
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 [2], social networks [17], genomics [12], and neuroimaging [26]. 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 [23, 22]. Tensor completion examines the minimum number of entries needed for a consistent recovery [8, 9]. 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 [13, 1], revitalizing classical methods such as CANDECOMP/PARAFAC (CP) decomposition [10] and Tucker decomposition [19]. 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 $\{\textit{very like}, \textit{like}, \textit{neutral}, \textit{dislike}, \textit{very dislike}\}$. 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, $\textit{very like} \prec \textit{like} \prec \textit{neutral} \prec \dots$. 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, $\textit{very like} \succ \textit{like} \succ \textit{neutral} \succ \dots$, 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 [13, 9] fails in the first aspect, whereas the binary tensor model [8] 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 ($K \geq 3$)	\times	\checkmark	\checkmark
Multi-level categories ($L \geq 3$)	\checkmark	\times	\checkmark
Error rate for tensor denoising	d^{-1} for $K = 2$	$d^{-(K-1)/2}$	$d^{-(K-1)}$
Optimality guarantee under low-rank models	unknown	\times	\checkmark
Sample complexity for tensor completion	d^K	Kd	Kd

Table 1: Comparison with previous work. 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 a low-rank estimation method and theory for tensors with ordinal-valued entries. Our main contributions are summarized in Table 1. We are among the first to establish the recovery theory for both signal tensor and the quantization operators from a limited number of highly discrete entries. We show that our mean squared error bound is minimax optimal under the low-rank model. 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 [5, 7, 4] and then extended to ordinal observations [3]. 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 [21, 11, 8]. We address ordinal tensor observations from two perspective. From statistical perspective, our proposed model generalizes the binary tensor model while preserving palindromic invariance [14] for ordinal observations. From algorithm perspective, we address regularized MLE with alternating optimization (Non-convex). The approximated (convex) algorithm is developed in the context of binary tensors [8]. However, we advocate our approach and provide strong evidence. We show the improvement of error bound from $\mathcal{O}(d^{-(K-1)/2})$ to $\mathcal{O}(d^{-K})$ in [8] and numerically compare the two approaches.

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_ω to denote the tensor entry indexed by ω , 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. 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_ω are (conditionally) independently distributed with cumulative probabilities:

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

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

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

80 spacing between the proportional odds:

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

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

86 3.2 Latent-variable interpretation

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

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

89 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
 90 continuous-valued tensor following an additive noise model:

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

91 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 $\mathbb{P}(\varepsilon)$.
 92 From the viewpoint of (4), the parameter tensor Θ can be interpreted as the latent signal tensor prior
 93 to contamination and quantization.

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

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

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

101 Other link functions are also possible, such as Laplace, Cauchy, etc [14]. All the models share the
 102 property that the ordered categories can be thought of as contiguous interval on some continuous scale.
 103 We point out that, although the latent-variable interpretation is incisive, our estimation procedure
 104 does not refer to the existence of \mathcal{Y}^* . Therefore, our model (1) is general and still valid in the absence
 105 of quantization process. More generally, we make the following assumptions about the link f .

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

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

109 3.3 Problem 1: Tensor denoising

110 The first question we aim to address is tensor denoising:

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

113 We focus on a class of “low-rank” and “flat” signal tensors, which is a plausible assumption in
 114 practical applications [26, 4]. Specifically, we consider the parameter space:

$$\mathcal{P} = \left\{ \Theta \in \mathbb{R}^{d_1 \times \dots \times d_K} : \text{rank}(\Theta) \leq \mathbf{r}, \|\Theta\|_\infty \leq \alpha \right\}, \text{ where } \mathbf{r} = (r_1, \dots, r_K) \text{ denotes the Tucker rank of } \Theta.$$

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

$$\Theta = \mathcal{C} \times_1 \mathbf{M}_1 \times_1 \dots \times_K \mathbf{M}_K, \quad (5)$$

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 [13]. 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 Θ are uniformly bounded in magnitude by a constant $\alpha \in \mathbb{R}_+$. In view of (4), we refer to α 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_ω 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 Θ based on the model (1)?

The answer to (P2) depends on the choice of Ω . We consider a general model on Ω that allows both uniform and non-uniform sampling. Let $\Pi = \{\pi_{i_1, \dots, i_K}\}$ denote a predefined 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 Ω is drawn with replacement using distribution Π . This sampling model relaxes the uniform sampling in literature.

4 Rank-constrained M-estimator

We present a general treatment to both problems mentioned above. With a little abuse of notation, we use Ω 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 Π (for the tensor completion). Define $f(-\infty) = 0$ and $f(\infty) = 1$. The log-likelihood associated with the observed entries is

$$\mathcal{L}_{\mathcal{Y}, \Omega}(\Theta, \mathbf{b}) = \sum_{\omega \in \Omega} \sum_{\ell \in [L]} \left\{ \mathbb{1}_{\{y_\omega = \ell\}} \log [f(b_\ell - \theta_\omega) - f(b_{\ell-1} - \theta_\omega)] \right\}. \quad (6)$$

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

$$\hat{\Theta} = \arg \max_{\Theta \in \mathcal{P}} \mathcal{L}_{\mathcal{Y}, \Omega}(\Theta, \mathbf{b}), \quad \text{where } \mathcal{P} = \{\Theta \in \mathbb{R}^{d_1 \times \dots \times d_K} : \text{rank}(\Theta) \leq \mathbf{r}, \|\Theta\|_\infty \leq \alpha\}. \quad (7)$$

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

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 all $\ell \in [L]$, and

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

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

4.1 Estimation error for tensor denoising

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

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

The next theorem establishes the upper bound for the MSE of the proposed $\hat{\Theta}$ in (7).

Theorem 4.1 (Statistical convergence). *Consider an ordinal tensor $\mathcal{Y} \in [L]^{d_1 \times \dots \times d_K}$ generated from model (1), with the link function f and the true coefficient tensor $\Theta^{\text{true}} \in \mathcal{P}$. Define $r_{\max} = \max_k r_k$. Then, with very high probability, the estimator in (7) satisfies*

$$\text{MSE}(\hat{\Theta}, \Theta^{\text{true}}) \leq \min \left(4\alpha^2, \frac{c_1 U_\alpha^2 r_{\max}^{K-1} \sum_k d_k}{L_\alpha^2 \prod_k d_k} \right), \quad (8)$$

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

157 Theorem 4.1 establishes the statistical convergence for the estimator (7). In fact, the proof of this
 158 theorem (see the Supplement) shows that the same statistical rate holds, not only for the global
 159 optimizer (7), but also for any local optimizer $\hat{\Theta}$ in the level set $\{\hat{\Theta} \in \mathcal{P} : \mathcal{L}_{\mathcal{Y}, \Omega}(\hat{\Theta}) \geq \mathcal{L}_{\mathcal{Y}, \Omega}(\Theta^{\text{true}})\}$.
 160 This suggests that the local optimality itself is not necessarily a severe concern in our context, as long
 161 as the convergent objective is large enough.

162 To gain insight into these bounds, we consider a special setting with equal dimension in all modes,
 163 i.e., $d_1 = \dots = d_K = d$. In such a case, our bound (8) reduces to

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

164 Our estimator achieves consistency with polynomial convergence rate. We compare the bound with
 165 existing literature. In the special case $L = 2$, [8] proposed a max-norm constrained estimator $\hat{\Theta}$ with
 166 $\text{MSE}(\hat{\Theta}, \Theta^{\text{true}}) \asymp d^{-(K-1)/2}$. In contrast, our estimator converges at a rate of $d^{-(K-1)}$, which is
 167 much faster than theirs. This provides an answer to the open question posed in [8] whether the square
 168 root in the bound is removable. The improvement stems from utilizing the exact low-rankness of Θ ,
 169 whereas the surrogate rank measure employed in [8] is scale-sensitive.

170 Our bound also generalizes the previous results on ordinal matrices. The convergence rate for
 171 rank-constrained matrix estimation was $\mathcal{O}(1/\sqrt{d})$ [3], which fits into our special case when $K = 2$.
 172 Furthermore, (8) reveals that the convergence becomes favorable as the order of data tensor increases.
 173 Intuitively, the sample size for tensor analysis is the number of entries, $\prod_k d_k$, and the number of free
 174 parameters is roughly on the order of $\sum_k d_k$, assuming $r_{\max} = \mathcal{O}(1)$. A higher tensor order implies
 175 more effective sample size per parameter, and thus has a faster convergence rate in high dimensions.

176 We next show the statistical optimality of our estimator $\hat{\Theta}$. The result is based on the information
 177 theory, and applies to all estimators in \mathcal{P} , including but not limited to $\hat{\Theta}$ in (7).

178 **Theorem 4.2** (Minimax lower bound). *Assume the same set-up as in Theorem 4.1, and $d_{\max} =$
 179 $\max_k d_k \geq 8$. Let $\inf_{\hat{\Theta}}$ denote the infimum over all estimators $\hat{\Theta} \in \mathcal{P}$ based on the ordinal tensor
 180 observation $\mathcal{Y} \in [L]^{d_1 \times \dots \times d_K}$. Then, under the model (1),*

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

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

182 We see that the lower bound matches the upper bound in (8) on the polynomial order of tensor
 183 dimension. Therefore, our estimator (7) is order-optimal.

184 4.2 Sample complexity for tensor completion

185 We now consider the tensor completion problem, when only a subset of entries Ω are observed. We
 186 consider a general sampling procedure induced by Π . The recovery accuracy is assessed by the
 187 weighted squared error:

$$\|\Theta - \hat{\Theta}\|_{F, \Pi}^2 \stackrel{\text{def}}{=} \frac{1}{|\Omega|} \mathbb{E}_{\Omega \sim \Pi} \|\Theta - \hat{\Theta}\|_F^2 = \sum_{\omega \in [d_1] \times \dots \times [d_K]} \pi_{\omega} (\Theta_{\omega} - \hat{\Theta}_{\omega})^2. \quad (9)$$

188 Note that the recovery error depends on the distribution Π . In particular, tensor entries with higher
 189 sampling probabilities have more influence on the recovery accuracy, compared to the ones with
 190 lower sampling probabilities.

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

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

193 then the error in (9) provides an upper bound for MSE:

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

194 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 Ω is chosen at random with replacement according to a probability distribution Π . Let $\hat{\Theta}$ be the solution to (7), and assume $r_{\max} = \mathcal{O}(1)$. Then, with very high probability,

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

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 [15], or $\mathcal{O}(d^{N/2})$ based on tensor nuclear-norm regularization [24]. Existing methods that achieve $\tilde{\mathcal{O}}(Kd)$ sample complexity require either a deterministic cross sampling design [25] or univariate measurements [8]. Our method extends the conclusions to multi-level measurements under a broader class of sampling schemes.

5 Numerical Implementation

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 (6) a non-convex problem. We employ the alternating optimization approach by utilizing the Tucker representation of Θ and make each update convex problem. Multiple initializations are utilized, the sub-step for solving tensor factors is paralleled, and can be boosted by stochastic gradient descent. The detailed algorithm is in the Supplements.

The problem (7) 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. 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.

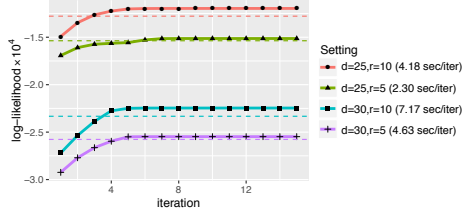


Figure 1: Trajectory of objective function with the input tensor generated from probit model (1) with $d_1 = d_2 = d_3 = d$ and $r_1 = r_2 = r_3 = r$. The dashed line is the objective value at the true parameter $\mathcal{L}_{\mathcal{Y}, \Omega}(\Theta^{\text{true}})$.

The algorithm 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 (5), where the core tensor entries are i.i.d. drawn from

236 $N(0, 1)$, and the factors \mathbf{M}_k are uniformly sampled from matrices with orthonormal columns. We
 237 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
 238 each simulation study, we report the summary statistics across $n_{\text{sim}} = 30$ replications.

239 6.1 Finite-sample performance

240 The first experiment examines the performance under complete observations. We assess the empirical
 241 relationship between the MSE and various model complexity, such as dimension d , rank r , and signal
 242 level $\alpha = \|\Theta\|_\infty$. Figure 2a plots the MSE versus the tensor dimension d for ranks $r \in \{3, 5, 8\}$.
 243 The decay in the error appears to behave on the order of d^{-2} , which is consistent with our theoretical
 244 results (8). We find that a higher rank leads to a larger error, as reflected by the upward shift of the
 245 curve as r increases. Indeed, a higher rank implies the more parameters to estimate, thus increasing
 246 the difficulty of the estimation. Figure 2b shows the MSE versus the signal level under $d = 20$.
 247 Interestingly, a larger estimation error is observed when the signal is either too small or too large. The
 248 non-monotonic behavior may seem surprising, but this is an intrinsic feature in the estimation with
 249 ordinal data. In view of the latent-variable interpretation (see Section 3.2), estimation from ordinal
 250 observation can be interpreted as an inverse problem of quantization. Therefore, the estimation error
 251 diverges in the absence of noise \mathcal{E} , because it is impossible to distinguish two different signal tensors,
 252 e.g., $\Theta_1 = \mathbf{a}_1 \otimes \mathbf{a}_2 \otimes \mathbf{a}_3$ and $\Theta_2 = \text{sign}(\mathbf{a}_1) \otimes \text{sign}(\mathbf{a}_2) \otimes \text{sign}(\mathbf{a}_3)$, from the quantized observations.
 253 This phenomenon [7, 18] is contrary to the classical continuous-valued tensor problem.

254 The second experiment investigates the incomplete observations. We consider L -level tensors with
 255 $d = 20$, $\alpha = 10$ and choose a subset of tensor entries via uniform sampling. Figure 2c shows the
 256 MSE of $\hat{\Theta}$ versus the fraction of observation $\rho = |\Omega|/d^K$. As expected, the error reduces with
 257 increased ρ or decreased r . Figure 2d evaluates the impact of ordinal levels L to estimation accuracy,
 258 under the setting $\rho = 0.5$. An improved performance is observed as L grows, especially from binary
 259 observations ($L = 2$) to multi-level ordinal observations ($L \geq 3$). The result showcases the benefit of
 multi-level observations compared to binary observations.

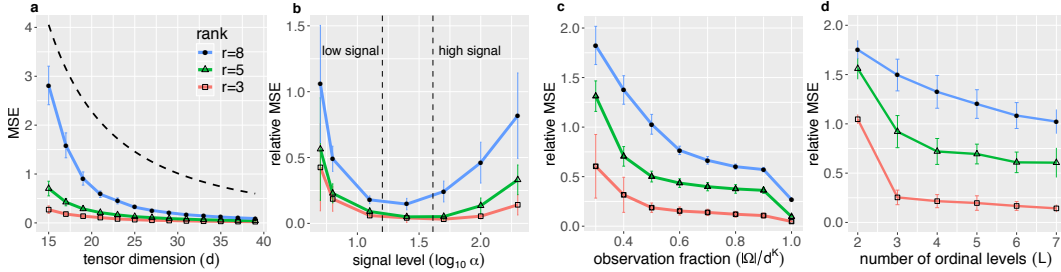


Figure 2: Empirical relationship between (relative) MSE versus (a) dimension d , (b) signal level α , (c) observation fraction ρ , and (d) number of ordinal levels L .

261 6.2 Comparison with alternative methods

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

- 263 • Continuous tensor decomposition (**Continuous-T**) [1] is a low-rank approximation method based
 264 on classical Tucker model.
- 265 • One-bit tensor completion (**1bit-T**) [8] is a max-norm penalized tensor learning method based on
 266 partial binary observations.
- 267 • Ordinal matrix completion (**Ordinal-M**) [3] is a rank-constrained matrix estimation method based
 268 on noisy, quantized observations.

269 We apply the each method to L -level ordinal tensors \mathcal{Y} generated from model (1). The **Continuous-T**
 270 is applied to \mathcal{Y} treating the L levels as continuous observations. The **Ordinal-M** is applied to the
 271 matrix $\mathcal{Y}_{(1)}$ obtained via 1-mode unfolding. The **1bit-T** is applied to \mathcal{Y} in two ways. The first
 272 approach (**1bit-sign-T**) follows from [8] that transforms \mathcal{Y} to a binary tensor, by taking the entrywise
 273 sign of the mean-adjusted tensor. The second approach (**1bit-category-T**) transforms the order-3
 274 ordinal tensor \mathcal{Y} to an order-4 binary tensor $\mathcal{Y}^\sharp = \llbracket y_{ijkl}^\sharp \rrbracket$ via dummy variable. We evaluate the
 275 method by mean absolute deviation (MAD) and misclassification rate (MCR).

276 Figure 3 compares the prediction accuracy under the setting $\alpha = 10$, $d = 20$, and $r = 5$. The problem
 277 size we considered is comparable to [8]. We find that our method outperforms the others in both
 278 MAD and MCR. In particular, methods built on multi-level observations (**Ordinal-T**, **Ordinal-M**,

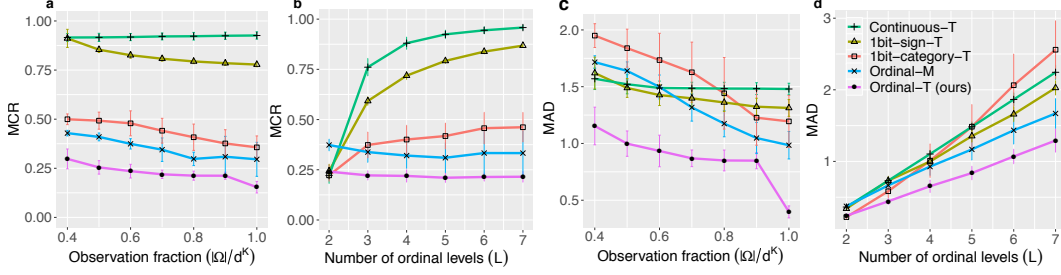


Figure 3: Performance comparison in MCR (a, b) and MAD (c, d). (b, d) Prediction errors versus the number of ordinal levels L when $\rho = 0.8$. (a, c) Prediction errors versus sample complexity $\rho = |\Omega|/d^K$ when $L = 5$.

1bit-category-T exhibit stable MCR over ρ and L , whereas the others 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**). 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) [20]. 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 $\{\hat{\mathbf{M}}_k\}$, we perform a clustering analysis via K-mean (see detailed procedure and results in Supplement). The 68 brain nodes are grouped into 11 clusters, and the clustering captures the spatial separation of brain nodes. 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 prediction performance with the alternative methods. 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 HCP and InCarMusic. Standard errors are reported in parentheses.

In the second application, we analyze ordinal tensor which records the ratings of 139 songs on a scale of 1 to 5 from 42 users on 26 contexts [2]. we conduct tensor completion on the $\mathcal{Y} \in [5]^{42 \times 139 \times 26}$ ordinal tensor with only 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.

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 over other alternative 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 [6]. The interplay between computational efficiency and statistical accuracy in general tensor problems warrants future research.

References

- [1] Evrim 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.
- [2] Linas Baltrunas, Marius Kaminskas, Bernd Ludwig, Omar Moling, Francesco Ricci, Aykan Aydin, Karl-Heinz Lücke, 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.
- [3] Sonia A Bhaskar. Probabilistic low-rank matrix completion from quantized measurements. *The Journal of Machine Learning Research*, 17(1):2131–2164, 2016.
- [4] 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.
- [5] 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.
- [6] 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.
- [7] 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.
- [8] Navid Ghadermarzy, Yaniv Plan, and Ozgur Yilmaz. Learning tensors from partial binary measurements. *IEEE Transactions on Signal Processing*, 67(1):29–40, 2018.
- [9] 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.
- [10] 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.
- [11] David Hong, Tamara G Kolda, and Jed A Duersch. Generalized canonical polyadic tensor decomposition. *SIAM Review*. In press. *arXiv:1808.07452*, 2019.
- [12] 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.
- [13] Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009.
- [14] Peter McCullagh. Regression models for ordinal data. *Journal of the Royal Statistical Society: Series B (Methodological)*, 42(2):109–127, 1980.
- [15] 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.
- [16] Sahand Negahban, Martin J Wainwright, et al. Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *The Annals of Statistics*, 39(2):1069–1097, 2011.
- [17] 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.
- [18] Pragya Sur and Emmanuel J Candès. A modern maximum-likelihood theory for high-dimensional logistic regression. *Proceedings of the National Academy of Sciences*, 116(29):14516–14525, 2019.

- 363 [19] Ledyard R Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31
364 (3):279–311, 1966.
- 365 [20] David C Van Essen, Stephen M Smith, Deanna M Barch, Timothy EJ Behrens, Essa Yacoub,
366 Kamil Ugurbil, Wu-Minn HCP Consortium, et al. The WU-Minn human connectome project:
367 an overview. *Neuroimage*, 80:62–79, 2013.
- 368 [21] Miaoyan Wang and Lexin Li. Learning from binary multiway data: Probabilistic tensor
369 decomposition and its statistical optimality. *arXiv preprint arXiv:1811.05076*, 2018.
- 370 [22] Miaoyan Wang and Yuchen Zeng. Multiway clustering via tensor block models. *Advances in*
371 *Neural Information Processing Systems 32 (NeurIPS 2019)*. In press. *arXiv:1906.03807*, 2019.
- 372 [23] Dong Xia and Fan Zhou. The sup-norm perturbation of hosvd and low rank tensor denoising.
373 *Journal of Machine Learning Research*, 20(61):1–42, 2019.
- 374 [24] Ming Yuan and Cun-Hui Zhang. On tensor completion via nuclear norm minimization. *Founda-*
375 *tions of Computational Mathematics*, 16(4):1031–1068, 2016.
- 376 [25] Anru Zhang et al. Cross: Efficient low-rank tensor completion. *The Annals of Statistics*, 47(2):
377 936–964, 2019.
- 378 [26] Hua Zhou, Lexin Li, and Hongtu Zhu. Tensor regression with applications in neuroimaging
379 data analysis. *Journal of the American Statistical Association*, 108(502):540–552, 2013.