(\mu_{r_1, \dots, r_K}, \sigma^2)$ is a special case of model (1) with the subgaussianity parameter σ equal to the error variance.*

Example 2 (Stochastic Block Model) *Let \mathcal{Y} be a binary-valued tensor. The multiway stochastic block model $y_{i_1, \dots, i_K} \sim_{i.i.d.} \text{Bernoulli}(\mu_{r_1, \dots, r_K})$ is a special case of model (1) with the subgaussianity parameter σ equal to $\frac{1}{4}$.*

More generally, our model also applied to hybrid error distributions in which different types of distribution can be allowed for different portions of the data. This scenario may happen, for example, when the data tensor \mathcal{Y} represents concatenated measurements from multiple data sources.

We consider a least-square approach for estimating model (1). Let $\Theta = \mathcal{C} \times_1 \mathbf{M}_1 \times_2 \dots \times_K \mathbf{M}_K$ denote the mean signal tensor with block structure. The mean tensor is assumed to belong to the following parameter space

$$\mathcal{P}_{R_1, \dots, R_K} = \left\{ \Theta \in \mathbb{R}^{d_1 \times \dots \times d_K} : \Theta = \mathcal{C} \times_1 \mathbf{M}_1 \times_2 \dots \times_K \mathbf{M}_K, \text{ with some} \right. \quad (4)$$

$$\text{membership matrices } \mathbf{M}_k \text{'s and a core tensor } \mathcal{C} \in \mathbb{R}^{R_1 \times \dots \times R_K} \left. \right\}. \quad (5)$$

As in most previous work on tensor clustering, we assume that clustering size $\mathbf{R} = (R_1, \dots, R_K)$ are known in our theoretical analysis and simply write \mathcal{P} for short. In practice, \mathbf{R} needs to be determined from data; we address this general case in Section 5.2. The least-square estimator for model (1) is

$$\hat{\Theta} = \arg \min_{\Theta \in \mathcal{P}} \left\{ -2\langle \mathcal{Y}, \Theta \rangle + \|\Theta\|_F^2 \right\}. \quad (6)$$

The objective is equal (ignoring constants) to the sum of squares $\|\mathcal{Y} - \Theta\|_F^2$ and hence the name of our estimator. Estimating Θ consists of finding both the core tensor \mathcal{C} and the membership matrix estimates \mathbf{M}_k 's. Before we discuss the properties of $\hat{\Theta}$, we present the identifiability of \mathbf{M}_k 's and \mathcal{C} from Θ .

The following irreducible assumption is necessary for the tensor block model to be identifiable.

Assumption 1 (Irreducible cores) *The core tensor \mathcal{C} is called irreducible if it cannot be written as a block tensor with the number of mode- k clusters smaller than R_k , for any $k \in [K]$.*

In the matrix case ($K = 2$), the assumption is equivalent to saying that \mathcal{C} has no two identical rows and no two identical columns. In the higher-order case, it requires that none of order- $(K-1)$ fibers of \mathcal{C} are identical. Note that the being irreducible is a weaker assumption than being full rank.

Proposition 1 (Identifiability) *Consider a Gaussian or Bernoulli tensor block model (2). Suppose the core tensor satisfies Assumption 1. Then every factor matrix \mathbf{M}_k is identifiable up to permutations of cluster labels.*

Our identifiability result is stronger than the classical Tucker model. In a classical Tucker model [2, 3] and many other factor analyses [4, 5], the factors are identifiable only up to orthogonal rotations. In those models, the (column) space spanned by \mathbf{M}_k can be recovered, but not the individual factors. In contrast, our model does not suffer from rotational invariance, and as we show in Section 4, every single factor can be consistently estimated in high dimensions. This brings a benefit to the interpretation of tensor factors in the block model.

4 Statistical convergence

In this section, we assess the estimation accuracy of the least-squares estimator (6). While the least squares corresponds to the maximum likelihood estimator (MLE) for Gaussian tensor model, the

96 same assertion does not hold for the other types of distribution such as stochastic tensor block model.
 97 Surprisingly, we will show that, with very high probability, a simple least-square estimator can
 98 achieve a convergence rate that is nearly the same as the optimal one in a general class of block
 99 tensors.

100 For the true signal tensor $\Theta_{\text{true}} \in \mathcal{P}$ and its estimator $\hat{\Theta} \in \mathcal{P}$, define

$$\text{Loss}(\Theta_{\text{true}}, \hat{\Theta}) = \frac{1}{\prod_k d_k} \|\Theta_{\text{true}} - \hat{\Theta}\|_F^2.$$

101 **Theorem 1 (Convergence rate)** *Let $\hat{\Theta}$ be the least-square estimator of Θ_{true} under model (1). There*
 102 *exist two constants $C_1, C_2 > 0$ such that, with very high probability,*

$$\text{Loss}(\Theta_{\text{true}}, \hat{\Theta}) \leq \frac{C_1 \sigma^2}{\prod_k d_k} \left(\prod_k R_k + \sum_k d_k \log R_k \right), \quad (7)$$

103 *holds uniformly over $\Theta_{\text{true}} \in \mathcal{P}_{\mathbf{R}}$ and all error distribution satisfying (3).*

104 The convergence rate in (7) consists of two parts. The first part $\prod_k R_k$ reflects the complexity for
 105 estimating the core tensor \mathcal{C} , while the second part $\sum_k d_k \log R_k$ results from the complexity for
 106 estimating the supports of \mathbf{M}_k 's. It is the price that one has to pay for not knowing the locations of
 107 the blocks.

108 We now compare our bound with existing literature. The classical Tucker tensor decomposition has a
 109 minimax convergence rate $\sum_k d_k R'_k$ [2], where R'_k is the multilinear rank at mode k . In the case of
 110 block model, this yields $\sum_k d_k R_k$, because the mode- k rank is bounded by the number of clusters
 111 in mode- k . Now, as both the dimension $d_{\min} = \min_k d_k$ and clustering size $R_{\min} = \min_k R_k$ tend
 112 to infinity, we have $\prod_k R_k + \sum_k d_k \log R_k \ll \sum_k d_k R_k$. Therefore, by fully exploiting the block
 113 structure, we obtain a better convergence rate than previously possible.

114 Our bound also generalizes the previous results on structured matrix estimation in network analysis [6,
 115 7]. The optimal convergence rate for estimating the (matrix) stochastic block model was $R_1 R_2 +$
 116 $d_1 \log R_1 + d_2 \log R_2$ [6], which fits into our special case when $K = 2$. Earlier work [7] suggests
 117 the following heuristics on the sample complexity for high-dimensional matrix problems:

$$\frac{(\text{number of parameters}) + \log(\text{complexity of models})}{\text{number of samples}}. \quad (8)$$

118 Our result supports this important principle for general $K \geq 2$. Note that, in tensor estimation, the
 119 total number of entries corresponds to the sample size $\prod_k d_k$, the number of parameters is $\prod_k R_k$,
 120 and combinatoric complexity for estimating block models is of order $\prod_k R_k^{d_k}$. The principle (8) thus
 121 provide an intuition for (7).

122 The optimality of the least-square estimator is safeguarded by the following information-theoretical
 123 lower bound.

124 **Theorem 2 (Information-theoretical bound)** *Under the Gaussian or Bernoulli tensor block mod-*
 125 *els, there exist some constants $\alpha_0 > 0, \beta_0 \in (0, 1)$, such that*

$$\inf_{\hat{\Theta}} \sup_{\Theta \in \mathcal{P}} \mathbb{P} \left\{ \|\hat{\Theta} - \Theta_{\text{true}}\|_F^2 > c_0 \sigma^2 \prod_k R_k + \sum_k d_k \log R_k \right\} > \beta_0.$$

126 *(add the proof)*

127 We next study the clustering consistency of our method. Let $\hat{\mathbf{M}}_k$ denote the estimated membership
 128 matrix corresponding to $\hat{\Theta}$, and $\mathbf{M}_{k,\text{true}}$ the true membership matrix. We define the clustering error
 129 rate as $\text{CER}(\hat{\mathbf{M}}_k, \mathbf{M}_{k,\text{true}}) = d_k^{-1} \sum_{i \in [d_k]} \mathbb{1}_{\{\hat{\mathbf{M}}_k(i) \neq \mathbf{M}_{k,\text{true}}(i)\}}$.

130 The following Theorem implies that our method achieves clustering consistence.

131 **Theorem 3 (Clustering consistency)** *Suppose the Assumption (1) holds. Let $\hat{\mathbf{M}}_k$'s be the estima-*
 132 *tors from (6). Then the proportions of misclassified indices goes to zero in probability; i.e. there exist*
 133 *permutation matrices \mathbf{P}_k 's such that*

$$\sum_k \text{CER}(\hat{\mathbf{M}}_k, \mathbf{P}_k \mathbf{M}_{k,\text{true}}) \rightarrow 0, \quad \text{in probability.}$$

(add the proof) Under stronger distribution assumptions on \mathcal{E} , we can establish the finite-sample convergence rate for the clustering accuracy. See more results in the Supplements.

5 Numerical Implementation

5.1 Alternating optimization

We introduce an alternating optimization for solving (6). Note that the optimization (6) can be written as

$$(\hat{\mathcal{C}}, \{\hat{\mathbf{M}}_k\}) = \arg \min_{\mathcal{C} \in \mathbb{R}^{R_1 \times \dots \times R_K}, \text{ membership matrices } \mathbf{M}_k \text{'s}} f(\mathcal{C}, \{\mathbf{M}_k\}),$$

where $f(\mathcal{C}, \{\mathbf{M}_k\}) = \|\mathcal{Y} - \mathcal{C} \times_1 \mathbf{M}_1 \times_2 \dots \times_K \mathbf{M}_K\|_F^2$.

The decision variables consists of $K + 1$ blocks of variables, one for the core tensor \mathcal{C} and K for the membership matrices \mathbf{M}_k 's. We notice that, if any K out of the $K + 1$ blocks of variables are known, then the last block of variables can be solved explicitly. This observation suggests that we can iteratively update one block of variables at a time while keeping other others fixed. Specifically, given the collection of $\hat{\mathbf{M}}_k$'s, the core tensor estimate $\hat{\mathcal{C}} = \arg \min_{\mathcal{C}} f(\mathcal{C}, \{\hat{\mathbf{M}}_k\})$ collects the sample averages within each tensor block. Given the block mean $\hat{\mathcal{C}}$ and $K - 1$ membership matrices, the last membership matrix can be solved using simple nearest neighbor search over only R_k discrete points. The full procedure is described in Algorithm 1.

Algorithm 1 Multiway clustering based on tensor block models

Input: Data tensor $\mathcal{Y} \in \mathbb{R}^{d_1 \times \dots \times d_K}$, clustering size $\mathbf{R} = (R_1, \dots, R_K)$.

Output: Block mean tensor $\hat{\mathcal{C}} \in \mathbb{R}^{R_1 \times \dots \times R_K}$, and the membership matrices $\hat{\mathbf{M}}_k$.

1: Initialize the marginal clustering by performing independent k -means on each of the K modes.

2: **repeat**

3: Update the core tensor $\hat{\mathcal{C}} = \llbracket \hat{c}_{r_1, \dots, r_K} \rrbracket$. Specifically, for each $(r_1, \dots, r_K) \in [R_1] \times \dots \times [R_K]$,

$$\hat{c}_{r_1, \dots, r_K} = \frac{1}{n_{r_1, \dots, r_K}} \sum_{\hat{\mathbf{M}}_1^{-1}(r_1) \times \dots \times \hat{\mathbf{M}}_K^{-1}(r_K)} \mathcal{Y}_{i_1, \dots, i_K}, \quad (9)$$

where $n_{r_1, \dots, r_K} = \prod_k |\hat{\mathbf{M}}_k^{-1}(r_k)|$ is the number of entries in the block indexed by (r_1, \dots, r_K) .

4: Update membership matrices $\hat{\mathbf{M}}_k$'s:

5: **for** k in $\{1, 2, \dots, K\}$ **do**

6: Update the mode- k membership matrix $\hat{\mathbf{M}}_k$. Specifically, for each $a \in [d_k]$, assign the cluster label $\hat{\mathbf{M}}_k(a) \in [R_k]$ for which

$$\hat{\mathbf{M}}_k(a) = \arg \min_{r \in [R_k]} \sum_{\mathbf{I}_{-k}} \left(c_{\hat{\mathbf{M}}_1(i_1), \dots, r, \dots, \hat{\mathbf{M}}_K(i_K)} - \mathcal{Y}_{i_1, \dots, a, \dots, i_K} \right)^2,$$

where $\mathbf{I}_{-k} = (i_1, \dots, i_{k-1}, i_{k+1}, \dots, i_K)$ denotes the coordinates except the k -th mode.

7: **end for**

8: **until** Convergence

The above algorithm can be viewed as a higher-order extension for the ordinary (one-way) k -means algorithm. The core tensor \mathcal{C} serves as the role of centroids. As each iteration reduces the value of the objective function, which is bounded below, convergence of the algorithm is guaranteed. On the other hand, obtaining the global optimizer for such a non-convex optimization is typically difficult, even for the one-way k -means [8]. As such, we run the algorithm multiple times, using random initializations of the independent one-way k -mean on each of the K modes.

5.2 Tuning parameter selection

Our algorithm 1 takes the number of clusters \mathbf{R} as an input. In practice such information is often unknown and \mathbf{R} needs to be estimated from the data \mathcal{Y} . We propose to select this tuning parameter

157 using Bayesian information criterion (BIC),

$$\text{BIC}(\mathbf{R}) = \log \left(\|\mathcal{Y} - \hat{\Theta}\|_F \right) + \frac{\sum_k \log d_k}{\prod_k d_k} p_e, \quad (10)$$

158 where p_e is the effective number of parameters in the model. In our case we take $p_e = \prod_k R_k +$
 159 $\sum_k d_k \log R_k$, which is inspired from (8). We choose $\hat{\mathbf{R}}$ that minimizes $\text{BIC}(\mathbf{R})$ via grid search.
 160 Our choice of BIC is based on heretics, and we test its empirical performance in Section 7.

161 6 Extension to regularized estimation

162 In some large-scale applications, not every block in a data tensor is of equally importance. Examples
 163 include genome expression data, in which only a few entries represent the signals while the majority
 164 comes from the background noise (see Figure 2). Another example is community detection in a large
 165 sparse network, where only a few blocks represent the communities of interest and others represent
 166 small, noisy groups with weak connection. Although our estimator (6) can still handle this scenario
 167 by assigning small value to some of the $\hat{c}_{r_1, \dots, r_K}$'s, the estimates may suffer from high variance. It is
 168 thus beneficial to utilize regularized estimates for improved interpretability and better bias-variance
 169 trade-off.

170 Here we illustrate the regularized estimation using *sparsity* on the block means for localizing important
 171 blocks in the data tensor. This problem can be cast into variable selection on the block parameters.
 172 We propose the following regularized least-square estimation:

$$\hat{\Theta}^{\text{sparse}} = \arg \min_{\Theta \in \mathcal{P}} \left\{ \|\mathcal{Y} - \Theta\|_F^2 + \lambda \|\mathcal{C}\|_\rho \right\},$$

173 where $\|\mathcal{C}\|_\rho$ is the penalty function with ρ being an index for the tensor norm, $\mathcal{C} \in R^{R_1 \times \dots \times R_K}$ is
 174 the block means, and λ is the penalty tuning parameter. Some widely used penalties include Lasso
 175 penalty ($\rho = 1$), sparse subset penalty ($\rho = 0$), ridge penalty ($\rho = \text{Frobenius norm}$), elastic net
 176 (linear combination of $\rho = 1$ and $\rho = \text{Frobenius norm}$), among many others.

177 For parsimony purpose, we only discuss the properties for Lasso and sparse subset penalties here;
 178 other penalizations can be derived similarly. Sparse estimation incurs slight changes to Algorithm 1.
 179 When updating the core tensor \mathcal{C} in (9), we fit a penalized least square problem:

$$\hat{\mathcal{C}}^{\text{sparse}} = \arg \min_{\mathcal{C} \in \mathbb{R}^{R_1 \times \dots \times R_K}} \|\mathcal{Y} - \mathcal{C} \times_1 \hat{\mathbf{M}}_1 \times \dots \times_K \hat{\mathbf{M}}_K\|_F^2 + \lambda \|\mathcal{C}\|_\rho.$$

180 The closed-form update for the entry-wise sparse estimate $\hat{c}_{r_1, \dots, r_K}^{\text{sparse}}$ can be found (See Lemma 1 in
 181 Appendix):

$$\hat{c}_{r_1, \dots, r_K}^{\text{sparse}} = \begin{cases} \hat{c}_{r_1, \dots, r_K}^{\text{ols}} \mathbb{1}_{\{|\hat{c}_{r_1, \dots, r_K}^{\text{ols}}| \geq \frac{2\sqrt{\lambda}}{\sqrt{n_{r_1, \dots, r_K}}}\}} & \text{if } \rho = 1, \\ \text{sign}(\hat{c}_{r_1, \dots, r_K}^{\text{ols}}) \left(\hat{c}_{r_1, \dots, r_K}^{\text{ols}} - \frac{2\lambda}{n_{r_1, \dots, r_K}} \right) & \text{if } \rho = 0. \end{cases} \quad (11)$$

182 where $\hat{c}_{r_1, \dots, r_K}^{\text{ols}}$ denotes the ordinary least-square estimate in (9), for all $(r_1, \dots, r_K) \in [d_1] \times \dots \times [d_K]$.
 183 From Equation (11), we note that the penalization on \mathcal{C} shrink some block means to zero, and the
 184 shrinkage depends on the block size. For two blocks with same mean $\hat{c}_{r_1, \dots, r_K}^{\text{ols}}$, the smaller block
 185 receives a more stringent thresholding value $\lambda/n_{r_1, \dots, r_K}$. As such, our penalization encourages the
 186 identification of large blocks with high means.

187 In practice, the choice of penalty function ρ often depends on the goals and interpretations in specific
 188 problems. Given a penalization, we choose to select the tuning parameter λ via BIC (10), where
 189 we modify p_e into $p_e^{\text{sparse}} = \|\hat{\mathcal{C}}^{\text{sparse}}\|_0 + \sum_k d_k \log R_k$. Here $\|\hat{\mathcal{C}}^{\text{sparse}}\|_0$ denotes the total number of
 190 non-zero blocks in the estimates. The empirical performance of this proposal will be evaluated in
 191 Section 7.

192 7 Experiments

193 In this section, we evaluate the finite-sample performance of our method. We consider both non-sparse
 194 and sparse tensors, and compare our clustering performance with the other popular tensor-based

195 methods. The other methods we consider are (i) CP decomposition and (ii) Tucker decomposition,
 196 followed by k -means clusterings on the factors on each of the modes. We refer our method as TBM
 197 (tensor block model), and the other two methods as CP+ k means, and Tucker+ k -means. In the sparse
 198 case, we also consider (iii) sparse CP decomposition and (iv) sparse Tucker decomposition, followed
 199 by k -means on the sparse factors.

200 There are the statistics we would use to evaluate the performance in different cases:

201 (1) MSE: measuring the difference between the underlying means and the estimated means.
 202 (2) CER (clustering error rate): the adjusted rand index between two partitions. This statistic
 203 measures the agreement between the true partition and estimated partition of the data tensor.
 204 In this case, we have three kinds of CER in total: CER of mode 1, CER of mode 2 and CER of mode 3;
 205

206 For sparsity tensor, we further consider the following three statistics: (3) Total Correct Rate: 1 - the
 207 proportion of misjudgment while determining whether the mean signal is zero.;

208 (4) Correct Zero Rate: the proportion of zero elements are correctly identified in the underlying mean
 209 tensor;

210 (5) Correct One Rate: the proportion of non-zero elements are correctly identified in the underlying
 211 mean tensor.
 212

213 **Non-sparse case.** For non-sparse tensor, we assess the following four aspects.

214 first, we assess the relationship between MSE and the data size;

215 second, we verify the clustering approach when true d_1, \dots, d_K is given;

216 third, we evaluate the accuracy for estimating the number of clusters;

217 fifth, we would check the performance of our approach in a different kind of tensor and compare it
 218 with other clustering approach.

219 As for sparse tensor, we would evaluate the whole process: selecting d_1, \dots, d_K , choosing λ , doing
 220 tensor clustering.

221 Here we give a brief elaboration on the main way we use to generate the data. As for non-sparse
 222 tensor, given the cluster numbers d_1, \dots, d_K and the size of the tensor $n_1 n_2 n_3$, we assign the labels
 223 to each modes randomly. Next we randomly select the mean signal of clusters from Unif(-3,3) and
 224 add noise which comes from normal distribution with given standard deviation. Then we get the
 225 non-sparse tensor. As for sparse tensor, we randomly assign 0 to the mean of some clusters with
 226 given sparsity rate (the proportion of 0 elements) and then follow the same steps. We name this kind
 227 of tensor as tensor with constant clusters.

228 **Non-sparse case.** We begin with verifying the relationship between MSE and the sample size. The
 229 theoretical result indicates that the boundary of $RMSE = \sqrt{MSE} = \sqrt{\frac{error}{n_1 n_2 n_3}}$ decreases with

230 respect to sample size. Here we let n_1 to take values from 20 to 70, and $n_2 = \frac{n_1 \log d_1}{\log d_2}$, $n_3 = \frac{n_1 \log d_1}{\log d_3}$.

231 As for d_1, d_2, d_3 , we take them from $\{(4, 4, 4), (4, 8, 8), (8, 8, 8), (8, 8, 12)\}$. We simulated each
 232 situation 50 times and get the average RMSE among those cases. According to the panel (a) of Figure
 233 3, obviously, with sample size going up, the RMSE goes down. Additionally, the panel (b) of Figure
 234 3 indicates the RMSE decreases roughly at the rate of $1/N$ where $N = \sqrt{n_2 n_3 / \log d_1}$ is the rescaled
 235 sample size.

236 In the second simulation, we generate 50 non-sparse tensors with the same noise, size and cluster
 237 numbers each time. We use our approach (Algorithm 1) to do the clustering and the result is shown as
 238 Table 1. In both data size: $40 \times 40 \times 40$ and $40 \times 40 \times 80$, the CER on all modes are 0 when the noise is 4.
 239 As the noise goes up, the CER is increased gradually. Furthermore, from the $d_1 = 3, d_2 = 5, d_3 = 4$
 240 cases, we notice that the CER of mode i seems to be smaller when the number of clusters in the i th
 241 mode is less.

242 To evaluate the performance of our approach on selecting the number of clusters, we generate 50
 243 non-sparse tensors with the same noise, size and cluster numbers in each case in Table 2 as the third
 244 simulation. The reason why we only evaluate the performance of estimation on cluster numbers on
 245 non-sparse tensor is in sparse case, the reasonable cluster numbers may not be unique. As expected,

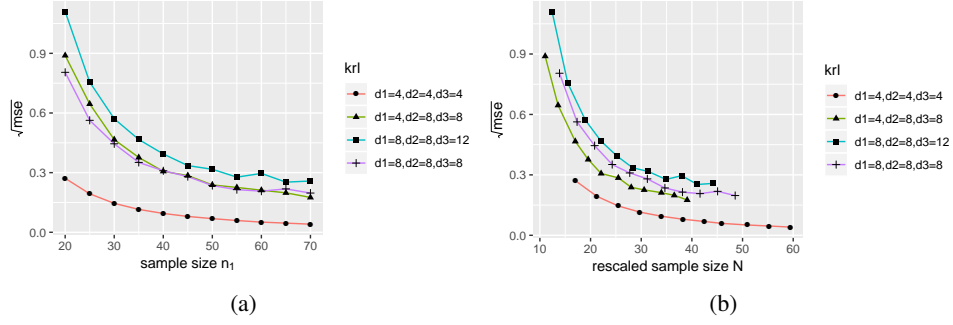


Figure 3: Plots of the root mean squared error (RMSE) versus sample size when using our tensor clustering algorithm. Each curve corresponds to a fixed (d_1, d_2, d_3) . (a): Plots of average RMSE over 50 simulations against n_1 ; (b): Plots of average RMSE over 50 simulations against $\sqrt{n_2 n_3} / \log d_1$.

n_1	n_2	n_3	noise	CER(mode 1)	CER(mode 2)	CER(mode3)
40	40	40	4	0(0)	0(0)	0(0)
40	40	40	8	0(0)	0.0136(0.0226)	0.0005(0.0036)
40	40	40	12	0.0365(0.0789)	0.12(0.0878)	0.0802(0.1009)
40	45	50	4	0(0)	0(0)	0(0)
40	45	50	8	0(0)	0.0027(0.0121)	0(0)
40	45	50	12	0.0158(0.0489)	0.0641(0.0629)	0.0336(0.0647)

Table 1: The CERs over 50 simulated tensors ($d_1 = 3, d_2 = 5, d_3 = 4$) each time.

we achieve 100% accuracy when the noise is 4, again. The overall accuracy goes down as the noise increased. Additionally, we notice that the smaller d_i is, the more accurate the estimated value is. There are two extremely low overall accuracy which appears when noise is 12 and the tensor size is $40 \times 40 \times 40$. However, the accuracy is improved quickly as the tensor size is enlarged to $40 \times 40 \times 80$. We found that not only the accuracy of mode 3 decreased after we add observations on mode 3, but also the accuracy of other two modes decreased a lot. The reason is the length of an observation in mode 1 and mode 2 is longer than before. Therefore, it is very important for us to get enough observations to guarantee the overall accuracy.

In the forth simulation, the true cluster numbers are not given, so we estimate them first and then use the estimated true cluster numbers to estimate the partition of clusters as well as underlying mean signals. We set the true cluster numbers to be $d_1 = 3, d_2 = 5, d_3 = 4$ specifically here, and the results are shown in Table 1. By looking into each mode separately, as the sample size of that mode increased, the CER of that mode decreased without any exception.

In the fifth simulation, we generate the tensor in two different ways. The new way to generate data is using $\sum_{s=1}^S d_s \mathbf{u}_s \mathbf{v}_s^T$ after giving the S . First we randomly choose \mathbf{u}_s and \mathbf{v}_s from normal distribution where $i = 1, \dots, S$. Then we add noise which is also sampled from normal distribution with given standard deviation. We name the tensor generated in this way as tensor with multiplicative clusters. In this simulation, we would check whether our method is still robust on the tensors with multiplicative cluster. Meanwhile, we would compare the performance of our approach with CPD k-means (do k-means clustering after using CP decomposition on the data tensor) under different noise and different kinds of tensor. To be more specific, when the data is a tensor with multiplicative clusters we just use the true rank that we use to generate the data tensor to do the CP k-means. As for the parameter d_1, d_2, d_3 in k-means, we also use the true parameter. We apply the true d_1, d_2, d_3 in our method, too. When the data is a tensor with constant clusters, we still use true d_1, d_2, d_3 in both methods. But for CPD k-means, we use BIC to select the best S according to the formula: $error + \frac{\log(n_1 n_2 n_3)(n_1 + n_2 + n_3 - 2)S}{n_1 n_2 n_3}$. Here we set $n_1 = 50, n_2 = 50, n_3 = 50$ and $d_1 = 4, d_2 = 4, d_3 = 4$ in all the cases, and add noise 0, 10, 20 to generate the tensors. Additionally, as for tensor with constant clusters, I randomly select underlying signals between -3 and 3; as for tensors with multiplicative clusters, I randomly select the elements of \mathbf{u}_s and \mathbf{v}_s between -3 and 3. Thus the mean signal of tensor with multiplicative clusters would be bigger than that of tensors with

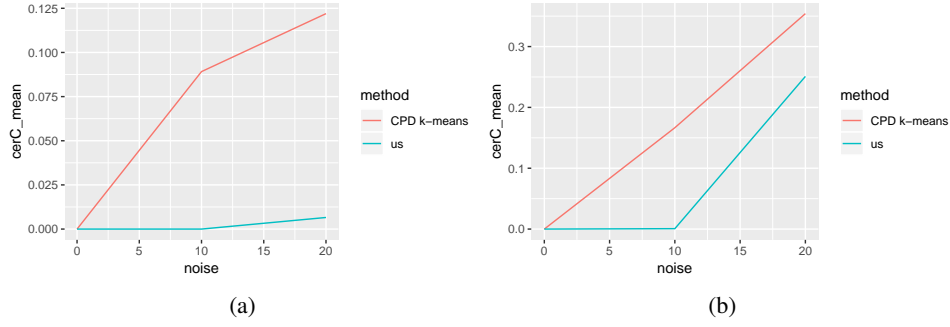


Figure 4: Plots of the mean CER C versus noise when using our tensor clustering algorithm and CPD k-means where $n_1 = n_2 = n_3 = 50, d_1 = d_2 = d_3 = 4$. Each curve corresponds to different method. (a): Plots of average CER in mode 1 over 50 simulations against noise while the data is a tensor with multiplicative clusters where $S = 3$; (b): Plots of average CER in mode 1 over 50 simulations against noise while the data is a tensor with constant clusters.

constant clusters. Therefore, the same noise would have bigger effect on the tensor with constant clusters. Our result is shown in Figure 4. To our surprise, our approach is not only better at working on tensor with constant clusters, but also better at working on tensor with multiplicative clusters. Furthermore, our approach shows more robustness under different level of noise, too.

Sparse case. We also test the performance of our approach under different λ when the data is sparse. The $\bar{\lambda}$ in Table 3 is the mean λ we choose across 50 simulations on the same sparsity rate. According to Table 3, the correct zero rate is increased with the increment on λ while the correct one rate is exactly the opposite. As for the λ we choose, it shows that the lowest total correct rate is 0.8586, which appears when noise is 8 and sparsity rate is 0.8. Overall, the λ we select by BIC is fairly good, which does works better than other λ .

8 Conclusion

Sparsity is only one form of regularization. In specific applications, prior knowledge often suggests various constraints among parameters, which may be exploited to regularize parameter estimates. For example, in the stochastic block model, sometimes it may be reasonable to impose symmetry on the parameters along certain subsets of modes, which further reduces the dimension of the problem. In some other applications, non-negativity of parameter of parameter values may be enforced. In our software, we implement the common penalizations but leave .. to further study.

References

- [1] Sara C Madeira and Arlindo L Oliveira. Biclustering algorithms for biological data analysis: a survey. *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, 1(1):24–45, 2004.
- [2] Anru Zhang and Dong Xia. Tensor SVD: Statistical and computational limits. *IEEE Transactions on Information Theory*, 2018.
- [3] Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009.
- [4] Robin A Darton. Rotation in factor analysis. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 29(3):167–194, 1980.
- [5] Hervé Abdi. Factor rotations in factor analyses.
- [6] Chao Gao, Yu Lu, Zongming Ma, and Harrison H Zhou. Optimal estimation and completion of matrices with biclustering structures. *The Journal of Machine Learning Research*, 17(1):5602–5630, 2016.

- 307 [7] Chao Gao and Zongming Ma. Minimax rates in network analysis: Graphon estimation, commu-
308 nity detection and hypothesis testing. *arXiv preprint arXiv:1811.06055*, 2018.
- 309 [8] Daniel Aloise, Amit Deshpande, Pierre Hansen, and Preyas Popat. Np-hardness of euclidean
310 sum-of-squares clustering. *Machine learning*, 75(2):245–248, 2009.