

---

# Tensor denoising and completion based on ordinal observations

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Higher-order tensors arise frequently in applications such as neuroimaging, recommendation system, social network analysis, and psychological studies. We consider the problem of low-rank tensor estimation from possibly incomplete, ordinal-valued observations. Two related problems are studied, one on tensor denoising and another on tensor completion. We propose a multi-linear cumulative link model, develop a rank-constrained M-estimator, and obtain theoretical accuracy guarantees. Our mean squared error bound enjoys a faster convergence rate than previous results, and we show that the proposed estimator is minimax optimal under the class of low-rank models. Furthermore, the procedure developed serves as an efficient completion method which guarantees consistent recovery of an order- $K$  ( $d, \dots, d$ )-dimensional low-rank tensor using only  $\tilde{O}(Kd)$  noisy, quantized observations. We demonstrate the outperformance of our approach over previous methods on the tasks of clustering and collaborative filtering.

## 1 Introduction

Multidimensional arrays, a.k.a. tensors, arise in a variety of applications including recommendation systems [Baltrunas et al., 2011], social networks [Nickel et al., 2011], genomics [Hore et al., 2016], and neuroimaging [Zhou et al., 2013]. Two main problems have gained increased attention for analyzing those noisy, high-dimensional datasets: tensor denoising, tensor completion. Tensor denoising aims to recover a signal tensor from its noisy entries [Xia and Zhou, 2019, Wang and Zeng, 2019]. Tensor completion examines the minimum number of entries needed for a consistent recovery [Ghadermarzy et al., 2018, 2019]. Low-rankness is often imposed to the signal tensor, thereby efficiently reducing the intrinsic dimension in both problems.

A number of low-rank tensor estimation methods have been proposed [Kolda and Bader, 2009, Acar et al., 2010], revitalizing classical methods such as CANDECOMP/PARAFAC (CP) decomposition [Hitchcock, 1927] and Tucker decomposition [Tucker, 1966]. These tensor methods treat the entries as continuous-valued. In many cases, however, we encounter datasets of which the entries are qualitative. For example, the Netflix problem records the ratings of users on movies over time. Each data is a rating on a nominal scale  $\{very\ like, like, neutral, dislike, very\ dislike\}$ . Another example is in the signal processing, where the digits are frequently rounded or truncated so that only integer values are available. Those qualitative observations take values in a limited set of categories, making the learning problem harder compared to continuous observations.

Ordinal entries are categorical variables with an ordering among the categories; for example,  $very\ like \prec like \prec neutral \prec \dots$ . The analyses of tensors with the ordinal entries are mainly complicated by two key properties needed for a reasonable model. First, the model should be invariant under a reversal of categories, say,  $very\ like \succ like \succ neutral \succ \dots$ , but not under arbitrary label permutations. Second, the parameter interpretations should be consistent under merging or splitting of contiguous categories. The classical continuous tensor model Kolda and Bader [2009], Ghadermarzy et al. [2019]

fails in the first aspect, whereas the binary tensor model Ghadermarzy et al. [2018] lacks the second property. An appropriate model for ordinal tensors has yet to be studied.

	Bhaskar [2016]	Ghadermarzy et al. [2018]	This paper
Higher-order tensors ( $K \geq 3$ )	✗	✓	✓
Multi-level categories ( $L \geq 3$ )	✓	✗	✓
Error rate for tensor denoising	$d^{-1}$ for $K = 2$	$d^{-(K-1)/2}$	$d^{-(K-1)}$
Optimality guarantee under low-rank models	unknown	✗	✓
Sample complexity for tensor completion	$d^K$	$Kd$	$Kd$

Table 1: Comparison with previous work. For ease of presentation, we summarize the error rate and sample complexity assuming equal tensor dimension in all modes.  $K$ : tensor order;  $L$ : number of ordinal levels;  $d$ : dimension at each mode.

**Our contribution.** This paper presents an efficient low-rank estimation method and theory for tensors with ordinal-valued entries. Our main contributions are summarized in Table 1. We propose a cumulative link model for higher-order tensors, develop a rank-constrained M-estimator, and obtain theoretical accuracy guarantees. We successfully recover true signal where neither the underlying signal nor the quantization operator is known. We show that our mean squared error bound has minimax optimal rate under the low-rank model. Our estimator enjoys a faster convergence rate  $\mathcal{O}(d^{-(K-1)/2})$  than  $\mathcal{O}(d^{-K})$  in Ghadermarzy et al. [2018], which is a substantial improvement as the order  $K$  increases. Proposed tensor completion algorithm guarantees consistent recovery of an order- $K$  ( $d, \dots, d$ )-dimensional low-rank tensor using only  $\tilde{\mathcal{O}}(Kd)$  noisy, quantized observations.

**Related work.** Our work is related to, but clearly distinctive from, several lines of existing literature. Matrix completion from quantized samples was firstly introduced for binary observations Cai and Zhou [2013], Davenport et al. [2014], Bhaskar and Javanmard [2015] and then extended to ordinal observations Bhaskar [2016]. As we show in Section 4, applying existing matrix methods to an ordinal tensor results in a suboptimal estimator with a slower convergence rate. Therefore, a full exploitation of the tensor structure is necessary; this is the focus of the current paper.

Our work is also connected to non-Gaussian tensor decomposition. Existing work focuses exclusively on univariate observations such as binary- or continuous-valued entries [Wang and Li, 2018, Hong et al., 2019, Ghadermarzy et al., 2018]. As we mentioned earlier, the ordinal observations add considerable challenges to the model formulation. We address the challenges from two perspectives. From statistical perspective, our proposed model generalizes the usual binary tensor model while preserving palindromic invariance McCullagh [1980] for ordinal observations. From algorithm perspective, our alternating optimization compares favorably to the approximated (convex) algorithm developed in the context of binary tensors Ghadermarzy et al. [2018]. We numerically compare the two approaches in Section 6.

## 2 Preliminaries

Let  $\mathcal{Y} \in \mathbb{R}^{d_1 \times \dots \times d_K}$  denote an order- $K$  ( $d_1, \dots, d_K$ )-dimensional tensor. We use  $y_\omega$  to denote the tensor entry indexed by  $\omega$ , where  $\omega \in [d_1] \times \dots \times [d_K]$ . The Frobenius norm of  $\mathcal{Y}$  is defined as  $\|\mathcal{Y}\|_F = \sum_\omega y_\omega^2$  and the infinity norm of  $\mathcal{Y}$  is defined as  $\|\mathcal{Y}\|_\infty = \max_\omega |y_\omega|$ . We use  $\mathcal{Y}_{(k)}$  to denote the unfolded matrix of size  $d_k$ -by- $\prod_{i \neq k} d_i$ , obtained by reshaping the tensor along the mode- $k$ . The Tucker rank of  $\mathcal{Y}$  is defined as a length- $K$  vector  $\mathbf{r} = (r_1, \dots, r_K)$ , where  $r_k$  is the rank of matrix  $\mathcal{Y}_{(k)}$  for all  $k \in [K]$ . We say that an event  $A$  occurs “with very high probability” if  $\mathbb{P}(A)$  tends to 1 faster than any polynomial of tensor dimension  $d_{\min} = \min\{d_1, \dots, d_K\} \rightarrow \infty$ . We use lower-case letters ( $a, b, c, \dots$ ) for scalars/vectors, upper-case boldface letters ( $\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$ ) for matrices, and calligraphy letters ( $\mathcal{A}, \mathcal{B}, \mathcal{C}, \dots$ ) for tensors of order three or greater. For ease of notation, we allow the basic arithmetic operators (e.g.,  $\leq, +, -$ ) to be applied to pairs of tensors in an element-wise manner. We use the shorthand  $[n]$  to denote the  $n$ -set  $\{1, \dots, n\}$  for  $n \in \mathbb{N}_+$ .

## 3 Model formulation and motivation

### 3.1 Observation model

Let  $\mathcal{Y}$  denote an order- $K$  ( $d_1, \dots, d_K$ )-dimensional data tensor. Suppose the entries of  $\mathcal{Y}$  are ordinal-valued, and the observation space is denoted by  $[L] := \{1, \dots, L\}$ . We propose a cumulative link model for the ordinal tensor  $\mathcal{Y} = \llbracket y_\omega \rrbracket \in [L]^{d_1 \times \dots \times d_K}$ . Specifically, assume the entries  $y_\omega$  are

81 (conditionally) independently distributed with cumulative probabilities:

$$\mathbb{P}(y_\omega \leq \ell) = f(b_\ell - \theta_\omega), \text{ for all } \ell \in [L - 1], \quad (1)$$

82 where  $\mathbf{b} = (b_1, \dots, b_{L-1})$  is a set of unknown scalars satisfying  $b_1 < \dots < b_{L-1}$ ,  $\Theta = \llbracket \theta_\omega \rrbracket \in$   
 83  $\mathbb{R}^{d_1 \times \dots \times d_K}$  is a continuous-valued parameter tensor satisfying certain low-rank structure (to be  
 84 specified later), and  $f(\cdot) : \mathbb{R} \mapsto [0, 1]$  is a known, strictly increasing function. We refer to  $\mathbf{b}$  as the  
 85 cut-off points and  $f$  the link function.

86 The formulation (1) imposes an additive model to the transformed probability of cumulative categories.  
 87 This modeling choice is to respect the ordering structure among the categories. For example, if  
 88 we choose the inverse link  $f^{-1}(x) = \log \frac{x}{1-x}$  to be the log odds, then the model (1) implies linear  
 89 spacing between the proportional odds:

$$\log \frac{\mathbb{P}(y_\omega \leq \ell)}{\mathbb{P}(y_\omega > \ell)} - \log \frac{\mathbb{P}(y_\omega \leq \ell - 1)}{\mathbb{P}(y_\omega > \ell - 1)} = b_\ell - b_{\ell-1}, \quad (2)$$

90 for all  $y_\omega$ . When there are only two categories in the observation space (e.g. binary tensors), the  
 91 cumulative model (1) is equivalent to the usual binomial link model. In general, however, when  
 92 the number of categories  $L \geq 3$ , the proportional odds assumption (2) is more parsimonious. The  
 93 ordered categories can be envisaged as contiguous intervals on the continuous scale, where the points  
 94 of division are exactly  $b_1 < \dots < b_{L-1}$ . This interpretation will be made explicit in the next section.

### 95 3.2 Latent-variable interpretation

96 The ordinal tensor model (1) with certain types of link  $f$  has the equivalent representation as an  
 97  $L$ -level quantization model on  $\mathcal{Y} = \llbracket y_\omega \rrbracket$ :

$$y_\omega = \sum_{\ell=1}^L \ell \mathbb{1}_{y_\omega^* \in (b_{\ell-1}, b_\ell]}, \quad (3)$$

98 for all  $\omega \in [d_1] \times \dots \times [d_K]$ . We define  $b_0 = -\infty$  and  $b_L = \infty$ . Here,  $\mathcal{Y}^* = \llbracket y_\omega^* \rrbracket$  is a latent  
 99 continuous-valued tensor following an additive noise model:

$$\underbrace{\mathcal{Y}^*}_{\text{latent continuous-valued tensor}} = \underbrace{\Theta}_{\text{signal tensor}} + \underbrace{\mathcal{E}}_{\text{i.i.d. noise}}, \quad (4)$$

100 where  $\mathcal{E} = \llbracket \varepsilon_\omega \rrbracket \in \mathbb{R}^{d_1 \times \dots \times d_K}$  is a noise tensor with i.i.d. entries according to distribution  $\mathbb{P}(\varepsilon)$ .  
 101 From the viewpoint of (4), the parameter tensor  $\Theta$  can be interpreted as the latent signal tensor prior  
 102 to contamination and quantization.

103 The equivalence between the latent-variable model (3) and the cumulative link model (1) is established  
 104 if the link  $f$  is chosen to be the cumulative distribution function of noise  $\varepsilon$ , i.e.,  $f(\theta) = \mathbb{P}(\varepsilon \leq \theta)$ .  
 105 We describe two common choices of link  $f$ , or equivalently, the distribution of  $\varepsilon$ .

106 **Example 1** (Logistic model). The logistic model is characterized by (1) with  $f(\theta) = (1 + e^{-\theta/\sigma})^{-1}$ ,  
 107 where  $\sigma > 0$  is the scale parameter. Equivalently, the noise  $\varepsilon_\omega$  in (3) follows i.i.d. logistic distribution  
 108 with scale parameter  $\sigma$ .

109 **Example 2** (Probit model). The probit model is characterized by (1) with  $f(\theta) = \mathbb{P}(z \leq \theta/\sigma)$ ,  
 110 where  $z \sim N(0, 1)$ . Equivalently, the noise  $\varepsilon_\omega$  in (3) follows i.i.d.  $N(0, \sigma^2)$ .

111 Other link functions are also possible, such as Laplace, Cauchy, etc McCullagh [1980]. All the  
 112 models share the property that the ordered categories can be thought of as contiguous interval on  
 113 some continuous scale. We point out that, although the latent-variable interpretation is incisive, our  
 114 estimation procedure does not refer to the existence of  $\mathcal{Y}^*$ . Therefore, our model (1) is general and  
 115 still valid in the absence of quantization process. More generally, we make the following assumptions  
 116 about the link  $f$ .

117 **Assumption 1.** The link function is assumed to satisfy:

- 118 1.  $f(\theta)$  is strictly increasing and twice-differentiable in  $\theta \in \mathbb{R}/\{0\}$ .
- 119 2.  $f'(\theta)$  is strictly log-concave and symmetric with respect to  $\theta = 0$ .

### 3.3 Problem 1: Tensor denoising

The first question we aim to address is tensor denoising:

(P1) Given the quantization process induced by  $f$  and the cut-off points  $\mathbf{b}$ , how accurately can we estimate the latent signal tensor  $\Theta$  from the ordinal observation  $\mathcal{Y}$ ?

Clearly, the problem (P1) cannot be solved uniformly for all possible  $\Theta$ . We focus on a class of “low-rank” and “flat” signal tensors, which is a plausible assumption in practical applications Zhou et al. [2013], Bhaskar and Javanmard [2015]. Specifically, we consider the parameter space:

$$\mathcal{P} = \left\{ \Theta \in \mathbb{R}^{d_1 \times \dots \times d_K} : \text{rank}(\Theta) \leq \mathbf{r}, \|\Theta\|_\infty \leq \alpha \right\}. \quad (5)$$

where  $\mathbf{r} = (r_1, \dots, r_K)$  denotes the Tucker rank of  $\Theta$ .

The parameter tensor of our interest satisfies two constraints. The first is that  $\Theta$  is a low-rank tensor, with  $r_k = \mathcal{O}(1)$  for all  $k \in [K]$ . Equivalently,  $\Theta$  admits the Tucker decomposition:

$$\Theta = \mathcal{C} \times_1 \mathbf{M}_1 \times_1 \dots \times_K \mathbf{M}_K, \quad (6)$$

where  $\mathcal{C} \in \mathbb{R}^{r_1 \times \dots \times r_K}$  is a core tensor,  $\mathbf{M}_k \in \mathbb{R}^{d_k \times r_k}$  are factor matrices with orthogonal columns, and  $\times_k$  denotes the tensor-by-matrix multiplication Kolda and Bader [2009]. The Tucker low-rankness is popularly imposed in tensor data analysis, and is shown to provide a reasonable tradeoff between model complexity and model flexibility. We choose Tucker representation for well-posedness of optimization and easy interpretation.

The second constraint is that the entries of  $\Theta$  are uniformly bounded in magnitude by a constant  $\alpha \in \mathbb{R}_+$ . In view of (4), we refer to  $\alpha$  as the signal level. The entry-wise bound assumption is a technical condition that avoids the degeneracy in probability estimation with ordinal observations.

### 3.4 Problem 2: Tensor completion

Motivated by applications in collaborative filtering, we also consider a more general setup when only a subset of tensor entries  $y_\omega$  are observed. Let  $\Omega \subset [d_1] \times \dots \times [d_K]$  denote the set of observed indices. The second question is stated as follows:

(P2) Given an incomplete set of ordinal observations  $\{y_\omega\}_{\omega \in \Omega}$ , how many sampled entries do we need to consistently recover  $\Theta$  based on the model (1)?

The answer to (P2) depends on the choice of  $\Omega$ . We consider a general model on  $\Omega$  that allows both uniform and non-uniform sampling. Specifically, let  $\Pi = \{\pi_{i_1, \dots, i_K}\}$  denote a predefine probability distribution over the index set such that  $\sum_{\omega \in [d_1] \times \dots \times [d_K]} \pi_\omega = 1$ . We assume that each index in  $\Omega$  is drawn with replacement using distribution  $\Pi$ . This sampling model relaxes the uniform sampling in literature and is arguably a better fit in applications.

We consider the same parameter space (5) for the completion problem. In addition to the reasons mentioned in Section 3.3, the entrywise bound assumption also serves as the incoherence requirement for completion. In classical matrix completion, the incoherence is often imposed on the singular vectors. This assumption is recently relaxed for “flat” matrices with bounded magnitude Negahban et al. [2011], Cai and Zhou [2013], Bhaskar and Javanmard [2015]. We adopt the same assumption for higher-order tensors.

## 4 Rank-constrained M-estimator

We present a general treatment to both problems mentioned above. With a little abuse of notation, we use  $\Omega$  to denote either the full index set  $\Omega = [d_1] \times \dots \times [d_K]$  (for the tensor denoising) or a random subset induced from the sampling distribution  $\Pi$  (for the tensor completion). Define  $f(-\infty) = 0$  and  $f(\infty) = 1$ . The log-likelihood associated with the observed entries is

$$\mathcal{L}_{\mathcal{Y}, \Omega}(\Theta, \mathbf{b}) = \sum_{\omega \in \Omega} \sum_{\ell \in [L]} \left\{ \mathbb{1}_{\{y_\omega = \ell\}} \log [f(b_\ell - \theta_\omega) - f(b_{\ell-1} - \theta_\omega)] \right\}. \quad (7)$$

We propose a rank-constrained maximum likelihood estimator (a.k.a. M-estimator) for  $\Theta$ :

$$\hat{\Theta} = \arg \max_{\Theta \in \mathcal{P}} \mathcal{L}_{\mathcal{Y}, \Omega}(\Theta, \mathbf{b}), \quad \text{where } \mathcal{P} = \left\{ \Theta \in \mathbb{R}^{d_1 \times \dots \times d_K} : \text{rank}(\Theta) \leq \mathbf{r}, \|\Theta\|_\infty \leq \alpha \right\}. \quad (8)$$

161 In practice, the cut-off points  $\mathbf{b}$  are unknown and should be jointly estimated with  $\Theta$ . For technical  
 162 convenience, we assume in this section that the cut-off points  $\mathbf{b}$  are known. The adaptation of  
 163 unknown  $\mathbf{b}$  is addressed in Section 5 and the Supplement.

164 We define a few key quantities that will be used in our theory. Let  $g_\ell = f(\theta + b_\ell) - f(\theta + b_{\ell-1})$  for  
 165 all  $\ell \in [L]$ , and

$$A_\alpha = \min_{\ell \in [L], |\theta| \leq \alpha} g_\ell(\theta), \quad U_\alpha = \max_{\ell \in [L], |\theta| \leq \alpha} \frac{|\dot{g}_\ell(\theta)|}{g_\ell(\theta)}, \quad L_\alpha = \min_{\ell \in [L], |\theta| \leq \alpha} \left[ \frac{\dot{g}_\ell^2(\theta)}{g_\ell^2(\theta)} - \frac{\ddot{g}_\ell(\theta)}{g_\ell(\theta)} \right],$$

166 where  $\dot{g}(\theta) = dg(\theta)/d\theta$ , and  $\alpha$  is the entrywise bound of  $\Theta$ . In view of equation (4), these quantities  
 167 characterize the geometry including flatness and convexity of the latent noise distribution. Under the  
 168 Assumption 1, all these quantities are strictly positive and independent of tensor dimension.

#### 169 4.1 Estimation error for tensor denoising

170 For the tensor denoising problem, we assume that the full set of tensor entries are observed. We  
 171 assess the estimation accuracy using the mean squared error (MSE):

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

172 The next theorem establishes the upper bound for the MSE of the proposed  $\hat{\Theta}$  in (8).

173 **Theorem 4.1** (Statistical convergence). *Consider an ordinal tensor  $\mathcal{Y} \in [L]^{d_1 \times \dots \times d_K}$  generated from*  
 174 *model (1), with the link function  $f$  and the true coefficient tensor  $\Theta^{\text{true}} \in \mathcal{P}$ . Define  $r_{\max} = \max_k r_k$ .*  
 175 *Then, with very high probability, the estimator in (8) satisfies*

$$\text{MSE}(\hat{\Theta}, \Theta^{\text{true}}) \leq \min \left( 4\alpha^2, \frac{c_1 U_\alpha^2 r_{\max}^{K-1} \sum_k d_k}{L_\alpha^2 \prod_k d_k} \right), \quad (9)$$

176 where  $c_1 > 0$  is a constant that depends only on  $K$ .

177 Theorem 4.1 establishes the statistical convergence for the estimator (8). In fact, the proof of this  
 178 theorem (see the Supplement) shows that the same statistical rate holds, not only for the global  
 179 optimizer (8), but also for any local optimizer  $\hat{\Theta}$  in the level set  $\{\hat{\Theta} \in \mathcal{P} : \mathcal{L}_{\mathcal{Y}, \Omega}(\hat{\Theta}) \geq \mathcal{L}_{\mathcal{Y}, \Omega}(\Theta^{\text{true}})\}$ .  
 180 This suggests that the local optimality itself is not necessarily a severe concern in our context, as long  
 181 as the convergent objective is large enough. In Section 5, we perform empirical studies to assess the  
 182 algorithmic stability.

183 A similar conclusion is obtained for the prediction error, measured in Kullback-Leibler (KL) diver-  
 184 gence, between the categorical distributions in the observation space.

185 **Corollary 1** (Prediction error). Assume the same set-up as in Theorem 4.1. Let  $\mathbb{P}_{\mathcal{Y}}$  and  $\hat{\mathbb{P}}_{\mathcal{Y}}$  denote  
 186 the distributions generating the  $L$ -level ordinal tensor  $\mathcal{Y}$ , given the true parameter  $\Theta$  and its estimator  
 187  $\hat{\Theta}$ , respectively. Assume  $L \geq 2$ . Then, with very high probability,

$$\text{KL}(\mathbb{P}_{\mathcal{Y}} \| \hat{\mathbb{P}}_{\mathcal{Y}}) \leq \frac{c_1 U_\alpha^2 r_{\max}^{K-1}}{L_\alpha^2} \frac{(4L-6)f^2(0) \sum_k d_k}{A_\alpha \prod_k d_k}, \quad (10)$$

188 where  $c_1 > 0$  is the same constants as in Theorem 4.1.

189 To gain insight into these bounds, we consider a special setting with equal dimension in all modes,  
 190 i.e.,  $d_1 = \dots = d_K = d$ . In such a case, our bound (9) reduces to

$$\text{MSE}(\hat{\Theta}, \Theta^{\text{true}}) \asymp d^{-(K-1)}, \quad \text{as } d \rightarrow \infty.$$

191 Hence, our estimator achieves consistency with polynomial convergence rate. We compare the bound  
 192 with existing literature. In the special case  $L = 2$ , Ghadermarzy et al. [2018] proposed a max-norm  
 193 constrained estimator  $\hat{\Theta}$  with  $\text{MSE}(\hat{\Theta}, \Theta^{\text{true}}) \asymp d^{-(K-1)/2}$ . In contrast, our estimator converges at a  
 194 rate of  $d^{-(K-1)}$ , which is substantially faster than theirs. This provides a positive answer to the open  
 195 question posed in Ghadermarzy et al. [2018] whether the square root in the bound is removable. The  
 196 improvement stems from utilizing the exact low-rankness of  $\Theta$ , whereas the surrogate rank measure  
 197 employed in Ghadermarzy et al. [2018] is scale-sensitive.

Our bound also generalizes the previous results on ordinal matrices. The convergence rate for rank-constrained matrix estimation was  $\mathcal{O}(1/\sqrt{d})$  [Bhaskar, 2016], which fits into our special case when  $K = 2$ . Furthermore, our results (9) and (10) reveal that the convergence becomes favorable as the order of data tensor increases. Intuitively, the sample size for tensor data analysis is the number of entries,  $\prod_k d_k$ , and the number of free parameters is roughly on the order of  $\sum_k d_k$ , assuming  $r_{\max} = \mathcal{O}(1)$ . A higher tensor order implies higher effective sample size per parameter, and thus exhibits a faster convergence rate in high dimensions.

We next show the statistical optimality of our estimator  $\hat{\Theta}$ . The result is based on the information theory, and applies to all estimators in  $\mathcal{P}$ , including but not limited to  $\hat{\Theta}$  in (8).

**Theorem 4.2** (Minimax lower bound). *Assume the same set-up as in Theorem 4.1, and  $d_{\max} = \max_k d_k \geq 8$ . Let  $\inf_{\hat{\Theta}}$  denote the infimum over all estimators  $\hat{\Theta} \in \mathcal{P}$  based on the ordinal tensor observation  $\mathcal{Y} \in [L]^{d_1 \times \dots \times d_K}$ . Then, under the model (1),*

$$\inf_{\hat{\Theta}} \sup_{\Theta^{true} \in \mathcal{P}} \mathbb{P} \left\{ \text{MSE}(\hat{\Theta}, \Theta^{true}) \geq c \min \left( \alpha^2, \frac{Cr_{\max} d_{\max}}{\prod_k d_k} \right) \right\} \geq \frac{1}{8},$$

where  $C = C(\alpha, L, f, \mathbf{b}) > 0$  and  $c > 0$  are constants independent of tensor dimension and the rank.

We see that the lower bound matches the upper bound in (9) on the polynomial order of tensor dimension. Therefore, our estimator (8) is order-optimal.

## 4.2 Sample complexity for tensor completion

We now consider the tensor completion problem, when only a subset of entries  $\Omega$  are observed. We consider a general sampling procedure induced by  $\Pi$ . The recovery accuracy is assessed by the weighted squared error:

$$\|\Theta - \hat{\Theta}\|_{F, \Pi}^2 \stackrel{\text{def}}{=} \frac{1}{|\Omega|} \mathbb{E}_{\Omega \sim \Pi} \|\Theta - \hat{\Theta}\|_F^2 = \sum_{\omega \in [d_1] \times \dots \times [d_K]} \pi_{\omega} (\Theta_{\omega} - \hat{\Theta}_{\omega})^2. \quad (11)$$

Note that the recovery error depends on the distribution  $\Pi$ . In particular, tensor entries with higher sampling probabilities have more influence on the recovery accuracy, compared to the ones with lower sampling probabilities.

**Remark 1.** If we assume each entry is sampled with strictly positive probability; i.e. there exists a constant  $\mu > 0$  s.t.

$$\pi_{\omega} \geq \frac{1}{\mu \prod_k d_k}, \quad \text{for all } \omega \in [d_1] \times \dots \times [d_K],$$

then the error in (11) provides an upper bound for MSE:

$$\|\Theta - \hat{\Theta}\|_{F, \Pi}^2 \geq \frac{\|\Theta - \hat{\Theta}\|_F^2}{\mu \prod_k d_k} = \frac{1}{\mu} \text{MSE}(\hat{\Theta}, \Theta^{true}).$$

The equality is attained under uniform sampling with  $\mu = 1$ .

**Theorem 4.3.** *Assume the same set-up as in Theorem 4.1. Suppose that we observe a subset of tensor entries  $\{y_{\omega}\}_{\omega \in \Omega}$ , where  $\Omega$  is chosen at random with replacement according to a probability distribution  $\Pi$ . Let  $\hat{\Theta}$  be the solution to (8), and assume  $r_{\max} = \mathcal{O}(1)$ . Then, with very high probability,*

$$\|\Theta - \hat{\Theta}\|_{F, \Pi}^2 \rightarrow 0, \quad \text{as } \frac{|\Omega|}{\sum_k d_k} \rightarrow \infty.$$

Theorem 4.3 shows that our estimator achieves consistent recovery using as few as  $\tilde{\mathcal{O}}(Kd)$  noisy, quantized observations from an order- $K$   $(d, \dots, d)$ -dimensional tensor. Note that  $\tilde{\mathcal{O}}(Kd)$  roughly matches the degree of freedom for an order- $K$  tensor of fixed rank  $r$ , suggesting the optimality of our sample requirement. This sample complexity substantially improves over earlier result  $\mathcal{O}(d^{\lceil K/2 \rceil})$  based on square matricization Mu et al. [2014], or  $\mathcal{O}(d^{N/2})$  based on tensor nuclear-norm regularization Yuan and Zhang [2016]. Existing methods that achieve  $\tilde{\mathcal{O}}(Kd)$  sample complexity require either a deterministic cross sampling design Zhang et al. [2019] or univariate measurements Ghadermarzy et al. [2018]. Our method extends the conclusions to multi-level measurements under a broader class of sampling schemes.

## 5 Numerical Implementation

We describe the algorithm to seek the optimizer of (7). In practice, the cut-off points  $\mathbf{b}$  are often unknown, so we choose to maximize  $\mathcal{L}_{\mathcal{Y},\Omega}$  jointly over  $(\Theta, \mathbf{b}) \in \mathcal{P} \times \mathcal{B}$ . The non convexity of feasible set  $\mathcal{P}$  makes the optimization (7) a non-convex problem. We employ the alternating optimization approach by utilizing the Tucker representation of  $\Theta$ . Specifically, based on (6) and (7), the objective function consists of  $K + 2$  blocks of variables, one for the cut-off points  $\mathbf{b}$ , one for the core tensor  $\mathcal{C}$ , and  $K$  for the factor matrices  $\mathbf{M}_k$ 's. The optimization is a simple convex problem if any  $K + 1$  out of the  $K + 2$  blocks are fixed. We update one block at a time while holding others fixed, and alternate the optimization throughout the iteration. The convergence is guaranteed whenever  $\mathcal{L}_{\mathcal{Y},\Omega}$  is bounded from above, since the alternating procedure monotonically increases the objective.

We comment on two implementation details before concluding this section. First, the problem (8) is non-convex, so the algorithm usually has no theoretical guarantee on global optimality. Nevertheless, as shown in Section 4.1, the desired rate holds not only for the global optimizer, but also for the local optimizer with  $\mathcal{L}_{\mathcal{Y},\Omega}(\hat{\Theta}) \geq \mathcal{L}_{\mathcal{Y},\Omega}(\Theta^{\text{true}})$ . In practice, we find the convergence point  $\hat{\Theta}$  upon random initialization is often satisfactory, in that the corresponding objective  $\mathcal{L}_{\mathcal{Y},\Omega}(\hat{\Theta})$  is close to and actually slightly larger than the objective evaluated at the true parameter  $\mathcal{L}_{\mathcal{Y},\Omega}(\Theta^{\text{true}})$ . Figure 5 shows the trajectory of the objective function from the algorithm with the input tensor generated from probit model (1) with  $d_1 = d_2 = d_3 = d$  and  $r_1 = r_2 = r_3 = r$ . The dashed line is the objective value at the true parameter  $\mathcal{L}_{\mathcal{Y},\Omega}(\Theta^{\text{true}})$ . The algorithm generally converges quickly to a desirable value in reasonable number of steps. The actual running time per iteration is shown in the plot legend.

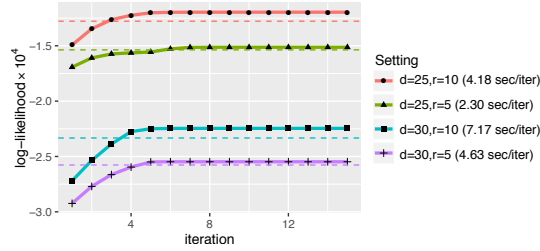


Figure 1: Trajectory of objective function with various  $d$  and  $r$ .

Second, the algorithm takes the rank  $r$  as an input. In practice, the rank  $r$  is hardly known and needs to be estimated from the data. We suggest to use Bayesian information criterion (BIC) and choose the rank that minimizes BIC; i.e.

$$\hat{r} = \arg \min_{r \in \mathbb{N}_+^K} \text{BIC}(r) = \arg \min_{r \in \mathbb{N}_+^K} \{-2\mathcal{L}_{\mathcal{Y}}(\hat{\Theta}(r), \hat{\mathbf{b}}(r)) + p_e(r) \log(\prod_k d_k)\},$$

where  $\hat{\Theta}(r), \hat{\mathbf{b}}(r)$  are the estimates given the rank  $r$ , and  $p_e(r) \stackrel{\text{def}}{=} \sum_k (d_k - r_k)r_k + \prod_k r_k$  is the effective number of parameters in the model. We select  $\hat{r}$  that minimizes BIC through a grid search. The choice of BIC is intended to balance between the goodness-of-fit for the data and the degrees of freedom in the population model.

## 6 Experiments

In this section, we evaluate the empirical performance of our method. We investigate both the complete and the incomplete settings, and compare the recovery accuracy with other tensor-based methods. Unless otherwise stated, the ordinal data tensors are generated from model (1) using standard probit link  $f$ . We consider the setting with  $K = 3$ ,  $d_1 = d_2 = d_3 = d$ , and  $r_1 = r_2 = r_3 = r$ . The parameter tensors are simulated based on (6), where the core tensor entries are i.i.d. drawn from  $N(0, 1)$ , and the factors  $\mathbf{M}_k$  are uniformly sampled from matrices with orthonormal columns. We set the cut-off points  $b_\ell = f^{-1}(\ell/L)$  for  $\ell \in [L]$ , such that  $f(b_\ell)$  are evenly spaced from 0 to 1. In each simulation study, we report the summary statistics across  $n_{\text{sim}} = 30$  replications.

### 6.1 Finite-sample performance

The first experiment examines the performance under complete observations. We assess the empirical relationship between the MSE and various aspects of model complexity, such as dimension  $d$ , rank  $r$ , and signal level  $\alpha = \|\Theta\|_\infty$ . Figure 2a plots the estimation error versus the tensor dimension

277  $d$  for ranks  $r \in \{3, 5, 8\}$ . The decay in the error appears to behave on the order of  $d^{-2}$ , which  
 278 is consistent with our theoretical results (9). We find that a higher rank leads to a larger error,  
 279 as reflected by the upward shift of the curve as  $r$  increases. Indeed, a higher rank implies the  
 280 more parameters to estimate, thus increasing the difficulty of the estimation. Figure 2b shows the  
 281 estimation error versus the signal level under  $d = 20$ . Interestingly, a larger estimation error is  
 282 observed when the signal is either too small or too large. The non-monotonic behavior may seem  
 283 surprising, but this is an intrinsic feature in the estimation with ordinal data. In view of the latent-  
 284 variable interpretation (see Section 3.2), estimation from ordinal observation can be interpreted as  
 285 an inverse problem of quantization. Therefore, the estimation error diverges in the absence of noise  
 286  $\mathcal{E}$ , because it is impossible to distinguish two different signal tensors, e.g.,  $\Theta_1 = \mathbf{a}_1 \otimes \mathbf{a}_2 \otimes \mathbf{a}_3$  and  
 287  $\Theta_2 = \text{sign}(\mathbf{a}_1) \otimes \text{sign}(\mathbf{a}_2) \otimes \text{sign}(\mathbf{a}_3)$ , from the quantized observations. This phenomenon Davenport  
 288 et al. [2014], Sur and Candès [2019] is contrary to the classical continuous-valued tensor problem.

289 The second experiment investigates the incomplete observations. We consider  $L$ -level tensors with  
 290  $d = 20$ ,  $\alpha = 10$  and choose a subset of tensor entries via uniform sampling. Figure 2c shows the  
 291 estimation error of  $\hat{\Theta}$  versus the fraction of observation  $\rho = |\Omega|/d^K$ . As expected, the error reduces  
 292 with increased  $\rho$  or decreased  $r$ . Figure 2d evaluates the impact of ordinal levels  $L$  to estimation  
 293 accuracy, under the setting  $\rho = 0.5$ . An improved performance is observed as  $L$  grows, especially  
 294 from binary observations ( $L = 2$ ) to multi-level ordinal observations ( $L \geq 3$ ). The result showcases  
 295 the benefit of multi-level observations compared to binary observations.

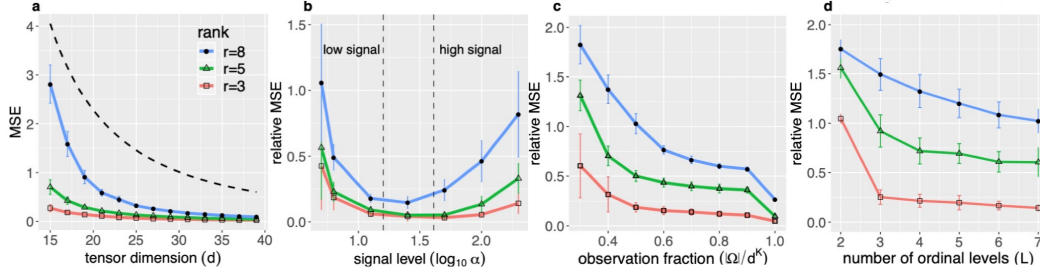


Figure 2: Empirical relationship between (relative) MSE versus (a) dimension  $d$ , (b) signal level  $\alpha$ , (c) observation fraction  $\rho$ , and (d) number of ordinal levels  $L$ . In panels (b)-(d), we plot the relative MSE  $= \|\hat{\Theta} - \Theta^{\text{true}}\|_F / \|\Theta^{\text{true}}\|_F$  for better visualization.

## 296 6.2 Comparison with alternative methods

297 Next, we compare our ordinal tensor method (**Ordinal-T**) with three popular low-rank methods:

- 298 • Continuous tensor decomposition (**Continuous-T**) Acar et al. [2010] is a low-rank approximation
- 299 method based on classical Tucker model.
- 300 • One-bit tensor completion (**1bit-T**) Ghadermarzy et al. [2018] is a max-norm penalized tensor
- 301 learning method based on partial binary observations.
- 302 • Ordinal matrix completion (**Ordinal-M**) Bhaskar [2016] is a rank-constrained matrix estimation
- 303 method based on noisy, quantized observations.

304 We apply each of the above methods to  $L$ -level ordinal tensors  $\mathcal{Y}$  generated from model (1). The  
 305 **Continuous-T** is applied directly to  $\mathcal{Y}$  by treating the  $L$  levels as continuous observations. The  
 306 **Ordinal-M** is applied to the matrix  $\mathcal{Y}_{(1)}$  obtained via 1-mode unfolding. The **1bit-T** is applied  
 307 to  $\mathcal{Y}$  in two ways. The first approach (**1bit-sign-T**) follows from Ghadermarzy et al. [2018] that  
 308 transforms  $\mathcal{Y}$  to a binary tensor, by taking the entrywise sign of the mean-adjusted tensor. The second  
 309 approach (**1bit-category-T**) transforms the order-3 ordinal tensor  $\mathcal{Y}$  to an order-4 binary tensor  
 310  $\mathcal{Y}^\# = \llbracket y_{ijkl}^\# \rrbracket$  via dummy variable encoding. We evaluate the methods by capabilities in predicting  
 311  $y_\omega^{\text{mode}} = \arg \max_\ell \mathbb{P}(y_\omega = \ell)$ . Two performance metrics are considered: mean absolute deviation,  
 312  $\text{MAD} = d^{-K} \sum_\omega |y_\omega^{\text{mode}} - \hat{y}_\omega^{\text{mode}}|$ , and misclassification rate,  $\text{MCR} = d^{-K} \sum_\omega \mathbb{1}_{\{y_\omega^{\text{mode}} \neq \text{round}(\hat{y}_\omega^{\text{mode}})\}}$ ,  
 313 where  $\text{round}(\cdot)$  gives the nearest integer of the prediction.

314 Figure 3 compares the prediction accuracy under the setting  $\alpha = 10$ ,  $d = 20$ , and  $r = 5$ . The  
 315 problem size we considered is comparable to Ghadermarzy et al. [2018]. We find that our method  
 316 outperforms the others in both MAD and MCR. In particular, methods built on multi-level observations  
 317 (**Ordinal-T**, **Ordinal-M**, **1bit-category-T**) exhibit stable MCR over  $\rho$  and  $L$ , whereas the others



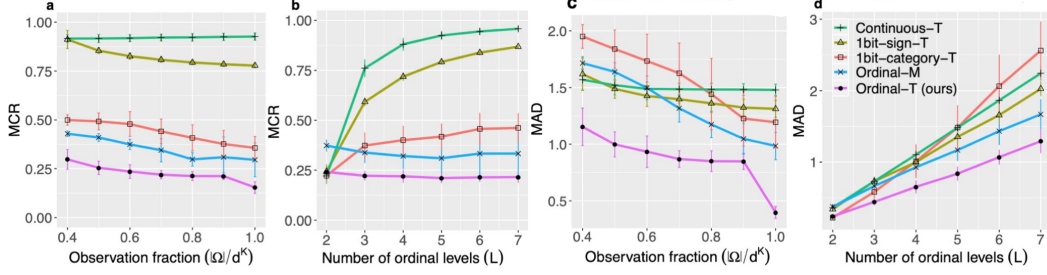


Figure 3: Performance comparison in MCR (a, b) and MAD (c, d). (b, d) Prediction errors versus the number of ordinal levels  $L$  when  $\rho = 0.8$ . (a, c) Prediction errors versus sample complexity  $\rho = |\Omega|/d^K$  when  $L = 5$ .

two methods (**Continuous-T**, **1bit-sign-T**) generally fail except for  $L = 2$  (Figures 3a-b). This observation highlights the necessity of modeling multi-level probabilities. Interestingly, although both **1bit-category-T** and our method **Ordinal-T** behave similarly for binary tensors ( $L = 2$ ), the improvement of our method is substantial as  $L$  increases (Figures 3a and 3c). One possible reason is that our method incorporates the intrinsic ordering among the  $L$  levels via proportional odds assumption (2), whereas **1bit-category-T** ignores the ordinal structure and dependence among the induced binary entries. Figures 3c-d assess the prediction accuracy with sample size. We see a clear advantage of our method (**Ordinal-T**) over the matricization (**Ordinal-M**) in both complete and non-complete observations. When the observation fraction is small, e.g.,  $|\Omega|/d^K = 0.4$ , the tensor-based completion shows  $\sim 30\%$  error reduction compared to the matricization.

## 7 Data Applications

We apply our method to two datasets. In the first application, we analyze an ordinal tensor consisting of structural connectivities among 68 brain regions from 136 individuals in Human Connectome Project (HCP) [Van Essen et al., 2013]. In the second application, we complete ordinal tensor data which records the ratings on a scale of 1 to 5 from 42 users to 139 songs on 26 contexts [Baltrunas et al., 2011] with missing values.

### 7.1 Human Connectome Project (HCP)

Each entry in the HCP dataset takes value on a nominal scale,  $\{high, moderate, low\}$ , indicating the strength level of fiber connection. We convert the dataset to a 3-level ordinal tensor  $\mathcal{Y} \in [3]^{68 \times 68 \times 136}$  and apply our method with a logistic link function. The BIC suggests  $\mathbf{r} = (23, 23, 8)$ . Based on the estimated factor matrices  $\{\hat{M}_k\}$ , we perform a clustering analysis via K-mean (see detailed procedure in the Supplement). The 68 brain nodes are grouped into 11 clusters, and the clustering captures the spatial separation of brain nodes (see the Supplement). In particular, the top three clusters represent the left, right and connection between two hemispheres; the smaller clusters represent local regions driving by similar nodes. We compare the goodness-of-fit of various tensor methods on the HPC data. Table 2 summarizes the prediction error via 5-fold stratified cross-validation averaged over 10 runs. Our method outperforms the others.

Method	Human Connectome Project (HCP) dataset		InCarMusic dataset	
	MAD	MCR	MAD	MCR
Ordinal-T (ours)	0.1607 (0.0005)	0.1606 (0.0005)	1.37 (0.039)	0.59 (0.009)
Continuous-T	0.2530 (0.0002)	0.1599 (0.0002)	2.39 (0.152)	0.94 (0.027)
1bit-sign-T	0.3566 (0.0010)	0.1563 (0.0010)	1.39 (0.003)	0.81 (0.005)

Table 2: Comparison of prediction error in the HPC and InCarMusic analyses. Standard errors are reported in parentheses.

### 7.2 InCarMusic recommendation system

We apply ordinal tensor completion to a recommendation system *InCarMusic*. *InCarMusic* is a mobile application that offers music recommendation to passengers based on contexts Baltrunas et al. [2011]. we conduct tensor completion on the  $42 \times 139 \times 26$  ordinal tensor with 2,844 observed entries. Table 2 shows the averaged prediction error via 5-fold cross validation. The high missing rate makes the accurate classification challenging. Nevertheless, our method achieves the best performance among the three.

## 8 Conclusions

We have developed a low-rank tensor estimation method based on possibly incomplete, ordinal-valued observations. A sharp error bound is established, and we demonstrate the outperformance of our approach compared to other methods. The work unlocks several directions of future research. One interesting question would be the inference problem, i.e., to assess the uncertainty of the obtained estimates and the imputation. Other directions include the trade-off between (non)convex optimization and statistical/computational efficiency. While convex relaxations are popular approach for matrix-based problem, they are often slow in practice Chen et al. [2019]. The interplay between computational efficiency and statistical accuracy in general tensor problems warrants future research.

## References

- Evrin Acar, Daniel M Dunlavy, Tamara G Kolda, and Morten Mørup. Scalable tensor factorizations with missing data. In *Proceedings of the 2010 SIAM international conference on data mining*, pages 701–712. SIAM, 2010.
- Linus Baltrunas, Marius Kaminskas, Bernd Ludwig, Omar Moling, Francesco Ricci, Aykan Aydin, Karl-Heinz Lüke, and Roland Schwaiger. Incarmusic: Context-aware music recommendations in a car. In *International Conference on Electronic Commerce and Web Technologies*, pages 89–100. Springer, 2011.
- Sonia A Bhaskar. Probabilistic low-rank matrix completion from quantized measurements. *The Journal of Machine Learning Research*, 17(1):2131–2164, 2016.
- Sonia A Bhaskar and Adel Javanmard. 1-bit matrix completion under exact low-rank constraint. In *2015 49th Annual Conference on Information Sciences and Systems (CISS)*, pages 1–6. IEEE, 2015.
- Tony Cai and Wen-Xin Zhou. A max-norm constrained minimization approach to 1-bit matrix completion. *The Journal of Machine Learning Research*, 14(1):3619–3647, 2013.
- Han Chen, Garvesh Raskutti, and Ming Yuan. Non-convex projected gradient descent for generalized low-rank tensor regression. *The Journal of Machine Learning Research*, 20(1):172–208, 2019.
- Mark A Davenport, Yaniv Plan, Ewout Van Den Berg, and Mary Wootters. 1-bit matrix completion. *Information and Inference: A Journal of the IMA*, 3(3):189–223, 2014.
- Navid Ghadermarzy, Yaniv Plan, and Ozgur Yilmaz. Learning tensors from partial binary measurements. *IEEE Transactions on Signal Processing*, 67(1):29–40, 2018.
- Navid Ghadermarzy, Yaniv Plan, and Özgür Yilmaz. Near-optimal sample complexity for convex tensor completion. *Information and Inference: A Journal of the IMA*, 8(3):577–619, 2019.
- Frank L Hitchcock. The expression of a tensor or a polyadic as a sum of products. *Journal of Mathematics and Physics*, 6(1-4):164–189, 1927.
- David Hong, Tamara G Kolda, and Jed A Duersch. Generalized canonical polyadic tensor decomposition. *SIAM Review*. In press. *arXiv:1808.07452*, 2019.
- Victoria Hore, Ana Viñuela, Alfonso Buil, Julian Knight, Mark I McCarthy, Kerrin Small, and Jonathan Marchini. Tensor decomposition for multiple-tissue gene expression experiments. *Nature genetics*, 48(9):1094, 2016.
- Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009.
- Peter McCullagh. Regression models for ordinal data. *Journal of the Royal Statistical Society: Series B (Methodological)*, 42(2):109–127, 1980.
- Cun Mu, Bo Huang, John Wright, and Donald Goldfarb. Square deal: Lower bounds and improved relaxations for tensor recovery. In *International Conference on Machine Learning*, pages 73–81, 2014.

398 Sahand Negahban, Martin J Wainwright, et al. Estimation of (near) low-rank matrices with noise and  
399 high-dimensional scaling. *The Annals of Statistics*, 39(2):1069–1097, 2011.

400 Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. A three-way model for collective learning  
401 on multi-relational data. In *International Conference on Machine Learning*, volume 11, pages  
402 809–816, 2011.

403 Pragma Sur and Emmanuel J Candès. A modern maximum-likelihood theory for high-dimensional  
404 logistic regression. *Proceedings of the National Academy of Sciences*, 116(29):14516–14525,  
405 2019.

406 Ledyard R Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):  
407 279–311, 1966.

408 David C Van Essen, Stephen M Smith, Deanna M Barch, Timothy EJ Behrens, Essa Yacoub, Kamil  
409 Ugurbil, Wu-Minn HCP Consortium, et al. The WU-Minn human connectome project: an overview.  
410 *Neuroimage*, 80:62–79, 2013.

411 Miaoyan Wang and Lexin Li. Learning from binary multiway data: Probabilistic tensor decomposition  
412 and its statistical optimality. *arXiv preprint arXiv:1811.05076*, 2018.

413 Miaoyan Wang and Yuchen Zeng. Multiway clustering via tensor block models. *Advances in Neural*  
414 *Information Processing Systems 32 (NeurIPS 2019)*. In press. *arXiv:1906.03807*, 2019.

415 Dong Xia and Fan Zhou. The sup-norm perturbation of hosvd and low rank tensor denoising. *Journal*  
416 *of Machine Learning Research*, 20(61):1–42, 2019.

417 Ming Yuan and Cun-Hui Zhang. On tensor completion via nuclear norm minimization. *Foundations*  
418 *of Computational Mathematics*, 16(4):1031–1068, 2016.

419 Anru Zhang et al. Cross: Efficient low-rank tensor completion. *The Annals of Statistics*, 47(2):  
420 936–964, 2019.

421 Hua Zhou, Lexin Li, and Hongtu Zhu. Tensor regression with applications in neuroimaging data  
422 analysis. *Journal of the American Statistical Association*, 108(502):540–552, 2013.