

# Beyond the Signs: Nonparametric Tensor Completion via Sign Series

By Chanwoo Lee

Supporting document for the preliminary examination of  
PhD degree in Statistics

## Abstract

We consider the problem of tensor estimation from noisy observations with possibly missing entries. A nonparametric approach to tensor completion is developed based on a new model which we coin as sign representable tensors. The model represents the signal tensor of interest using a series of structured sign tensors. Unlike earlier methods, the sign series representation effectively addresses both low- and high-rank signals, while encompassing many existing tensor models—including CP models, Tucker models, single index models, structured tensors with repeating entries—as special cases. We provably reduce the tensor estimation problem to a series of structured classification tasks, and we develop a learning reduction machinery to empower existing low-rank tensor algorithms for more challenging high-rank estimation. Excess risk bounds, estimation errors, and sample complexities are established. We demonstrate the outperformance of our approach over previous methods on two datasets, one on human brain connectivity networks and the other on topic data mining.

## 1 Introduction

Higher-order tensors have recently received much attention in enormous fields including social networks (Anandkumar et al., 2014), neuroscience (Wang et al., 2017), and genomics (Hore et al., 2016). Tensor methods provide effective representation of the hidden structure in multiway data. In this paper we consider the signal plus noise model,

$$\mathcal{Y} = \Theta + \mathcal{E}, \tag{1}$$

where  $\mathcal{Y} \in \mathbb{R}^{d_1 \times \dots \times d_K}$  is an order- $K$  data tensor,  $\Theta$  is an unknown signal tensor of interest, and  $\mathcal{E}$  is a noise tensor. Our goal is to accurately estimate  $\Theta$  from the incomplete, noisy observation of  $\mathcal{Y}$ . In particular, we focus on the following two problems:

Q1 [Nonparametric tensor estimation]. How to flexibly estimate  $\Theta$  under a wide range of structures, including both low-rankness and high-rankness?

Q2 [Complexity of tensor completion]. How many observed tensor entries do we need to consistently estimate the signal  $\Theta$ ?

**Inadequacies of low-rank models.** The signal plus noise model (2) is popular in tensor literature. Existing methods estimate the signal tensor based on low-rankness of  $\Theta$  (Jain and Oh, 2014; Montanari and Sun, 2018). Common low-rank models include Canonical Polyadic (CP) tensors (Hitchcock, 1927), Tucker tensors (De Lathauwer et al., 2000), and block tensors (Wang and Zeng, 2019). While these methods have shown great success in theory, tensors in applications

often violate the low-rankness. Here we provide two examples to illustrate the limitation of classical models.

The first example reveals the sensitivity of tensor rank to order-preserving transformations. Let  $\mathcal{Z} \in \mathbb{R}^{30 \times 30 \times 30}$  be an order-3 tensor with  $\text{rank}(\mathcal{Z}) = 3$  (formal definition is deferred to the end of this section). Suppose a monotonic transformation  $f(z) = (1 + \exp(-cz))^{-1}$  is applied to  $\mathcal{Z}$  entrywise, and we let the signal  $\Theta$  in model (1) be the tensor after transformation. Figure 1a plots the numerical rank (see Section 7.1) of  $\Theta$  versus  $c$ . As we see, the rank increases rapidly with  $c$ , rendering traditional low-rank tensor methods ineffective in the presence of mild order-preserving nonlinearities. In digital processing (Ghadermarzy et al., 2018) and genomics analysis (Hore et al., 2016), the tensor of interest often undergoes unknown transformation prior to measurements. The sensitivity to transformation makes the low-rank model less desirable in practice.

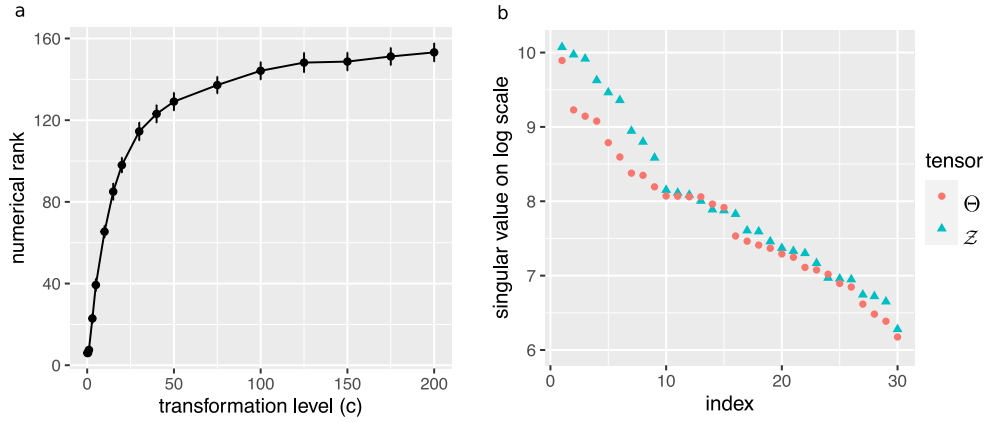


Figure 1: (a) Tensor rank vs.  $c$  in the first example. (b) Top  $d = 30$  tensor singular values in the second example.

The second example demonstrates the inadequacy of classical low-rankness in representing special structures. Here we consider the signal tensor of the form  $\Theta = \log(1 + \mathcal{Z})$ , where  $\mathcal{Z} \in \mathbb{R}^{d \times d \times d}$  is an order-3 tensor with entries  $\mathcal{Z}(i, j, k) = d^{-1} \max(i, j, k)$  for  $i, j, k \in \{1, \dots, d\}$ . The matrix analogy of  $\Theta$  was studied in graphon analysis Chan and Airolidi (2014). In this case neither  $\Theta$  nor  $\mathcal{Z}$  is low-rank; in fact, the rank is no smaller than the dimension  $d$  as illustrated in Figure 1b. Again, classical low-rank models fail to address this type of tensor structure.

In the above and many other examples, the signal tensors  $\Theta$  of interest have high rank. Classical low-rank models will miss these important structures. The observations have motivated us to develop more flexible tensor modeling.

**Our contributions.** We develop a new model called sign representable tensors to address the aforementioned challenges. Figure 2 illustrates our main idea. Our approach is built on the sign series representation of the signal tensor, and we propose to estimate the sign tensors through a series of weighted classifications. In contrast to existing methods, our method is guaranteed to recover a wide range of low- and high-rank signals. We highlight two main contributions that set our work apart from earlier literature.

Statistically, the problem of high-rank tensor estimation is challenging. Existing estimation theory (Anandkumar et al., 2014; Montanari and Sun, 2018; Cai et al., 2019) exclusively focuses on the regime of fixed  $r$  growing  $d$ . However, such premise fails in high-rank tensors, where the rank may grow with, or even exceed, the dimension. A proper notion of nonparametric complexity is crucial. We show that, somewhat surprisingly, the sign tensor series not only preserves all information in

the original signals, but also brings the benefits of flexibility and accuracy over classical low-rank models. The results fill the gap between parametric (low-rank) and nonparametric (high-rank) tensors, thereby greatly enriching the tensor model literature.

Computationally, a number of polynomial-time algorithms are readily available for 1-bit tensor estimation Wang and Li (2020); Han et al. (2020); Ghadermarzy et al. (2018). These algorithms enjoy computational efficiency while being restricted to binary inputs. Our work is orthogonal to these algorithm development, and we show that the high-rank tensor estimate is provably reducible to a series of binary tensor problems with carefully-designed weights. This reduction provides a generic engine to empower existing algorithms for a wider range of structured tensor problems. We use a divide-and-concur approach to combine efficient base algorithms, thereby achieving computational accuracy without the need to reinvent the wheel. The flexibility to import and adapt existing tensor algorithms is one advantage of our method.

We also highlight the challenges associated with tensors compared to matrices. High-rank matrix estimation is recently studied under nonlinear models (Ganti et al., 2015) and subspace clustering (Ongie et al., 2017; Fan and Udell, 2019). However, high-rank tensor problems is more challenging, because the tensor rank often exceeds the dimension when order  $K$  greater than two (Anandkumar et al., 2017). This is in sharp contrast to matrices ( $K = 2$ ). We show that, applying matrix methods to higher-order tensors results in suboptimal estimates. A full exploitation of the higher-order structure is needed; this is another challenge we address in this paper.

**Notation.** We use  $[n] = \{1, \dots, n\}$  for  $n$ -set with  $n \in \mathbb{N}_+$ ,  $a_n \lesssim b_n$  if  $\lim_{n \rightarrow \infty} a_n/b_n \leq c$  for some constant  $c > 0$ , and  $a_n \asymp b_n$  if  $c_1 \leq \lim_{n \rightarrow \infty} a_n/b_n \leq c_2$  for some constants  $c_1, c_2 > 0$ . We use  $\mathcal{O}(\cdot)$  to denote the big-O notation,  $\tilde{\mathcal{O}}(\cdot)$  the variant hiding logarithmic factors. Let  $\Theta \in \mathbb{R}^{d_1 \times \dots \times d_K}$  denote an order- $K$   $(d_1, \dots, d_K)$ -dimensional tensor, and  $\Theta(\omega) \in \mathbb{R}$  denote the tensor entry indexed by  $\omega \in [d_1] \times \dots \times [d_K]$ . An event  $A$  is said to occur “with very high probability” if  $\mathbb{P}(A)$  tends to 1 faster than any polynomial of tensor dimension  $d := \min_k d_k \rightarrow \infty$ . The tensor rank (Hitchcock, 1927) is defined by  $\text{rank}(\Theta) = \min\{r \in \mathbb{N} : \Theta = \sum_{s=1}^r \mathbf{a}_s^{(1)} \otimes \dots \otimes \mathbf{a}_s^{(K)}\}$ , where  $\mathbf{a}_s^{(k)} \in \mathbb{R}^{d_k}$  are vectors for  $k \in [K], s \in [r]$ , and  $\otimes$  denotes the outer product of vectors. We use  $\text{sgn}(\cdot) : \mathbb{R} \rightarrow \{-1, 1\}$  to denote the sign function, where  $\text{sgn}(y) = 1$  if  $y \geq 0$  and  $-1$  otherwise. We allow univariate functions, such as  $\text{sgn}(\cdot)$  and general  $f : \mathbb{R} \rightarrow \mathbb{R}$ , to be applied to tensors in an element-wise manner.

## 2 Model and proposal overview

Let  $\mathcal{Y}$  be an order- $K$   $(d_1, \dots, d_K)$ -dimensional tensor generated from the model

$$\mathcal{Y} = \Theta + \mathcal{E}, \quad (2)$$

where  $\Theta \in \mathbb{R}^{d_1 \times \dots \times d_K}$  is an unknown signal tensor, and  $\mathcal{E}$  is a noise tensor consisting of zero-mean, independent but not necessarily identically distributed entries. We allow heterogenous noise, in that the marginal distribution of noise entry  $\mathcal{E}(\omega)$  may depend on  $\omega$ . For a cleaner exposition, we assume the noise is bounded and the range of  $Y$  is in  $[-1, 1]$ ; the extension to a sub-Gaussian noise is provided in Section 7.4. Our observation is an incomplete data tensor from (2), denoted  $\mathcal{Y}_\Omega$ , where  $\Omega \subset [d_1] \times \dots \times [d_K]$  is the index set of observed entries. We consider a general model on  $\Omega$  that allows both uniform and non-uniform samplings. Specifically, let  $\Pi = \{p_\omega\}$  be an arbitrarily predefined probability distribution over the full index set with  $\sum_{\omega \in [d_1] \times \dots \times [d_K]} p_\omega = 1$ . We use  $\omega \sim \Pi$  to denote the sampling rule, meaning  $\omega$  in  $\Omega$  are i.i.d. draws with replacement from distribution  $\Pi$ . The goal is to estimate  $\Theta$  from  $\mathcal{Y}_\Omega$ . Note that  $\Theta$  is not necessarily low-rank.

**Proposal intuition.** Before describing our main results, we provide the intuition behind our method. In the two examples in Section 1, the high-rankness in the signal  $\Theta$  makes the estimation challenging. Now let us examine the sign of the  $\pi$ -shifted signal  $\text{sgn}(\Theta - \pi)$  for any given  $\pi \in [-1, 1]$ . It turns out that, these sign tensors share the same sign patterns as low-rank tensors. Indeed, the signal tensor in the first example has the same sign pattern as a rank-4 tensor, since  $\text{sgn}(\Theta - \pi) = \text{sgn}(\mathcal{Z} - f^{-1}(\pi))$ . The signal tensor in the second example has the same sign pattern as a rank-2 tensor, since  $\text{sgn}(\Theta(i, j, k) - \pi) = \text{sgn}(\max(i, j, k) - d(e^\pi - 1))$  (see Example 5 in Section 3).

The above observation suggests a general framework to estimate both low- and high-rank signal tensors. Figure 2 illustrates the main crux of our method. We propose to estimate the signal tensor  $\Theta$  by taking the average over structured sign tensors

$$\hat{\Theta} = \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{sgn}(\hat{\mathcal{Z}}_\pi), \text{ where } \hat{\mathcal{Z}}_\pi = \arg \min_{\text{low rank tensor } \mathcal{Z}} \text{Weighted-Loss}(\text{sgn}(\mathcal{Z}), \text{sgn}(\mathcal{Y}_\Omega - \pi)). \quad (3)$$

Here  $\text{sgn}(\hat{\mathcal{Z}}_\pi) \in \{-1, 1\}^{d_1 \times \dots \times d_K}$  is the sign tensor estimated at a series of  $\pi \in \mathcal{H} = \{-1, \dots, -1/H, 0, 1/H, \dots, 1\}$ , and  $\text{Weighted-Loss}(\cdot, \cdot)$  denotes a classification objective function with an entry-specific weight to each tensor entry; its specific form will be described in Section 3.2. To obtain  $\text{sgn}(\hat{\mathcal{Z}}_\pi)$  for a given  $\pi$ , we propose to dichotomize the data tensor into a sign tensor  $\text{sgn}(\mathcal{Y}_\Omega - \pi)$  and estimate the de-noised sign by performing weighted classification.

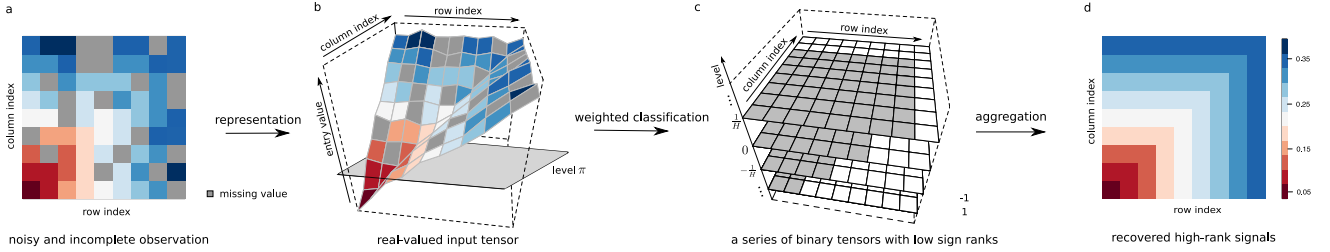


Figure 2: Illustration of our method in the context of an order-2 tensor (a.k.a. matrix). (a): a noisy, incomplete tensor input. (b)-(c): Estimation of sign tensor series  $\text{sgn}(\Theta - \pi)$  for  $\pi \in \{-1, \dots, -1/H, 0, 1/H, \dots, 1\}$ . (d): recovered signal  $\hat{\Theta}$ . The depicted signal is a full-rank matrix based on Example 5 in Section 3.

Our approach is built on the nonparametric sign representation of signal tensors. We show that a careful aggregation of dichotomized data not only preserves all information in the original signals, but also brings benefits of accuracy and flexibility over classical low-rank models. Unlike traditional methods, the sign representation is guaranteed to recover both low- and high-rank signals. In addition, a total of  $H = \text{poly}(d)$  dichotomized problems suffice to recover  $\Theta$  under the considered model. The method therefore enjoys both statistical effectiveness and computational efficiency.

### 3 Oracle properties of sign representable tensors

This section develops sign representable tensor models for  $\Theta$  in (2). We characterize the algebraic and statistical properties of sign tensor series, which serves the foundation for our method.

#### 3.1 Sign-rank and sign tensor series

Let  $\Theta$  be the tensor of interest, and  $\text{sgn}(\Theta)$  the corresponding sign pattern. The sign patterns induce an equivalence relationship between tensors. Two tensors are called sign equivalent, denoted  $\simeq$ , if they have the same sign pattern.

**Definition 1** (Sign-rank). The sign-rank of a tensor  $\Theta \in \mathbb{R}^{d_1 \times \dots \times d_K}$  is defined by the minimal rank among all tensors that share the same sign pattern as  $\Theta$ ; i.e.,

$$\text{srnk}(\Theta) = \min\{\text{rank}(\Theta') : \Theta' \simeq \Theta, \Theta' \in \mathbb{R}^{d_1 \times \dots \times d_K}\}.$$

This concept is important in combinatorics (Cohn and Umans, 2013), complexity theory (Alon et al., 2016), and quantum mechanics (De Wolf, 2003); we extend the notion to continuous-valued tensors. Note that the sign-rank concerns only the sign pattern but discards the magnitude information of  $\Theta$ . In particular,  $\text{srnk}(\Theta) = \text{srnk}(\text{sgn}\Theta)$ .

Like most tensor problems (Hillar and Lim, 2013), determining the sign-rank is NP hard in the worst case (Alon et al., 2016). Fortunately, tensors arisen in applications often possess special structures that facilitate analysis. The sign-rank is upper bounded by tensor rank. More generally, we show the following properties.

**Proposition 1** (Upper bounds of the sign-rank).

- (a) [Upper bounds] For any strictly monotonic function  $g: \mathbb{R} \rightarrow \mathbb{R}$  with  $g(0) = 0$ , we have  $\text{srnk}(\Theta) \leq \text{rank}(g(\Theta))$ .
- (b) [Broadness] For every order  $K \geq 2$  and dimension  $d$ , there exist tensors  $\Theta \in \mathbb{R}^{d \times \dots \times d}$  such that  $\text{rank}(\Theta) \geq d$  but  $\text{srnk}(\Theta - \pi) \leq 2$  for all  $\pi \in \mathbb{R}$ .

Proposition 1 demonstrates the strict broadness of low sign-rank family over the usual low-rank family. In particular, the sign-rank can be much smaller than the tensor rank, as we have shown in the two examples of Section 1. We provide several additional examples in Section 7.2 in which the tensor rank grows with dimension  $d$  but the sign-rank remains a constant. The results highlight the advantages of using sign-rank in the high-dimensional tensor analysis.

We now introduce a tensor family, which we coin as “sign representable tensors”.

**Definition 2** (Sign representable tensors). Fix a level  $\pi \in [-1, 1]$ . A tensor  $\Theta$  is called  $(r, \pi)$ -sign representable, if the tensor  $(\Theta - \pi)$  has sign-rank bounded by  $r$ . A tensor  $\Theta$  is called  $r$ -sign (globally) representable, if  $\Theta$  is  $(r, \pi)$ -sign representable for all  $\pi \in [-1, 1]$ . The collection  $\{\text{sgn}(\Theta - \pi) : \pi \in [-1, 1]\}$  is called the sign tensor series. We use  $\mathcal{P}_{\text{sgn}}(r) = \{\Theta : \max_{\pi \in [-1, 1]} \text{srnk}(\Theta - \pi) \leq r\}$  to denote the  $r$ -sign representable tensor family.

We next show that the  $r$ -sign representable tensor family is a general model that incorporates most existing tensor models, including low-rank tensors, single index models, GLM models, and structured tensors with repeating entries.

**Example 1** (CP/Tucker low-rank models). The CP and Tucker low-rank tensors are the two most popular tensor models (Kolda and Bader, 2009). Let  $\Theta$  be a low-rank tensor with CP rank  $r$ . We see that  $\Theta$  belongs to the sign representable family; i.e.,  $\Theta \in \mathcal{P}_{\text{sgn}}(r + 1)$  (the constant 1 is due to  $\text{rank}(\Theta - \pi) \leq r + 1$ ). Similar results hold for Tucker low-rank tensors  $\Theta \in \mathcal{P}_{\text{sgn}}(r + 1)$ , where  $r = \prod_k r_k$  with  $r_k$  being the  $k$ -th mode Tucker rank of  $\Theta$ .

**Example 2** (Tensor block models (TBMs)). Tensor block model (Wang and Zeng, 2019; Chi et al., 2020) assumes a checkerboard structure among tensor entries under marginal index permutation. The signal tensor  $\Theta$  takes at most  $r$  distinct values, where  $r$  is the total number of multiway blocks. Our model incorporates TBM because  $\Theta \in \mathcal{P}_{\text{sgn}}(r)$ .

**Example 3** (Generalized linear models (GLMs)). Let  $\mathcal{Y}$  be a binary tensor from a logistic model (Wang and Li, 2020) with mean  $\Theta = \text{logit}(\mathcal{Z})$ , where  $\mathcal{Z}$  is a latent low-rank tensor. Notice that

$\Theta$  itself may be high-rank (see Figure 1a). By definition,  $\Theta$  is a low-rank sign representable tensor. Same conclusion holds for general exponential-family models with a (known) link function (Hong et al., 2020).

**Example 4** (Single index models (SIMs)). Single index model is a flexible semiparametric model proposed in economics (Robinson, 1988) and high-dimensional statistics (Balabdaoui et al., 2019; Ganti et al., 2017). The SIM assumes the existence of a (unknown) monotonic function  $g: \mathbb{R} \rightarrow \mathbb{R}$  such that  $g(\Theta)$  has rank  $r$ . We see that  $\Theta$  belongs to the sign representable family; i.e.,  $\Theta \in \mathcal{P}_{\text{sgn}}(r+1)$ .

**Example 5** (Structured tensors with repeating entries). Here we revisit the model introduced in Figure 1b of Section 1. Let  $\Theta$  be an order- $K$  tensor with entries  $\Theta(i_1, \dots, i_K) = \log(1 + \max_k x_{i_k}^{(k)})$ , where  $x_{i_k}^{(k)}$  are given numbers in  $[0, 1]$  for all  $i_k \in [d_k], k \in [K]$ . We conclude that  $\Theta \in \mathcal{P}_{\text{sgn}}(2)$ , because the sign tensor  $\text{sgn}(\Theta - \pi)$  with an arbitrary  $\pi \in (0, \log 2)$  is a block tensor with at most two blocks (see Figure 2c). Similar results extend to structured tensors with entries  $\Theta(i_1, \dots, i_K) = g(\max_k x_{i_k}^{(k)})$ , where  $g(\cdot)$  is a polynomial of degree  $r$ . In this case,  $\Theta$  is a high-rank tensor with at most  $d_{\max}$  distinct entries but we have  $\Theta \in \mathcal{P}_{\text{sgn}}(2r)$  (see proofs in Section 7.2).

### 3.2 Statistical characterization of sign tensors via weighted classification

We now provide the explicit form of the weighted loss introduced in (3), and show that sign tensors are characterized by weighted classification. The results bridge the algebraic and statistical properties of sign representable tensors.

For a given  $\pi \in [-1, 1]$ , define a  $\pi$ -shifted data tensor  $\bar{\mathcal{Y}}_\Omega$  with entries  $\bar{\mathcal{Y}}(\omega) = (\mathcal{Y}(\omega) - \pi)$  for  $\omega \in \Omega$ . We propose a weighted classification objective function

$$L(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) = \frac{1}{|\Omega|} \sum_{\omega \in \Omega} \underbrace{|\bar{\mathcal{Y}}(\omega)|}_{\text{weight}} \times \underbrace{|\text{sgn} \mathcal{Z}(\omega) - \text{sgn} \bar{\mathcal{Y}}(\omega)|}_{\text{classification loss}}, \quad (4)$$

where  $\mathcal{Z} \in \mathbb{R}^{d_1 \times \dots \times d_K}$  is the decision variable to be optimized,  $|\bar{\mathcal{Y}}(\omega)|$  is the entry-specific weight equal to the distance from the tensor entry to the target level  $\pi$ . The entry-specific weights incorporate the magnitude information into classification, where entries far away from the target level are penalized more heavily in the objective. In the special case of binary tensor  $\mathcal{Y} \in \{-1, 1\}^{d_1 \times \dots \times d_K}$  and target level  $\pi = 0$ , the loss (4) reduces to usual classification loss.

Our proposed weighted classification function (4) is important for characterizing  $\text{sgn}(\Theta - \pi)$ . Define the weighted classification risk

$$\text{Risk}(\mathcal{Z}) = \mathbb{E}_{\mathcal{Y}_\Omega} L(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega), \quad (5)$$

where the expectation is taken with respect to  $\mathcal{Y}_\Omega$  under model (2) and the sampling distribution  $\omega \sim \Pi$ . The form of  $\text{Risk}(\cdot)$  implicitly depends on  $\pi$ ; we suppress  $\pi$  when no confusion arises.

**Proposition 2** (Global optimum of weighted risk). *Suppose the data  $\mathcal{Y}_\Omega$  is generated from model (2) with  $\Theta \in \mathcal{P}_{\text{sgn}}(r)$ . Then, for all  $\bar{\Theta}$  that are sign equivalent to  $\text{sgn}(\Theta - \pi)$ ,*

$$\text{Risk}(\bar{\Theta}) = \inf\{\text{Risk}(\mathcal{Z}): \mathcal{Z} \in \mathbb{R}^{d_1 \times \dots \times d_K}\} = \inf\{\text{Risk}(\mathcal{Z}): \text{rank}(\mathcal{Z}) \leq r\}.$$

The results show that the sign tensor  $\text{sgn}(\Theta - \pi)$  optimizes the weighted classification risk. This fact suggests a practical procedure to estimate  $\text{sgn}(\Theta - \pi)$  via empirical risk optimization of  $L(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega)$ . In order to establish the recovery guarantee, we shall address the uniqueness (up to sign equivalence) for the optimizer of  $\text{Risk}(\cdot)$ . The local behavior of  $\Theta$  around  $\pi$  plays a key role in the accuracy.

Some additional notation is needed for stating the results in full generality. Let  $d_{\text{total}} = \prod_{k=1}^K d_k$  denote the total number of tensor entries, and  $\Delta s = 1/d_{\text{total}}$  a small tolerance. We quantify the distribution of tensor entries  $\Theta(\omega)$  using a pseudo density, i.e., histogram with bin width  $2\Delta s$ . Let  $G(\pi) := \mathbb{P}_{\omega \sim \Pi}[\Theta(\omega) \leq \pi]$  denote the cumulative distribution function (CDF) of  $\Theta(\omega)$  under  $\omega \sim \Pi$ . We partition  $[-1, 1] = \mathcal{N}^c \cup \mathcal{N}$ , such that the pseudo density based on  $2\Delta$ -bin is uniformly bounded over  $\mathcal{N}^c$ ; i.e.,

$$\mathcal{N}^c = \left\{ \pi \in [-1, 1]: \frac{G(\pi + \Delta s) - G(\pi - \Delta s)}{\Delta s} \leq C \right\}, \text{ for some universal constant } C > 0,$$

and  $\mathcal{N}$  otherwise. Both  $\Theta$  and its induced CDF  $G$  implicitly depend on the tensor dimension.

**Assumption 1** ( $\alpha$ -smoothness). Fix  $\pi \in \mathcal{N}^c$ . Assume there exist constants  $\alpha = \alpha(\pi) > 0, c = c(\pi) > 0$ , independent of tensor dimension, such that,

$$\sup_{\Delta s \leq t < \rho(\pi, \mathcal{N})} \frac{G(\pi + t) - G(\pi - t)}{t^\alpha} \leq c, \quad (6)$$

where  $\rho(\pi, \mathcal{N}) := \min_{\pi' \in \mathcal{N}} |\pi - \pi'| + \Delta s$  denotes the adjusted distance from  $\pi$  to the nearest point in  $\mathcal{N}$ . The largest possible  $\alpha = \alpha(\pi)$  in (6) is called the smoothness index at level  $\pi$ . We make the convention that  $\alpha = \infty$  if the numerator in (6) is zero. A tensor  $\Theta$  is called  $\alpha$ -globally smooth, if (6) holds with global constants  $\alpha > 0, c > 0$  for all  $\pi \in \mathcal{N}^c$ .

The smoothness index  $\alpha$  quantifies the intrinsic hardness of recovering  $\text{sgn}(\Theta - \pi)$  from  $\text{Risk}(\cdot)$ . The value of  $\alpha$  depends on both the sampling distribution  $\omega \sim \Pi$  and the behavior of  $\Theta(\omega)$ . The recovery is easier at levels where points are less concentrated around  $\pi$  with a large value of  $\alpha > 1$ , or equivalently, when  $G(\pi)$  remains almost flat around  $\pi$ . A small value of  $\alpha < 1$  indicates the nonexistent (infinite) density at level  $\pi$ , or equivalently, when the  $G(\pi)$  jumps by greater than the tolerance  $\Delta s$  at  $\pi$ . Table 1 illustrates the  $G(\pi)$  for various models of  $\Theta$  (see Section 5).

We now reach the main theorem in this section. For two tensors  $\Theta_1, \Theta_2$ , define the mean absolute error (MAE) as  $\text{MAE}(\Theta_1, \Theta_2) = \mathbb{E}_{\omega \sim \Pi} |\Theta_1(\omega) - \Theta_2(\omega)|$ .

**Theorem 1** (Identifiability). Assume  $\Theta \in \mathcal{P}_{\text{sgn}}(r)$  is  $\alpha$ -globally smooth. Then, for all  $\pi \in \mathcal{N}^c$  and tensors  $\bar{\Theta} \simeq \text{sgn}(\Theta - \pi)$ , we have

$$\text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) \lesssim C(\pi) [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})]^{\alpha/(\alpha+1)} + \Delta s, \quad \text{for all } \mathcal{Z} \in \mathbb{R}^{d_1 \times \dots \times d_K},$$

where  $C(\pi) > 0$  is independent of  $\mathcal{Z}$ .

The result establishes the recovery stability of sign tensors  $\text{sgn}(\Theta - \pi)$  using optimization with population risk (5). The bound immediately shows the uniqueness of the optimizer for  $\text{Risk}(\cdot)$  up to a  $\Delta s$ -measure set under  $\Pi$ . We find that a higher value of  $\alpha$  implies more stable recovery, as intuition would suggest. Similar results hold for optimization with sample risk (4) (see Section 4).

We conclude this section by applying Assumption 1 to the examples described in Section 3.1. For simplicity, suppose  $\Pi$  is the uniform sampling. The tensor block model is  $\infty$ -globally smooth. This is because the set  $\mathcal{N}$  consists of finite  $2\Delta s$ -bin's covering the distinct block means in  $\Theta$ . Furthermore, we have  $\alpha = \infty$  for all  $\pi \in \mathcal{N}^c$ , since the numerator in (6) is zero. Similarly, the high-rank  $(d, d, d)$ -dimensional tensor  $\Theta(i, j, k) = \log(1 + \frac{1}{d} \max(i, j, k))$  is  $\infty$ -globally smooth because  $\alpha = \infty$  for all  $\pi$  except those in  $\mathcal{N}$ , where  $\mathcal{N}$  collects  $d$  many  $2\Delta s$ -bin's covering  $\log(1 + i/d)$  for all  $i \in [d]$ .

## 4 Nonparametric tensor completion via sign series

In previous sections we have established the sign series representation and its relationship to classification. In this section, we present our learning reduction proposal in details (Figure 2). We provide the estimation error bound and address the empirical implementation of the method.

### 4.1 Statistical error and sample complexity

Given a noisy incomplete tensor observation  $\mathcal{Y}_\Omega$  from model (2), we cast the problem of estimating  $\Theta$  into a series of weighted classifications. Specifically, we propose the signal tensor estimate using averaged structured sign tensors

$$\hat{\Theta} = \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{sgn} \hat{\mathcal{Z}}_\pi, \quad \text{with} \quad \hat{\mathcal{Z}}_\pi = \arg \min_{\mathcal{Z}: \text{rank} \mathcal{Z} \leq r} L(\mathcal{Z}, \mathcal{Y}_\Omega - \pi), \quad (7)$$

where  $\mathcal{H} = \{-1, \dots, -1/H, 0, 1/H, \dots, 1\}$  is the series of levels to aggregate,  $L(\cdot, \cdot)$  denotes the weighted classification objective defined in (4), and the rank constraint on  $\mathcal{Z}$  follows from Proposition 2. For the theory, we assume the true  $r$  is known; in practice,  $r$  could be chosen in a data adaptive fashion via cross-validation or elbow method (Hastie et al., 2009).

The next theorem establishes the statistical convergence for the sign tensor estimate (7).

**Theorem 2** (Sign tensor estimation). *Suppose  $\Theta \in \mathcal{P}_{\text{sgn}}(r)$  and  $\Theta(\omega)$  is  $\alpha$ -globally smooth under  $\omega \sim \Pi$ . Let  $\hat{\mathcal{Z}}_\pi$  be the estimate in (7),  $d = \max_{k \in [K]} d_k$ , and  $t_d = \frac{dr \log |\Omega|}{|\Omega|} \lesssim 1$ . Then, for all  $\pi \in \mathcal{N}^c$ , with very high probability over  $\mathcal{Y}_\Omega$ ,*

$$\text{MAE}(\text{sgn} \hat{\mathcal{Z}}_\pi, \text{sgn}(\Theta - \pi)) \lesssim t_d^{\alpha/(\alpha+2)} + \frac{1}{\rho^2(\pi, \mathcal{N})} t_d. \quad (8)$$

Theorem 2 provides the error bound for the sign tensor estimation. Compared to the population results in Theorem 1, we explicitly reveal the dependence of accuracy on the sample complexity and the level  $\pi$ . The result demonstrates the polynomial decay of sign errors with  $|\Omega|$ . Our sign estimate achieves consistent recovery using as few as  $\tilde{O}(dr)$  noisy entries.

Recall that  $\mathcal{N}$  collects the levels for which the sign tensor is possibly nonrecoverable. Let  $|\mathcal{N}|$  be the covering number of  $\mathcal{N}$  with  $2\Delta s$ -bin's, i.e.,  $|\mathcal{N}| = \lceil \text{Leb}(\mathcal{N})/2\Delta s \rceil$ , where  $\text{Leb}(\cdot)$  is the Lebesgue measure and  $\lceil \cdot \rceil$  is the ceiling function. Combining the sign representability of the signal tensor and the sign estimation accuracy, we obtain our main results on nonparametric tensor estimation.

**Theorem 3** (Tensor estimation error). *Consider the same conditions of Theorem 2. Let  $\hat{\Theta}$  be the estimate in (7). For any resolution parameter  $H \in \mathbb{N}_+$ , with very high probability over  $\mathcal{Y}_\Omega$ ,*

$$\text{MAE}(\hat{\Theta}, \Theta) \lesssim (t_d \log H)^{\frac{\alpha}{\alpha+2}} + \frac{1 + |\mathcal{N}|}{H} + t_d H \log H. \quad (9)$$

*In particular, setting  $H = (1 + |\mathcal{N}|)^{1/2} t_d^{-1/2} \asymp \text{poly}(d)$  yields the tightest upper bound in (52).*

Theorem 3 demonstrates the convergence rate of our tensor estimation. The bound (52) reveals three sources of errors: the estimation error for sign tensors, the bias from sign series representations, and the variance thereof. The resolution parameter  $H$  controls the bias-variance tradeoff. We remark that the signal estimation error (52) is generally no better than the corresponding sign error (8). This is to be expected, since magnitude estimation is a harder problem than sign estimation.



In the special case of full observation with equal dimension  $d_1 = \dots = d_K = d$  and bounded  $|\mathcal{N}| \leq C$ , our signal estimate achieves convergence

$$\text{MAE}(\hat{\Theta}, \Theta) \lesssim r d^{-(K-1) \min(\frac{\alpha}{\alpha+2}, \frac{1}{2})} \log^2 d,$$

by setting  $H \asymp d^{(K-1)/2}$ . Compared to earlier methods, our estimation accuracy applies to both low- and high-rank signal tensors. The rate depends on the sign complexity  $\Theta \in \mathcal{P}_{\text{sgn}}(r)$ , and this  $r$  is often much smaller than the usual tensor rank (see Section 3.1). Our result also reveals that the convergence becomes favorable as the order of data tensor increases.

We apply our method to the main examples in Section 3.1, and compare the results with existing literature. The numerical comparison is provided in Section 5.

**Example 2** (TBMs). Consider a tensor block model with  $r$  multiway blocks. Our result implies a rate  $\tilde{\mathcal{O}}(d^{-(K-1)/2})$  by taking  $\alpha = \infty$  and  $|\mathcal{N}| \leq r^K \lesssim \mathcal{O}(1)$ . This rate agrees with the previous root-mean-square error (RMSE) for block tensor estimation (Wang and Zeng, 2019).

**Example 3** (GLMs). Consider a GLM tensor  $\Theta = g(\mathcal{Z})$ , where  $g$  is a known link function and  $\mathcal{Z}$  is a latent low-rank tensor. Suppose the CDF of  $\Theta(\omega)$  is uniformly bounded as  $d \rightarrow \infty$ . Applying our results with  $\alpha = 1$  and finite  $|\mathcal{N}|$  yields  $\tilde{\mathcal{O}}(d^{-(K-1)/3})$ . This rate is slightly slower than the parametric RMSE rate (Zhang and Xia, 2018; Wang and Li, 2020), as expected. The reason is that our estimate remains valid for unknown  $g$  and general high-rank tensors. The nonparametric rate is the price one has to pay for not knowing the form  $\Theta = g(\mathcal{Z})$  as a priori.

**Example 4** (SIMs). The earlier example has shown the nonparametric rate  $\tilde{\mathcal{O}}(d^{-(K-1)/3})$  when applying our method to single index tensor model. In the matrix case with  $K = 2$ , our result yields error rate  $\tilde{\mathcal{O}}(d^{-1/3})$ , which is faster than the RMSE rate  $\mathcal{O}(d^{-1/4})$  obtained by Ganti et al. (2015).

**Example 5** (Structured tensors with repeating entries). We consider a more general model than that in Section 1. Consider a  $r$ -sign representable tensor  $\Theta \in \mathcal{P}_{\text{sgn}}(r)$  with at most  $d$  distinct entries with repetition pattern. Applying our results with  $\alpha = \infty$  and  $|\mathcal{N}| = d$  yields the rate  $\tilde{\mathcal{O}}(d^{-(K-2)/2})$ .

The following corollary reveal the sample complexity for nonparametric tensor completion.

**Corollary 1** (Sample complexity for nonparametric completion). *Assume the same conditions of Theorem 3 and bounded  $|\mathcal{N}|$ . Then, with high probability over  $\mathcal{Y}_\Omega$ ,*

$$\text{MAE}(\hat{\Theta}, \Theta) \rightarrow 0, \quad \text{as} \quad \frac{|\Omega|}{dr \log^2 |\Omega|} \rightarrow \infty.$$

Our result improves earlier work (Yuan and Zhang, 2016; Ghadermarzy et al., 2019; Lee and Wang, 2020) by allowing both low- and high-rank signals. Interestingly, the sample requirements depend only on the sign complexity  $dr$  but not the nonparametric complexity  $\alpha$ . Note that  $\tilde{\mathcal{O}}(dr)$  roughly matches the degree of freedom of sign tensors, suggesting the optimality of our sample requirements.

## 4.2 Implementation via learning reduction

This section addresses the practical implementation of our estimation (7). We take a learning reduction approach by dividing the full procedure into a meta algorithm and  $2H + 1$  base algorithms. The meta algorithm takes the average of  $(2H + 1) \asymp \text{poly}(d)$  sign tensors, whereas each base

---

**Algorithm 1** Nonparametric tensor completion via learning reduction

---

**Input:** Noisy and incomplete data tensor  $\mathcal{Y}_\Omega$ , rank  $r$ , resolution parameter  $H$ , ridge penalty  $\lambda$ .

- 1: **for**  $\pi \in \mathcal{H} = \{-1, \dots, -\frac{1}{H}, 0, \frac{1}{H}, \dots, 1\}$  **do**
- 2:     Define  $\pi$ -shifted tensor  $\bar{\mathcal{Y}} = \mathcal{Y} - \pi$  and corresponding sign tensor  $\text{sgn}(\bar{\mathcal{Y}}) = \text{sgn}(\mathcal{Y} - \pi)$ .
- 3:     Perform 1-bit tensor estimation algorithm (Ghadermarzy et al., 2018; Wang and Li, 2020; Hong et al., 2020; Alquier et al., 2019) on  $\bar{\mathcal{Y}}_\Omega$  and obtain

$$\hat{\mathcal{Z}}_\pi \leftarrow \arg \min_{\text{low-rank } \mathcal{Z}} \sum_{\omega \in \Omega} |\bar{\mathcal{Y}}(\omega)| F(\mathcal{Z}(\omega) \text{sgn} \bar{\mathcal{Y}}(\omega)) + \lambda \|\mathcal{Z}\|_F^2,$$

where  $F(\cdot)$  is the large-margin loss and  $\lambda$  is the penalty parameter.

4: **end for**

**Output:** Estimated signal tensor  $\hat{\Theta}_F = \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{sgn}(\hat{\mathcal{Z}}_\pi)$ .

---

algorithm estimates the tensor  $\text{sgn}(\Theta - \pi)$  given binary input  $\text{sgn}(\mathcal{Y} - \pi)$  and a target rank  $r$ . The full procedure is described in Algorithm 1 and Figure 2.

The base algorithm reduces to a low-rank 1-bit tensor estimation problem. Following the common practice in classification Bartlett et al. (2006), we replace the 0-1 loss  $\ell(z, y) = |\text{sgn} z - \text{sgn} y|$  in (4) with a continuous large-margin loss  $F(m)$  where  $m = z \text{sgn}(y)$  is the margin. Examples of large-margin loss are hinge loss  $F(m) = (1 - m)_+$ , logistic loss  $F(m) = \log(1 + e^{-m})$ , and  $\psi$ -loss  $F(m) = 2 \min(1, (1 - m)_+)$  with  $m_+ = \max(m, 0)$ . A number of polynomial-time algorithms with convergence guarantees are readily available for this problem (Ghadermarzy et al., 2018; Wang and Li, 2020; Hong et al., 2020; Alquier et al., 2019). We implement hinge loss (Alquier et al., 2019; Genzel and Stollenwerk, 2020; He et al., 2017) which maintains desirable statistical properties as in 0-1 loss. Here  $\text{Risk}_F(\cdot)$  is defined similarly as in (5) with hinge loss in place of 0-1 loss.

**Assumption 2** (Assumptions on surrogate loss).

- (a) (Approximation error) For any given  $\pi \in [-1, 1]$ , assume there exist a sequence of tensors  $\mathcal{Z}_\pi^{(n)} \in \mathcal{P}_{\text{sgn}}(r)$ , such that  $\text{Risk}_F(\mathcal{Z}_\pi^{(n)}) - \text{Risk}_F(\bar{\Theta}) \leq a_n$ , for some sequence  $a_n \rightarrow 0$  as  $n \rightarrow \infty$ . Furthermore, assume  $\|\mathcal{Z}_\pi^{(n)}\|_F \leq J$  for some constant  $J > 0$ .
- (b)  $F(z) = (1 - z)_+$  is hinge loss.

Assumption 2(a) quantifies the representation capability of and  $\mathcal{P}_{\text{sgn}}(r)$ . Assumption 2(b) implies the Fisher consistency bound for the weighted risk (Scott, 2011),

$$\text{Risk}(\mathcal{Z}) - \text{Risk}(\Theta - \pi) \lesssim \text{Risk}_F(\mathcal{Z}) - \text{Risk}_F(\Theta - \pi), \text{ for all } \pi \in [-1, 1] \text{ and all } \mathcal{Z}.$$

The Fisher consistency enable us to relate the excess risk of the large margin loss to that of 0-1 loss. Under Assumption 2, the resulting estimate from Algorithm 1 enjoys both statistical and computational efficiency.

**Theorem 4** (Large-margin estimation). *Consider the same setup as in Theorem 3, and denote  $t_n = \frac{d_{\max} r K \log n}{n}$ . Suppose the surrogate loss  $F$  satisfies Assumption 2 with  $a_n \lesssim t_n^{(\alpha+1)/(\alpha+2)}$ . Set  $\lambda \asymp t_n^{(\alpha+1)/(\alpha+2)} + t_n/\rho(\pi, \mathcal{N})$  in (46). Then, with high probability at least  $1 - \exp(-nt_n)$ , we have:*

- (a) (Sign tensor estimation). For all  $\pi \in [-1, 1]$  except for a finite number of levels,

$$\text{MAE}(\text{sgn}(\hat{\mathcal{Z}}_\pi), \text{sgn}(\Theta - \pi)) \lesssim t_n^{\frac{\alpha}{2+\alpha}} + \frac{1}{\rho^2(\pi, \mathcal{N})} t_n. \quad (10)$$

(b) (Tensor estimation).

$$\text{MAE}(\hat{\Theta}_F, \Theta) \lesssim (t_n \log H)^{\frac{\alpha}{2+\alpha}} + \frac{1 + |\mathcal{N}|}{H} + t_n H \log H. \quad (11)$$

In particular, setting  $H \asymp (1 + |\mathcal{N}|)^{1/2} t_n^{-1/2}$  yields the tightest upper bound in (11).

In principle, users can choose their own favorite large-margin losses, as long as the base algorithms are sample efficient. The comparison between various large-margin losses has been studied before Bartlett et al. (2006). Note that, instead of using  $\hat{\mathcal{Z}}_\pi$  as in existing 1-bit tensor algorithms Ghadermarzy et al. (2018); Wang and Li (2020), we use  $\text{sgn}(\hat{\mathcal{Z}}_\pi)$  for more challenging nonparametric estimation. The sign aggregation brings the benefits of flexibility and accuracy over classical low-rank models.

## 5 Simulations

In this section, we compare our nonparametric tensor method (**NonParaT**) with two alternative approaches: low-rank tensor CP decomposition (**CPT**), and the matrix version of our method applied to tensor unfolding (**NonParaM**). We assess the performance under both complete and incomplete observations. The signal tensors are generated based on four models listed in Table 1. The simulation covers a wide range of complexity, including block tensors, transformed low rank tensors, min/max hypergraphon with logarithm and exponential functions. We consider order-3 tensors of equal dimension  $d_1 = d_2 = d_3 = d$ , and set  $d \in \{15, 20, \dots, 55, 60\}$ ,  $r = 2$ ,  $H = 10 + (d - 15)/5$  in Algorithm 1. For **NonParaM**, we apply Algorithm 1 to each of the three unfolded matrices and report the average error. All summary statistics are averaged across 30 replicates.

Simulation	Signal Tensor $\Theta$	Rank	Sign Rank	$\alpha$	$ \mathcal{N} $	CDF	Noise
1	$\mathcal{C} \times \mathbf{M}_1 \times \mathbf{M}_2 \times \mathbf{M}_3$	$3^3$	$\leq 3^3$	$\infty$	$\leq 3^3$		Uniform $[-0.3, 0.3]$
2	$ \mathbf{a} \otimes \mathbf{1} \otimes \mathbf{1} - \mathbf{1} \otimes \mathbf{a} \otimes \mathbf{1} $	$d$	$\leq 3$	1	0		Normal $\mathcal{N}(0, 0.15)$
3	$\log(0.5 + Z_{\max})$	$\geq d$	2	$\infty$	$d$		Uniform $[-0.1, 0.1]$
4	$2.5 - \exp(\mathcal{Z}_{\min}^{1/3})$	$\geq d$	2	$\infty$	$d$		Normal $\mathcal{N}(0, 0.15)$

Table 1: Simulation models used for comparison. We use  $\mathbf{M}_k \in \{0, 1\}^{d \times 3}$  to denote membership matrices,  $\mathcal{C} \in \mathbb{R}^{3 \times 3 \times 3}$  the block means,  $\mathbf{a} = \frac{1}{d}(1, 2, \dots, d)^T \in \mathbb{R}^d$ ,  $\mathcal{Z}_{\max}$  and  $\mathcal{Z}_{\min}$  are order-3 tensors with entries  $\frac{1}{d} \max(i, j, k)$  and  $\frac{1}{d} \min(i, j, k)$ , respectively.

Figure 3 compares the estimation error under full observation. The MAE decreases with tensor dimension for all three methods. We find that our method **NonParaT** achieves the best performance in all scenarios, whereas the second best method is **CPT** for models 1-2, and **NonParaM** for models 3-4. One possible reason is that models 1-2 have controlled multilinear tensor rank, which makes tensor methods **NonParaT** and **CPT** more accurate than matrix methods. For models 3-4, the rank exceeds the tensor dimension, and therefore, the two nonparametric methods **NonParaT** and **NonparaM** exhibit the greater advantage for signal recovery.

Figure 4 shows the completion error against observation fraction. We fix  $d = 40$  and gradually increase the observation fraction  $\frac{|\Omega|}{d^3}$  from 0.3 to 1. We find that **NonParaT** achieves the lowest error among all methods. Our simulation covers a reasonable range of complexities; for example, model 1 has  $3^3$  jumps in the CDF of signal  $\Theta$ , and models 2 and 4 have unbounded noise. Nevertheless, our

method shows good performance in spite of model misspecification. This robustness is appealing in practice because the structure of underlying signal tensor is often unknown.

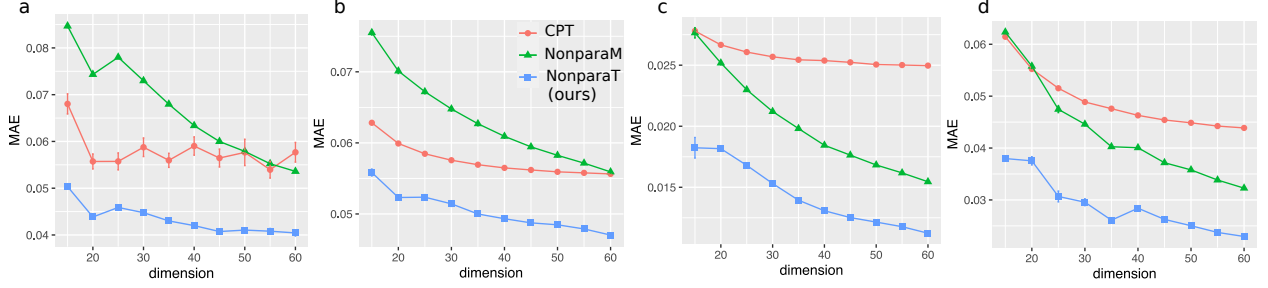


Figure 3: Estimation error versus tensor dimension. Panels (a)-(d) correspond to simulation models 1-4 in Table 1.

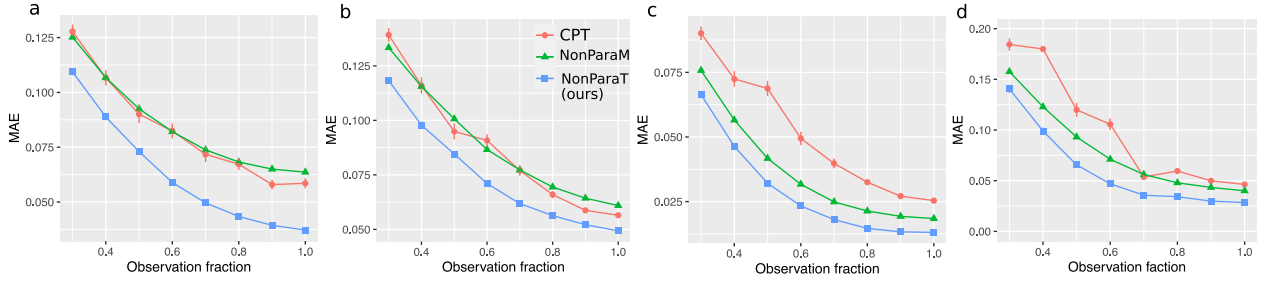


Figure 4: Completion error versus observation fraction. Panels (a)-(d) correspond to simulation models 1-4 in Table 1.

## 6 Data applications

We apply our method to two tensor datasets, the MRN-114 human brain connectivity data (Wang et al., 2017), and NIPS word occurrence data (Globerson et al., 2007).

### 6.1 Brain connectivity analysis

The brain dataset records the structural connectivity among 68 brain regions for 114 individuals along with their Intelligence Quotient (IQ) scores. We organize the connectivity data into an order-3 tensor, where entries encode the presence or absence of fiber connections between brain regions across individuals.

Figure 5 shows the MAE based on 5-fold cross-validations with  $r = 3, 6, \dots, 15$  and  $H = 20$ . We find that our method outperforms CPT in all combinations of ranks and missing rates. The achieved error reduction appears to be more profound as the missing rate increases. This trend highlights the applicability of our method in tensor completion tasks. In addition, our method exhibits a smaller standard error in cross-validation experiments as shown in Figure 5 and Table 2, demonstrating the stability over CPT. One possible reason is that that our estimate is guaranteed to be in  $[0, 1]$  (for binary tensor problem where  $\mathcal{Y} \in \{0, 1\}^{d_1 \times \dots \times d_K}$ ) whereas CPT estimation may fall outside the valid range  $[0, 1]$ .

We next investigate the pattern in the estimated signal tensor. Figure 6a shows the identified top edges associated with IQ scores. Specifically, we first obtain a denoised tensor  $\hat{\Theta} \in \mathbb{R}^{68 \times 68 \times 114}$  using

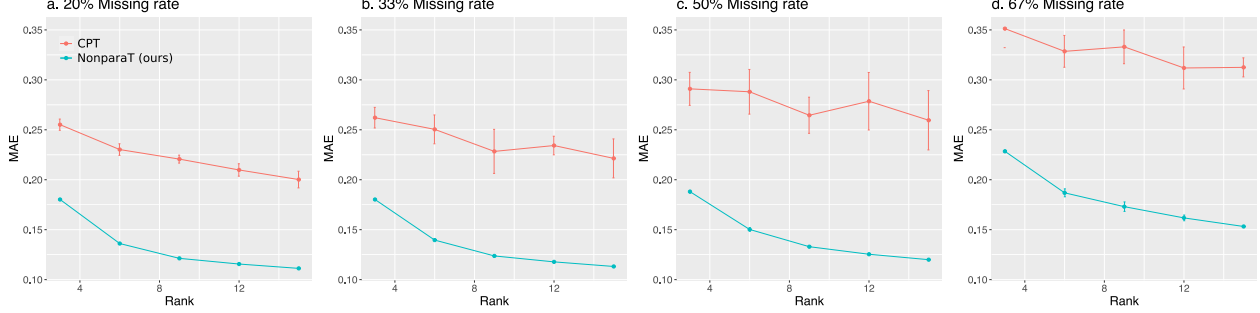


Figure 5: Estimation error versus rank under different missing rate. Panels (a)-(d) correspond to missing rate 20%, 33%, 50%, and 67%, respectively. Error bar represents the standard error over 5-fold cross-validations.

MRN-114 brain connectivity dataset					
Method	$r = 3$	$r = 6$	$r = 9$	$r = 12$	$r = 15$
NonparaT (Ours)	<b>0.18</b> (0.001)	<b>0.14</b> (0.001)	<b>0.12</b> (0.001)	<b>0.12</b> (0.001)	<b>0.11</b> (0.001)
Low-rank CPT	0.26(0.006)	0.23(0.006)	0.22(0.004)	0.21(0.006)	0.20(0.008)
NIPS word occurrence dataset					
Method	$r = 3$	$r = 6$	$r = 9$	$r = 12$	$r = 15$
NonparaT (Ours)	<b>0.18</b> (0.002)	<b>0.16</b> (0.002)	<b>0.15</b> (0.001)	<b>0.14</b> (0.001)	<b>0.13</b> (0.001)
Low-rank CPT	0.22(0.004)	0.20(0.007)	0.19(0.007)	0.17(0.007)	0.17(0.007)
Naive imputation (Baseline)	0.32(.001)				

Table 2: MAE comparison in the brain data and NIPS data analysis. Reported MAEs are averaged over five runs of cross-validation, with 20% entries for testing and 80% for training, with standard errors in parentheses. Bold numbers indicate the minimal MAE among three methods. For low-rank CPT, we use R function `rTensor` with default hyperparameters, and for our method, we set  $H = 20$ .

our method with  $r = 10$  and  $H = 20$ . Then, we perform a regression analysis of  $\hat{\Theta}(i, j, : ) \in \mathbb{R}^{144}$  against the normalized IQ score across the 144 individuals. The regression model is repeated for each edge  $(i, j) \in [68] \times [68]$ . We find that top edges represent the interhemispheric connections in the frontal lobes. The result is consistent with recent research on brain connectivity with intelligence (Li et al., 2009; Wang et al., 2017).

## 6.2 NIPS data analysis

The NIPS dataset consists of word occurrence counts in papers published from 1987 to 2003. We focus on the top 100 authors, 200 most frequent words, and normalize each word count by log transformation with pseudo-count 1. The resulting dataset is an order-3 tensor with entry representing the log counts of words by authors across years.

Table 2 compares the prediction accuracy of different methods. We find that our method substantially outperforms the low-rank CP method for every configuration under consideration. Further increment of rank appears to have little effect on the performance. The comparison highlights the advantage of our method in achieving accuracy while maintaining low complexity. In addition, we also perform naive imputation where the missing values are predicted using the sample average. Both our method and CPT outperform the naive imputation, implying the necessity of incorporating tensor structure in the analysis.

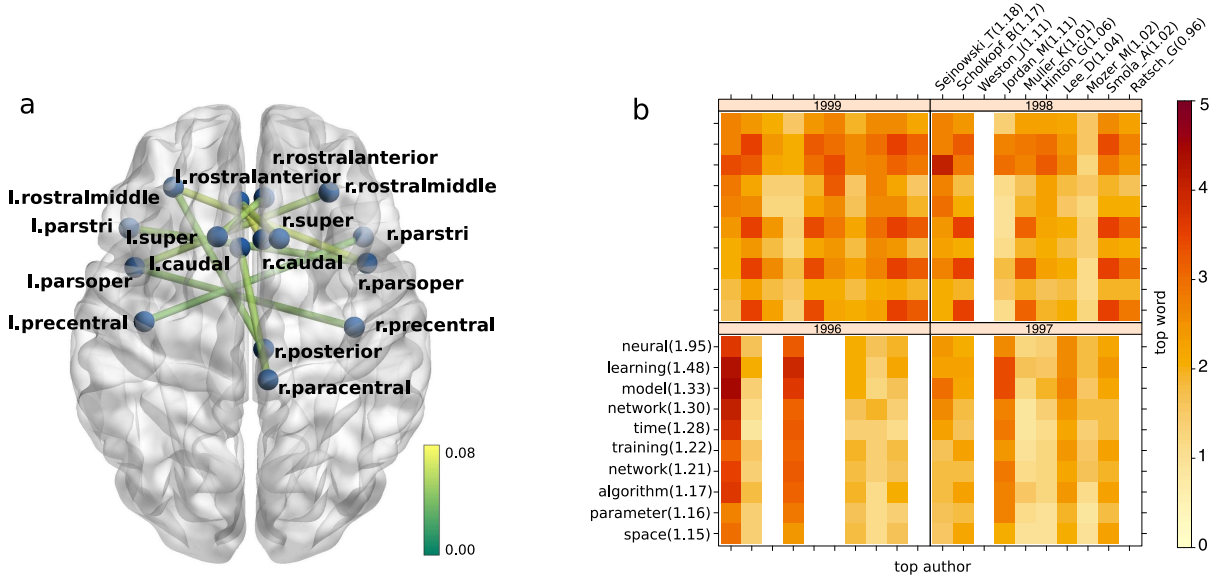


Figure 6: Estimated signal tensors in the data analysis. (a) top edges associated with IQ scores in the brain connectivity data. The color indicates the estimated IQ effect size. (b) top authors and words for years 1996-1999 in the NIPS data. Authors and words are ranked by marginal averages based on  $\hat{\Theta}$ , where the marginal average is denoted in the parentheses.

We next examine the estimated signal tensor  $\hat{\Theta}$  from our method. Figure 6b illustrates the results from NIPS data, where we plot the entries in  $\hat{\Theta}$  corresponding to top authors and most-frequent words (after excluding generic words such as *figure*, *results*, etc). The identified pattern is consistent with the active topics in the NIPS publication. Among the top words are *neural* (marginal mean = 1.95), *learning* (1.48), and *network* (1.21), whereas top authors are *T. Sejnowski* (1.18), *B. Scholkopf* (1.17), *M. Jordan* (1.11), and *G. Hinton* (1.06). We also find strong heterogeneity among word occurrences across authors and years. For example, *training* and *algorithm* are popular words for *B. Scholkopf* and *A. Smola* in 1998-1999, whereas *model* occurs more often in *M. Jordan* and in 1996. The detected pattern and achieved accuracy demonstrate the applicability of our method.

## 7 Additional results and proofs

In this section, we provide additional results not covered in previous sections. Section 7.1 gives detailed explanation to the examples mentioned in Section 1. Section 7.2 supplements Section 3.1 by providing more theoretical results on sign rank and its relationship to tensor rank. Section 7.3 collects the proofs for theorems in the main texts. Lastly, Section 7.4 extends all results to unbounded observations with sub-Gaussian noise.

### 7.1 Sensitivity of tensor rank to monotonic transformations

In Section 1 of the main paper, we have provided a motivating example to show the sensitivity of tensor rank to monotonic transformations. Here, we describe the details of the example set-up.

The step 1 is to generate a rank-3 tensor  $\mathcal{Z}$  based on the CP representation

$$\mathcal{Z} = \mathbf{a}^{\otimes 3} + \mathbf{b}^{\otimes 3} + \mathbf{c}^{\otimes 3},$$

where  $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathbb{R}^{30}$  are vectors consisting of  $N(0, 1)$  entries, and the shorthand  $\mathbf{a}^{\otimes 3} = \mathbf{a} \otimes \mathbf{a} \otimes \mathbf{a}$  denotes the Kronecker power. We then apply  $f(z) = (1 + \exp(-cz))^{-1}$  to  $\mathcal{Z}$  entrywise, and obtain a transformed tensor  $\Theta = f(\mathcal{Z})$ .

The step 2 is to determine the rank of  $\Theta$ . Unlike matrices, the exact rank determination for tensors is NP hard. Therefore, we choose to compute the numerical rank of  $\Theta$  as an approximation. The numerical rank is determined as the minimal rank for which the relative approximation error is below 0.1, i.e.,

$$\hat{r}(\Theta) = \min \left\{ s \in \mathbb{N}_+ : \min_{\hat{\Theta} : \text{rank}(\hat{\Theta}) \leq s} \frac{\|\Theta - \hat{\Theta}\|_F}{\|\Theta\|_F} \leq 0.1 \right\}.$$

We compute  $\hat{r}(\Theta)$  by searching over  $s \in \{1, \dots, 30^2\}$ , where for each  $s$ , we (approximately) solve the least-square minimization using CP function in R package **rTensor**. We repeat steps 1-2 ten times, and plot the averaged numerical rank of  $\Theta$  versus transformation level  $c$  in Figure 1a.

## 7.2 Tensor rank and sign-rank

In the main paper, we have provided several tensor examples with high tensor rank but low sign-rank. This section provides more examples and their proofs. Unless otherwise specified, let  $\Theta$  be an order- $K$   $(d, \dots, d)$ -dimensional tensor.

**Example 6** (Structured tensors with repeating entries). Suppose the tensor  $\Theta$  takes the form

$$\Theta(i_1, \dots, i_K) = \log \left( 1 + \frac{1}{d} \max(i_1, \dots, i_K) \right), \text{ for all } (i_1, \dots, i_K) \in [d]^K.$$

Then

$$\text{rank}(\Theta) \geq d, \quad \text{and} \quad \text{srnk}(\Theta - \pi) \leq 2 \text{ for all } \pi \in \mathbb{R}.$$

*Proof of Example 6.* We first prove the results for  $K = 2$ . The full-rankness of  $\Theta$  is verified from elementary row operations as follows

$$\begin{pmatrix} (\Theta_2 - \Theta_1)/(\log(1 + \frac{2}{d}) - \log(1 + \frac{1}{d})) \\ (\Theta_3 - \Theta_2)/(\log(1 + \frac{3}{d}) - \log(1 + \frac{2}{d})) \\ \vdots \\ (\Theta_d - \Theta_{d-1})/(\log(1 + \frac{d}{d}) - \log(1 + \frac{d-1}{d})) \\ \Theta_d / \log(1 + \frac{d}{d}) \end{pmatrix} = \begin{pmatrix} 1 & 0 & \ddots & \ddots & 0 \\ 1 & 1 & \ddots & \ddots & \ddots \\ \vdots & \vdots & \ddots & \ddots & \ddots \\ 1 & 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 \end{pmatrix},$$

where  $\Theta_i$  denotes the  $i$ -th row of  $\Theta$ . Now it suffices to show  $\text{srnk}(\Theta - \pi) \leq 2$  for  $\pi$  in the feasible range  $(\log(1 + \frac{1}{d}), \log 2)$ . In this case, there exists an index  $i^* \in \{2, \dots, d\}$ , such that  $\log(1 + \frac{i^*-1}{d}) < \pi \leq \log(1 + \frac{i^*}{d})$ . By definition, the sign matrix  $\text{sgn}(\Theta - \pi)$  takes the form

$$\text{sgn}(\Theta(i, j) - \pi) = \begin{cases} -1, & \text{both } i \text{ and } j \text{ are smaller than } i^*; \\ 1, & \text{otherwise.} \end{cases} \quad (12)$$

Therefore, the matrix  $\text{sgn}(\Theta - \pi)$  is a rank-2 block matrix, which implies  $\text{srnk}(\Theta - \pi) = 2$ .

We now extend the results to  $K \geq 3$ . By definition of the tensor rank, the rank of a tensor is lower bounded by the rank of its matrix slice. So we have  $\text{rank}(\Theta) \geq \text{rank}(\Theta(:, :, 1, \dots, 1)) = d$ . For the

sign rank with feasible  $\pi$ , notice that the sign tensor  $\text{sgn}(\Theta - \pi)$  takes the similar form as in (12),

$$\text{sgn}(\Theta(i_1, \dots, i_K) - \pi) = \begin{cases} -1, & i_k < i^* \text{ for all } k \in [K]; \\ 1, & \text{otherwise,} \end{cases} \quad (13)$$

where  $i^*$  denotes the index that satisfies  $\log(1 + \frac{i^*-1}{d}) < \pi \leq \log(1 + \frac{i^*}{d})$ . The equation (13) implies that  $\text{sgn}(\Theta - \pi) = -2\mathbf{a}^{\otimes K} + 1$ , where  $\mathbf{a} = (1, \dots, 1, 0, \dots, 0)^T$  takes 1 on the  $i$ -th entry if  $i < i^*$  and 0 otherwise. Henceforth  $\text{srnk}(\Theta - \pi) = 2$ .  $\square$

In fact, Example 6 is a special case of the following proposition.

**Proposition 3** (Structured tensors with repeating entries). *Let  $g: \mathbb{R} \rightarrow \mathbb{R}$  be a continuous function such that  $g(z) = 0$  has at most  $r \geq 1$  distinct real roots. For given numbers  $x_{i_k}^{(k)} \in [0, 1]$  for all  $i_k \in [d_k]$ , define a tensor  $\Theta \in \mathbb{R}^{d_1 \times \dots \times d_K}$  with entries*

$$\Theta(i_1, \dots, i_K) = g(\max(x_{i_1}^{(1)}, \dots, x_{i_K}^{(K)})), \quad (i_1, \dots, i_K) \in [d_1] \times \dots \times [d_K]. \quad (14)$$

*Then, the sign rank of  $(\Theta - \pi)$  satisfies*

$$\text{srnk}(\Theta - \pi) \leq 2r.$$

*The same conclusion holds if we use min in place of max in (14).*

*Proof of Proposition 3.* We reorder the tensor indices along each mode such that  $x_1^{(k)} \leq \dots \leq x_{d_k}^{(k)}$  for all  $k \in [K]$ . Based on the construction of  $\mathcal{Z}_{\max}$ , the reordering does not change the rank of  $\mathcal{Z}_{\max}$  or  $(\Theta - \pi)$ . Let  $z_1 < \dots < z_r$  be the  $r$  distinct real roots for the equation  $g(z) = \pi$ . We separate the proof for two cases,  $r = 1$  and  $r \geq 2$ .

- When  $r = 1$ . The continuity of  $g(\cdot)$  implies that the function  $(g(z) - \pi)$  has at most one sign change point. Using similar proof as in Example 6, we have

$$\text{sgn}(\Theta - \pi) = 1 - 2\mathbf{a}^{(1)} \otimes \dots \otimes \mathbf{a}^{(K)} \quad \text{or} \quad \text{sgn}(\Theta - \pi) = 2\mathbf{a}^{(1)} \otimes \dots \otimes \mathbf{a}^{(K)} - 1,$$

where  $\mathbf{a}^{(k)}$  are binary vectors defined by

$$\mathbf{a}^{(k)} = (\underbrace{1, \dots, 1}_{\text{positions for which } x_{i_k}^{(k)} < z_1}, 0, \dots, 0)^T, \quad \text{for } k \in [K].$$

Therefore,  $\text{srnk}(\Theta - \pi) \leq \text{rank}(\text{sgn}(\Theta - \pi)) = 2$ .

- When  $r \geq 2$ . By continuity, the function  $(g(z) - \pi)$  is non-zero and remains an unchanged sign in each of the intervals  $(z_s, z_{s+1})$  for  $1 \leq s \leq r - 1$ . Define the index set

$$\mathcal{I} = \{s \in \mathbb{N}_+ : \text{the interval } (z_s, z_{s+1}) \text{ in which } g(z) < \pi\}.$$

We now prove that the sign tensor  $\text{sgn}(\Theta - \pi)$  has rank bounded by  $2r - 1$ . To see this, consider the tensor indices for which  $\text{sgn}(\Theta - \pi) = -1$ ,

$$\begin{aligned} \{\omega : \Theta(\omega) - \pi < 0\} &= \{\omega : g(\mathcal{Z}_{\max}(\omega)) < \pi\} \\ &= \cup_{s \in \mathcal{I}} \{\omega : \mathcal{Z}_{\max}(\omega) \in (z_s, z_{s+1})\} \end{aligned}$$



$$= \cup_{s \in \mathcal{I}} \left( \{ \omega : x_{i_k}^{(k)} < z_{s+1} \text{ for all } k \in [K] \} \cap \{ \omega : x_{i_k}^{(k)} \leq z_s \text{ for all } k \in [K] \}^c \right). \quad (15)$$

The equation (15) is equivalent to

$$\mathbb{1}(\Theta(i_1, \dots, i_K) < \pi) = \sum_{s \in \mathcal{I}} \left( \prod_k \mathbb{1}(x_{i_k}^{(k)} < z_{s+1}) - \prod_k \mathbb{1}(x_{i_k}^{(k)} \leq z_s) \right), \quad (16)$$

for all  $(i_1, \dots, i_K) \in [d_1] \times \dots \times [d_K]$ , where  $\mathbb{1}(\cdot) \in \{0, 1\}$  denotes the indicator function. The equation (16) implies the low-rank representation of  $\text{sgn}(\Theta - \pi)$ ,

$$\text{sgn}(\Theta - \pi) = 1 - 2 \sum_{s \in \mathcal{I}} \left( \mathbf{a}_{s+1}^{(1)} \otimes \dots \otimes \mathbf{a}_{s+1}^{(K)} - \bar{\mathbf{a}}_s^{(1)} \otimes \dots \otimes \bar{\mathbf{a}}_s^{(K)} \right), \quad (17)$$

where  $\mathbf{a}_{s+1}^{(k)}, \bar{\mathbf{a}}_s^{(k)}$  are binary vectors defined by

$$\mathbf{a}_{s+1}^{(k)} = (\underbrace{1, \dots, 1}_{\text{positions for which } x_{i_k}^{(k)} < z_{s+1}}, 0, \dots, 0)^T, \quad \text{and} \quad \bar{\mathbf{a}}_s^{(k)} = (\underbrace{1, \dots, 1}_{\text{positions for which } x_{i_k}^{(k)} \leq z_s}, 0, \dots, 0)^T.$$

Therefore, by (17) and the assumption  $|\mathcal{I}| \leq r - 1$ , we conclude that

$$\text{srnk}(\Theta - \pi) \leq 1 + 2(r - 1) = 2r - 1.$$

Combining two cases yields that  $\text{srnk}(\Theta - \pi) \leq 2r$  for any  $r \geq 1$ .  $\square$

We next provide several additional examples such that  $\text{rank}(\Theta) \geq d$  whereas  $\text{srnk}(\Theta) \leq c$  for a constant  $c$  independent of  $d$ . We state the examples in the matrix case, i.e,  $K = 2$ . Similar conclusion extends to  $K \geq 3$ , by the following proposition.

**Proposition 4** (Rank relationship between matrices and tensors). *Let  $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$  be a matrix. For any given  $K \geq 3$ , define an order- $K$  tensor  $\Theta \in \mathbb{R}^{d_1 \times \dots \times d_K}$  by*

$$\Theta = \mathbf{M} \otimes \mathbf{1}_{d_3} \otimes \dots \otimes \mathbf{1}_{d_K},$$

where  $\mathbf{1}_{d_k} \in \mathbb{R}^{d_k}$  denotes an all-one vector, for  $3 \leq k \leq K$ . Then we have

$$\text{rank}(\Theta) = \text{rank}(\mathbf{M}), \quad \text{and} \quad \text{srnk}(\Theta - \pi) = \text{srnk}(\mathbf{M} - \pi) \text{ for all } \pi \in \mathbb{R}.$$

*Proof of Proposition 4.* The conclusion directly follows from the definition of tensor rank.  $\square$

**Example 7** (Stacked banded matrices). Let  $\mathbf{a} = (1, 2, \dots, d)^T$  be a  $d$ -dimensional vector, and define a  $d$ -by- $d$  banded matrix  $\mathbf{M} = |\mathbf{a} \otimes \mathbf{1} - \mathbf{1} \otimes \mathbf{a}|$ . Then

$$\text{rank}(\mathbf{M}) = d, \quad \text{and} \quad \text{srnk}(\mathbf{M} - \pi) \leq 3, \quad \text{for all } \pi \in \mathbb{R}.$$

*Proof of Example 7.* Note that  $\mathbf{M}$  is a banded matrix with entries

$$\mathbf{M}(i, j) = |i - j|, \quad \text{for all } (i, j) \in [d]^2.$$

Elementary row operation shows that  $\mathbf{M}$  is full rank as follows,

$$\begin{pmatrix} (\mathbf{M}_1 + \mathbf{M}_d)/(d-1) \\ \mathbf{M}_1 - \mathbf{M}_2 \\ \mathbf{M}_2 - \mathbf{M}_3 \\ \vdots \\ \mathbf{M}_{d-1} - \mathbf{M}_d \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & \cdots & 1 & 1 \\ -1 & 1 & 1 & \cdots & 1 & 1 \\ -1 & -1 & 1 & \cdots & 1 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ -1 & -1 & -1 & \cdots & -1 & 1 \end{pmatrix}.$$

We now show  $\text{srnk}(\mathbf{M} - \pi) \leq 3$  by construction. Define two vectors  $\mathbf{b} = (2^{-1}, 2^{-2}, \dots, 2^{-d})^T \in \mathbb{R}^d$  and  $\text{rev}(\mathbf{b}) = (2^{-d}, \dots, 2^{-1})^T \in \mathbb{R}^d$ . We construct the following matrix

$$\mathbf{A} = \mathbf{b} \otimes \text{rev}(\mathbf{b}) + \text{rev}(\mathbf{b}) \otimes \mathbf{b}. \quad (18)$$

The matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is banded with entries

$$\mathbf{A}(i, j) = \mathbf{A}(j, i) = \mathbf{A}(d-i, d-j) = \mathbf{A}(d-j, d-i) = 2^{-d-1} (2^{j-i} + 2^{i-j}), \text{ for all } (i, j) \in [d]^2.$$

Furthermore, the entry value  $\mathbf{A}(i, j)$  decreases with respect to  $|i - j|$ ; i.e.,

$$\mathbf{A}(i, j) \geq \mathbf{A}(i', j'), \quad \text{for all } |i - j| \geq |i' - j'|. \quad (19)$$

Notice that for a given  $\pi \in \mathbb{R}$ , there exists  $\pi' \in \mathbb{R}$  such that  $\text{sgn}(\mathbf{A} - \pi') = \text{sgn}(\mathbf{M} - \pi)$ . This is because both  $\mathbf{A}$  and  $\mathbf{M}$  are banded matrices satisfying monotonicity (19). By definition (18),  $\mathbf{A}$  is a rank-2 matrix. Henceforce,  $\text{srnk}(\mathbf{M} - \pi) = \text{srnk}(\mathbf{A} - \pi') \leq 3$ .  $\square$

**Remark 1.** The tensor analogy of banded matrices  $\Theta = |\mathbf{a} \otimes \mathbf{1} \otimes \mathbf{1} - \mathbf{1} \otimes \mathbf{a} \otimes \mathbf{1}|$  is used as simulation model 3 in the main paper.

**Example 8** (Stacked identity matrices). Let  $\mathbf{I}$  be a  $d$ -by- $d$  identity matrix. Then

$$\text{rank}(\mathbf{I}) = d, \quad \text{and} \quad \text{srnk}(\mathbf{I} - \pi) \leq 3 \text{ for all } \pi \in \mathbb{R}.$$

*Proof of Proposition 8.* Depending on the value of  $\pi$ , the sign matrix  $\text{sgn}(\mathbf{I} - \pi)$  falls into one of the two cases:

- (a)  $\text{sgn}(\mathbf{I} - \pi)$  is a matrix of all 1, or of all  $-1$ ;
- (b)  $\text{sgn}(\mathbf{I} - \pi) = 2\mathbf{I} - \mathbf{1}_d \otimes \mathbf{1}_d$ .

The first cases are trivial, so it suffices to show  $\text{srnk}(\mathbf{I} - \pi) \leq 3$  in the third case.

Based on Example 7, the rank-2 matrix  $\mathbf{A}$  in (18) satisfies

$$\mathbf{A}(i, j) \begin{cases} = 2^{-d}, & i = j, \\ \geq 2^{-d} + 2^{-d-2}, & i \neq j. \end{cases}$$

Therefore,  $\text{sgn}(2^{-d} + 2^{-d-3} - \mathbf{A}) = 2\mathbf{I} - \mathbf{1}_d \otimes \mathbf{1}_d$ . We conclude that  $\text{srnk}(\mathbf{I} - \pi) \leq \text{rank}(2^{-d} + 2^{-d-3} - \mathbf{A}) = 3$ .  $\square$

## 7.3 Proofs

### 7.3.1 Proofs of Propositions 1-2

*Proof of Proposition 1.*

Part (a). The strictly monotonicity of  $g$  implies that the inverse function  $g^{-1}: \mathbb{R} \rightarrow \mathbb{R}$  is well-defined. When  $g$  is strictly increasing, the mapping  $x \mapsto g(x)$  is sign preserving. Specifically, if  $x \geq 0$ , then  $g(x) \geq g(0) = 0$ . Conversely, if  $g(x) \geq 0 = g(0)$ , then applying  $g^{-1}$  to both sides gives  $x \geq 0$ . When  $g$  is strictly decreasing, the mapping  $x \mapsto g(x)$  is sign reversing. Specifically, if  $x \geq 0$ , then  $g(x) \leq g(0) = 0$ . Conversely, if  $g(x) \leq 0 = g(0)$ , then applying  $g^{-1}$  to both sides gives  $x \leq 0$ . Therefore,  $\Theta \simeq g(\Theta)$ , or  $\Theta \simeq -g(\Theta)$ . Since constant multiplication does not change the tensor rank, we have  $\text{srnk}(\Theta) = \text{srnk}(g(\Theta)) \leq \text{rank}(g(\Theta))$ .

Part (b). See Section 7.2 for constructive examples. □

*Proof of Proposition 2.* Fix  $\pi \in [-1, 1]$ . Based on the definition of classification loss  $L(\cdot, \cdot)$ , the function  $\text{Risk}(\cdot)$  relies only on the sign pattern of the tensor. Therefore, without loss of generality, we assume both  $\bar{\Theta}, \mathcal{Z} \in \{-1, 1\}^{d_1 \times \dots \times d_K}$  are binary tensors. We evaluate the excess risk

$$\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) = \mathbb{E}_{\omega \sim \Pi} \underbrace{\mathbb{E}_{\mathcal{Y}(\omega)} \{ |\mathcal{Y}(\omega) - \pi| [ |\mathcal{Z}(\omega) - \text{sgn}(\bar{\mathcal{Y}}(\omega))| - |\bar{\Theta}(\omega) - \text{sgn}(\bar{\mathcal{Y}}(\omega))| ] \}}_{\stackrel{\text{def}}{=} I(\omega)}. \quad (20)$$

Denote  $y = \mathcal{Y}(\omega)$ ,  $z = \mathcal{Z}(\omega)$ ,  $\bar{\theta} = \bar{\Theta}(\omega)$ , and  $\theta = \Theta(\omega)$ . The expression of  $I(\omega)$  is simplified as

$$\begin{aligned} I(\omega) &= \mathbb{E}_{y|\omega} [(y - \pi)(\bar{\theta} - z)\mathbf{1}(y \geq \pi) + (\pi - y)(z - \bar{\theta})\mathbf{1}(y < \pi)] \\ &= \mathbb{E}_{y|\omega} [(\bar{\theta} - z)(y - \pi)] \\ &= [\text{sgn}(\theta - \pi) - z](\theta - \pi) \\ &= |\text{sgn}(\theta - \pi) - z||\theta - \pi| \geq 0, \end{aligned} \quad (21)$$

where the third line uses the fact  $\mathbb{E}y = \theta$  and  $\bar{\theta} = \text{sgn}(\theta - \pi)$ , and the last line uses the assumption  $z \in \{-1, 1\}$ . The equality (21) is attained when  $z = \text{sgn}(\theta - \pi)$  or  $\theta = \pi$ . Combining (21) with (20), we conclude that, for all  $\mathcal{Z} \in \{-1, 1\}^{d_1 \times \dots \times d_K}$ ,

$$\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) = \mathbb{E}_{\omega \sim \Pi} |\text{sgn}(\Theta(\omega) - \pi) - \mathcal{Z}(\omega)||\Theta(\omega) - \pi| \geq 0. \quad (22)$$

In particular, setting  $\mathcal{Z} = \bar{\Theta} = \text{sgn}(\Theta - \pi)$  in (22) yields the minimum. Therefore,

$$\text{Risk}(\bar{\Theta}) = \min\{\text{Risk}(\mathcal{Z}): \mathcal{Z} \in \mathbb{R}^{d_1 \times \dots \times d_K}\} \leq \min\{\text{Risk}(\mathcal{Z}): \text{rank}(\mathcal{Z}) \leq r\}.$$

Since  $\text{srnk}(\Theta - \pi) \leq r$  by assumption, the last inequality becomes equality. The proof is complete. □

### 7.3.2 Proof of Theorem 1

*Proof of Theorem 1.* Fix  $\pi \notin \mathcal{N}$ . Based on (22) in Proposition 2, we have

$$\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) = \mathbb{E} [|\text{sgn}\mathcal{Z} - \text{sgn}\bar{\Theta}||\bar{\Theta}|]. \quad (23)$$

The Assumption 1 states that

$$\mathbb{P}(|\bar{\Theta}| \leq t) \leq \begin{cases} ct^\alpha, & \text{for all } \Delta s \leq t < \rho(\pi, \mathcal{N}), \\ C\Delta s, & \text{for all } 0 \leq t < \Delta s. \end{cases} \quad (24)$$

Without further specification, all relevant probability statements, such as  $\mathbb{E}$  and  $\mathbb{P}$ , are with respect to  $\omega \sim \Pi$ .

We divide the proof into two cases:  $\alpha > 0$  and  $\alpha = \infty$ .

- Case 1:  $\alpha > 0$ .

By (23), for all  $0 \leq t < \rho(\pi, \mathcal{N})$ ,

$$\begin{aligned} \text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) &\geq t\mathbb{E}(|\text{sgn}\mathcal{Z} - \text{sgn}\bar{\Theta}|\mathbb{1}\{|\bar{\Theta}| > t\}) \\ &\geq 2t\mathbb{P}(\text{sgn}\mathcal{Z} \neq \text{sgn}\bar{\Theta} \text{ and } |\bar{\Theta}| > t) \\ &\geq 2t\left\{\mathbb{P}(\text{sgn}\mathcal{Z} \neq \text{sgn}\bar{\Theta}) - \mathbb{P}(|\bar{\Theta}| \leq t)\right\} \\ &\geq t\left\{\text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) - C\Delta s - 2ct^\alpha\right\}, \end{aligned} \quad (25)$$

where the last line follows from the definition of MAE and (24). We maximize the lower bound (25) with respect to  $t$ , and obtain the optimal  $t_{\text{opt}}$ ,

$$t_{\text{opt}} = \begin{cases} \rho(\pi, \mathcal{N}), & \text{if } \text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) > \text{cut-off}, \\ \left[\frac{1}{2c(1+\alpha)}(\text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) - C\Delta s)\right]^{1/\alpha}, & \text{if } \text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) \leq \text{cut-off}. \end{cases}$$

where we have denoted the cut-off  $= 2c(1+\alpha)\rho^\alpha(\pi, \mathcal{N}) + C\Delta s$ . The corresponding lower bound of the inequality (25) becomes

$$\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) \geq \begin{cases} c_1\rho(\pi, \mathcal{N}) [\text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) - C\Delta s], & \text{if } \text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) > \text{cut-off}, \\ c_2 [\text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) - C\Delta s]^{\frac{1+\alpha}{\alpha}}, & \text{if } \text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) \leq \text{cut-off}, \end{cases}$$

where  $c_1, c_2 > 0$  are two constants independent of  $\mathcal{Z}$ . Combining both cases gives

$$\begin{aligned} \text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) &\lesssim [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})] + \Delta s \\ &\leq C(\pi)[\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})]^{\frac{\alpha}{1+\alpha}} + \Delta s, \end{aligned}$$

where  $C(\pi) > 0$  is a multiplicative factor independent of  $\mathcal{Z}$ .

- Case 2:  $\alpha = \infty$ . The inequality (25) now becomes

$$\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) \geq t [\text{MAE}(\text{sgn}\bar{\Theta}, \text{sgn}\mathcal{Z}) - C\Delta s], \quad \text{for all } 0 \leq t < \rho(\pi, \mathcal{N}). \quad (26)$$

The conclusion follows by taking  $t = \frac{\rho(\pi, \mathcal{N})}{2}$  in the inequality (26).

□

**Remark 2.** The proof of Theorem 1 shows that, under global  $\alpha$ -smoothness of  $\Theta$ ,

$$\text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) \lesssim [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})] + \Delta s, \quad (27)$$

for all  $\mathcal{Z} \in \mathbb{R}^{d_1 \times \dots \times d_K}$ . For fixed  $\pi$ , the second term is absorbed into the first term.

### 7.3.3 Proof of Theorem 2

The following lemma provides the variance-to-mean relationship implied by the  $\alpha$ -smoothness of  $\Theta$ . The relationship plays a key role in determining the convergence rate based on empirical process theory (Shen and Wong, 1994); also see Theorem 5.

**Lemma 1** (Variance-to-mean relationship). *Consider the same setup as in Theorem 2. Fix  $\pi \notin \mathcal{N}$ . Let  $L(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega)$  be the  $\pi$ -weighted classification loss*

$$\begin{aligned} L(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) &= \frac{1}{|\Omega|} \sum_{\omega \in \Omega} \underbrace{|\bar{\mathcal{Y}}(\omega)|}_{\text{weight}} \times \underbrace{|\text{sgn}\mathcal{Z}(\omega) - \text{sgn}\bar{\mathcal{Y}}(\omega)|}_{\text{classification loss}} \\ &= \frac{1}{|\Omega|} \sum_{\omega \in \Omega} \ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}), \end{aligned} \quad (28)$$

where we have denoted the function  $\ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}) \stackrel{\text{def}}{=} |\bar{\mathcal{Y}}(\omega)| |\text{sgn}\mathcal{Z}(\omega) - \text{sgn}\bar{\mathcal{Y}}(\omega)|$ . Under Assumption 1 of the  $\alpha$ -smoothness of  $\Theta$ , we have

$$\text{Var}[\ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) - \ell_\omega(\bar{\Theta}, \bar{\mathcal{Y}}_\Omega)] \lesssim [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})] + \Delta s, \quad (29)$$

for all tensors  $\mathcal{Z} \in \mathbb{R}^{d_1 \times \dots \times d_K}$ . Here the expectation and variance are taken with respect to both  $\mathcal{Y}$  and  $\omega \sim \Pi$ .

*Proof of Lemma 1.* We expand the variance by

$$\begin{aligned} \text{Var}[\ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) - \ell_\omega(\bar{\Theta}, \bar{\mathcal{Y}}_\Omega)] &\lesssim \mathbb{E}|\ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) - \ell_\omega(\bar{\Theta}, \bar{\mathcal{Y}}_\Omega)|^2 \\ &\lesssim \mathbb{E}|\ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) - \ell_\omega(\bar{\Theta}, \bar{\mathcal{Y}}_\Omega)| \\ &\leq \mathbb{E}|\text{sgn}\mathcal{Z} - \text{sgn}\bar{\Theta}| = \text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}), \end{aligned} \quad (30)$$

where the second line comes from the boundedness of classification loss  $L(\cdot, \cdot)$ , and the third line comes from the inequality  $||a - b| - |c - b|| \leq |a - b|$  for  $a, b, c \in \{-1, 1\}$ , together with the boundedness of classification weight  $|\bar{\mathcal{Y}}(\omega)|$ . Here we have absorbed the constant multipliers in  $\lesssim$ . The conclusion (29) then directly follows by applying Remark 2 to (30).  $\square$

*Proof of Theorem 2.* Fix  $\pi \notin \mathcal{N}$ . For notational simplicity, we suppress the subscript  $\pi$  and write  $\hat{\mathcal{Z}}$  in place of  $\hat{\mathcal{Z}}_\pi$ . Denote  $n = |\Omega|$  and  $\rho = \rho(\pi, \mathcal{N})$ .

Because the classification loss  $L(\cdot, \cdot)$  is scale-free, i.e.,  $L(\mathcal{Z}, \cdot) = L(c\mathcal{Z}, \cdot)$  for every  $c > 0$ , we consider the estimation subject to  $\|\mathcal{Z}\|_F \leq 1$  without loss of generality. Specifically, let

$$\hat{\mathcal{Z}} = \arg \min_{\mathcal{Z}: \text{rank}(\mathcal{Z}) \leq r, \|\mathcal{Z}\|_F \leq 1} L(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega). \quad (31)$$

We next apply the empirical process theory to bound  $\hat{\mathcal{Z}}$ . To facilitate the analysis, we view the data  $\bar{\mathcal{Y}}_\Omega = \{\bar{\mathcal{Y}}(\omega): \omega \in \Omega\}$  as a collection of  $n$  independent random variables where the randomness is from both  $\bar{\mathcal{Y}}$  and  $\omega \sim \Pi$ . Write the index set  $\Omega = \{1, \dots, n\}$ , so the loss function (28) becomes

$$L(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) = \frac{1}{n} \sum_{i=1}^n \ell_i(\mathcal{Z}, \bar{\mathcal{Y}}).$$

We use  $f_{\mathcal{Z}}: [d_1] \times \cdots \times [d_n] \rightarrow \mathbb{R}$  to denote the function induced by tensor  $\mathcal{Z}$  such that  $f_{\mathcal{Z}}(\omega) = \mathcal{Z}(\omega)$  for  $\omega \in [d_1] \times \cdots \times [d_K]$ . Under this set-up, the quantity of interest

$$L(\mathcal{Z}, \bar{\mathcal{Y}}_{\Omega}) - L(\bar{\Theta}, \bar{\mathcal{Y}}_{\Omega}) = \frac{1}{n} \sum_{i=1}^n \underbrace{[\ell_i(\mathcal{Z}, \bar{\mathcal{Y}}) - \ell_i(\bar{\Theta}, \bar{\mathcal{Y}})]}_{\stackrel{\text{def}}{=} \Delta_i(f_{\mathcal{Z}}, \bar{\Theta})},$$

is an empirical process induced by function  $f_{\mathcal{Z}} \in \mathcal{F}_{\mathcal{T}}$  where  $\mathcal{T} = \{\mathcal{Z}: \text{rank}(\mathcal{Z}) \leq r, \|\mathcal{Z}\|_F \leq 1\}$ . Note that there is an one-to-one correspondence between sets  $\mathcal{F}_{\mathcal{T}}$  and  $\mathcal{T}$ .

Let  $L_n$  denote the desired convergence rate to seek. By definition of  $\hat{\mathcal{Z}}$  in (31), we have,

$$L(\hat{\mathcal{Z}}, \bar{\mathcal{Y}}_{\Omega}) - L(\bar{\Theta}, \bar{\mathcal{Y}}_{\Omega}) = \frac{1}{n} \sum_{i=1}^n \Delta_i(f_{\hat{\mathcal{Z}}}, \bar{\Theta}) \leq 0.$$

Therefore, we have the following inclusion of probability events,

$$\begin{aligned} & \left\{ (\omega, \mathcal{Y}_{\omega}): \text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta}) \geq L_n \right\} \\ & \subset \left\{ (\omega, \mathcal{Y}_{\omega}): \exists \mathcal{Z} \text{ s.t. } \text{rank}(\mathcal{Z}) \leq r, \text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) \geq L_n, \text{ and } \frac{1}{n} \sum_{i=1}^n \Delta_i(f_{\mathcal{Z}}, \bar{\Theta}) \leq 0 \right\} \\ & \subset \left\{ (\omega, \mathcal{Y}_{\omega}): \sup_{\substack{\text{rank}(\mathcal{Z}) \leq r \\ \text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) \geq L_n}} -\frac{1}{n} \sum_{i=1}^n \Delta_i(f_{\mathcal{Z}}, \bar{\Theta}) \geq 0 \right\} \\ & \subset \bigcup_{\ell=1}^{\infty} \left\{ (\omega, \mathcal{Y}_{\omega}): \sup_{\mathcal{Z} \in A_{\ell}} -\frac{1}{n} \sum_{i=1}^n \Delta_i(f_{\mathcal{Z}}, \bar{\Theta}) \geq 0 \right\}, \end{aligned} \quad (32)$$

where we have partitioned  $\{\mathcal{Z}: \text{rank}(\mathcal{Z}) \leq r \text{ and } \text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) \geq L_n\}$  in to union of  $A_{\ell}$  with

$$A_{\ell} = \{\mathcal{Z}: \text{rank}(\mathcal{Z}) \leq r \text{ and } \ell L_n \leq \text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) < (\ell + 1)L_n\},$$

for  $\ell = 1, 2, \dots$ . Let  $\Gamma$  denote the target probability for the first line in (32). To bound  $\Gamma$ , we bound the sum of probability over the sets  $A_{\ell}$ . For each  $A_{\ell}$ , we consider the centered empirical process,

$$v_n(f_{\mathcal{Z}}) := -\frac{1}{n} \sum_{i=1}^n (\Delta_i(f_{\mathcal{Z}}, \bar{\Theta}) - \mathbb{E} \Delta_i(f_{\mathcal{Z}}, \bar{\Theta})). \quad (33)$$

Notice  $(\ell + 1)L_n \geq \mathbb{E} \Delta_i(f_{\mathcal{Z}}, \bar{\Theta}) = \text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) \geq \ell L_n$  for all  $\mathcal{Z} \in A_{\ell}$ . Combining (32), (33) and union bound yields

$$\Gamma \leq \sum_{\ell=1}^{\infty} \mathbb{P} \left\{ \sup_{\mathcal{Z} \in A_{\ell}} v_n(f_{\mathcal{Z}}) \geq \ell L_n =: M(\ell) \right\}. \quad (34)$$

Notice that, based on Lemma 1, the variance of empirical process is bounded by

$$\begin{aligned} \sup_{\mathcal{Z} \in A_{\ell}} \text{Var} \Delta_i(f_{\mathcal{Z}}, \bar{\Theta}) & \lesssim \sup_{\mathcal{Z} \in A_{\ell}} \left( [\mathbb{E} \Delta_i(f_{\mathcal{Z}}, \bar{\Theta})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho} \mathbb{E} \Delta_i(f_{\mathcal{Z}}, \bar{\Theta}) \right) + \Delta s \\ & \leq M(\ell + 1)^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho} M(\ell + 1) + \Delta s =: V(\ell). \end{aligned}$$

We next bound the right-hand side of (34) by choosing  $L_n$  that satisfies conditions in Theorem 5 (The specification of  $L_n$  is deferred to the next paragraph). One such  $L_n$  is chosen, Theorem 5 gives us

$$\begin{aligned}\Gamma &\lesssim \sum_{\ell=1}^{\infty} \exp\left(-\frac{nM^2(\ell)}{V(\ell) + 2M(\ell)}\right) \\ &\lesssim \sum_{\ell=1}^{\infty} \exp(-\rho \ell n L_n) \\ &\leq \left(\frac{e^{-n\rho L_n}}{1 - e^{-n\rho L_n}}\right).\end{aligned}\tag{35}$$

Now, we specify  $L_n$  that satisfies the condition of Theorem 5. The quantity  $L_n$  is determined by the solution to the following inequality,

$$\sup_{\ell \geq 1} \frac{1}{x} \int_x^{\sqrt{x^{\alpha/(\alpha+1)} + x/\rho + \Delta s}} \sqrt{\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}_{\mathcal{T}}, \|\cdot\|_2)} d\varepsilon \lesssim n^{1/2}, \quad \text{where } x = \ell L_n.\tag{36}$$

In particular, the smallest  $L_n$  satisfying (36) yields the best upper bound of the error rate. Here  $\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}_{\mathcal{T}}, \|\cdot\|_2)$  denotes the  $L_2$ -norm,  $\varepsilon$ -bracketing number (c.f. Definition 3) for function family  $\mathcal{F}_{\mathcal{T}}$ .

Based on Lemma 2, the inequality (36) is satisfied with the choice

$$L_n \asymp t_n^{(\alpha+1)/(\alpha+2)} + \frac{t_n}{\rho}, \quad \text{where } t_n = \left(\frac{d_{\max} r K \log n}{n}\right) \text{ and } d_{\max} := \max_{k \in [K]} d_k.$$

Finally, it follows from Theorem 5 and (35) that

$$\begin{aligned}\mathbb{P}\left\{\text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta}) \geq L_n\right\} &\lesssim \left(\frac{e^{-n\rho L_n}}{1 - e^{-n\rho L_n}}\right) \\ &\lesssim e^{-nt_n},\end{aligned}$$

where the last inequality uses the fact that  $\rho L_n \gtrsim t_n \gtrsim \frac{1}{n}$  by our choice of  $L_n$  and  $t_n$ .

Inserting the above bound into (27) gives that, with high probability at least  $1 - \exp(-nt_n)$ ,

$$\begin{aligned}\text{MAE}(\text{sgn}\hat{\mathcal{Z}}, \text{sgn}\bar{\Theta}) &\lesssim [\text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta})]^{\alpha/(\alpha+1)} + \frac{1}{\rho}[\text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta})] + \Delta s \\ &\lesssim t_n^{\alpha/(\alpha+2)} + \frac{1}{\rho^{\alpha/(\alpha+1)}} t_n^{\alpha/(\alpha+1)} + \frac{1}{\rho} t_n^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho^2} t_n \\ &\leq 4t_n^{\alpha/(\alpha+2)} + \frac{4}{\rho^2} t_n,\end{aligned}\tag{37}$$

where the second line uses the fact that  $\Delta s \ll t_n$ , and the last line follows from the fact that  $a(b^2 + b^{(\alpha+2)/(\alpha+1)} + b + 1) \leq 4a(b^2 + 1)$  with  $a = \frac{t_n}{\rho^2}$  and  $b = \rho t_n^{-1/(\alpha+2)}$ . We plug  $t_n$  into (37) and absorb the term  $K$  into the constant. The conclusion is then proved by noting  $n = |\Omega|$  by definition.  $\square$

**Definition 3** (Bracketing number). Consider a family of functions  $\mathcal{F}$ , and let  $\varepsilon > 0$ . Let  $\mathcal{X}$  denote the domain space equipped with measure  $\Pi$ . We call  $\{(f_m^l, f_m^u)\}_{m=1}^M$  an  $L_2$ -metric,  $\varepsilon$ -bracketing function set of  $\mathcal{F}$ , if for every  $f \in \mathcal{F}$ , there exists an  $m \in [M]$  such that

$$f_m^l(x) \leq f(x) \leq f_m^u(x), \quad \text{for all } x \in \mathcal{X},$$

and

$$\|f_m^l - f_m^u\|_2 \stackrel{\text{def}}{=} \sqrt{\mathbb{E}_{x \sim \Pi} |f_m^l(x) - f_m^u(x)|^2} \leq \varepsilon, \quad \text{for all } m = 1, \dots, M.$$

The bracketing number with  $L_2$ -metric, denoted  $\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}, \|\cdot\|_2)$ , is the logarithm of the smallest cardinality of the  $\varepsilon$ -bracketing function set of  $\mathcal{F}$ .

**Lemma 2** (Bracketing complexity of low-rank tensors). *Define the family of rank- $r$  bounded tensors  $\mathcal{T} = \{\mathcal{Z} \in \mathbb{R}^{d_1 \times \dots \times d_K} : \text{rank}(\mathcal{Z}) \leq r, \|\mathcal{Z}\|_F \leq 1\}$  and the induced function family  $\mathcal{F}_{\mathcal{T}} = \{f_{\mathcal{Z}} : \mathcal{Z} \in \mathcal{T}\}$ . Set*

$$L_n \asymp \left( \frac{d_{\max} r K \log n}{n} \right)^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho(\pi, \mathcal{N})} \left( \frac{d_{\max} r K \log n}{n} \right), \quad \text{where } d_{\max} = \max_{k \in [K]} d_k.$$

Then, the following inequality is satisfied provided that  $\Delta s \lesssim n^{-1}$ ,

$$\sup_{\ell \geq 1} \frac{1}{\ell L_n} \int_{\ell L_n}^{\sqrt{\ell L_n^{\alpha/(\alpha+1)} + \frac{\ell L_n}{\rho(\pi, \mathcal{N})} + \Delta s}} \sqrt{\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}_{\mathcal{T}}, \|\cdot\|_2)} d\varepsilon \leq C n^{1/2}, \quad (38)$$

where  $C > 0$  is a constant independent of  $r, K$  and  $d_{\max}$ .

*Proof of Lemma 2.* To simplify the notation, we denote  $\rho = \rho(\pi, \mathcal{N})$ . Notice that

$$\|f_{\mathcal{Z}_1} - f_{\mathcal{Z}_2}\|_2 \leq \|f_{\mathcal{Z}_1} - f_{\mathcal{Z}_2}\|_{\infty} \leq \|\mathcal{Z}_1 - \mathcal{Z}_2\|_F \quad \text{for all } \mathcal{Z}_1, \mathcal{Z}_2 \in \mathcal{T}.$$

It follows from Kosorok (2007, Theorem 9.22) that the  $L_2$ -metric,  $(2\varepsilon)$ -bracketing number of  $\mathcal{F}_{\mathcal{T}}$  is bounded by

$$\mathcal{H}_{[\cdot]}(2\varepsilon, \mathcal{F}_{\mathcal{T}}, \|\cdot\|_2) \leq \mathcal{H}(\varepsilon, \mathcal{T}, \|\cdot\|_F) \leq C d_{\max} r K \log \frac{K}{\varepsilon}.$$

The last inequality is from the covering number bounds for rank- $r$  bounded tensors; see Mu et al. (2014, Lemma 3). Inserting the bracketing number into (38) gives

$$g(L, \ell) = \frac{1}{\ell L} \int_{\ell L}^{\sqrt{\ell L^{\alpha/(\alpha+1)} + \rho^{-1} \ell L + \Delta s}} \sqrt{d_{\max} r K \log \left( \frac{K}{\varepsilon} \right)} d\varepsilon. \quad (39)$$

Define  $g(L) := \sup_{\ell \geq 1} g(L, \ell)$ . By the monotonicity the integrand in (39), we bound  $g(L)$  by

$$\begin{aligned} g(L) &\leq \sup_{\ell \geq 1} \frac{\sqrt{d_{\max} r K}}{\ell L} \int_{\ell L}^{\sqrt{\ell L^{\alpha/(\alpha+1)} + \rho^{-1} \ell L + n^{-1}}} \sqrt{\log \left( \frac{K}{\varepsilon} \right)} d\varepsilon \\ &\leq \sup_{\ell \geq 1} \sqrt{d_{\max} r K \log \left( \frac{K}{\ell L} \right)} \left( \frac{(\ell L)^{\alpha/(2\alpha+2)} + \sqrt{\rho^{-1} \ell L + n^{-1}}}{\ell L} - 1 \right) \\ &\lesssim \sqrt{d_{\max} r K \log(1/L)} \left[ \frac{1}{L^{(\alpha+2)/(2\alpha+2)}} + \frac{1}{\sqrt{\rho L}} \left( 1 + \frac{\rho}{2nL} \right) \right], \end{aligned} \quad (40)$$



where the second line follows from  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for  $a, b > 0$  and the last line comes from the fact that the bound achieves maximum when  $\ell = 1$ . It remains to verify that  $g(L_n) \leq Cn^{1/2}$  for  $L_n$  specified in (38). Plugging  $L_n$  into the last line of (40) gives

$$\begin{aligned} g(L_n) &\leq \sqrt{d_{\max} r K \log(1/L_n)} \left( \frac{1}{L_n^{(\alpha+2)/(2\alpha+2)}} + \frac{2}{\sqrt{\rho L_n}} \right) \\ &\leq \sqrt{d_{\max} r K \log n} \left( \left[ \left( \frac{d_{\max} r K \log n}{n} \right)^{\frac{\alpha+1}{\alpha+2}} \right]^{-\frac{\alpha+2}{2\alpha+2}} + \left[ 2\rho \left( \frac{d_{\max} r K \log n}{\rho n} \right) \right]^{-\frac{1}{2}} \right) \\ &\leq Cn^{1/2}, \end{aligned}$$

where  $C > 0$  is a constant independent of  $r, K$  and  $d_{\max}$ . The proof is therefore complete.  $\square$

**Theorem 5** (Theorem 3 in Shen and Wong (1994)). *Let  $\mathcal{F}$  be a class of functions defined on  $\mathcal{X}$  with  $\sup_{f \in \mathcal{F}} \|f\|_{\infty} \leq T$ . Let  $(\mathbf{X}_i)_{i=1}^n$  be i.i.d. random variables with distribution  $\mathbb{P}_{\mathbf{X}}$  over  $\mathcal{X}$ . Set  $\sup_{f \in \mathcal{F}} \text{Var} f(\mathbf{X}) = V < \infty$ . Define the empirical process  $\hat{\mathbb{E}}f = \frac{1}{n} \sum_{i=1}^n f(\mathbf{X}_i)$ . Define  $x_n^*$  to be the solution to the following inequality*

$$\frac{1}{x} \int_x^{\sqrt{V}} \sqrt{\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}, \|\cdot\|_2)} d\varepsilon \lesssim \sqrt{n}.$$

Suppose  $\sqrt{V} \leq T$  and

$$x_n^* \lesssim \frac{V}{T}, \quad \text{and} \quad \mathcal{H}_{[\cdot]}(\sqrt{V}, \mathcal{F}, \|\cdot\|_2) \lesssim \frac{n(x_n^*)^2}{V}.$$

Then, we have

$$\mathbb{P} \left( \sup_{f \in \mathcal{F}} \hat{\mathbb{E}}f - \mathbb{E}f \geq x_n^* \right) \lesssim \exp \left( -\frac{n(x_n^*)^2}{V + T x_n^*} \right).$$

### 7.3.4 Proof of Theorem 3

*Proof of Theorem 3.* By definition of  $\hat{\Theta}$ , we have

$$\begin{aligned} \text{MAE}(\hat{\Theta}, \Theta) &= \mathbb{E} \left| \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{sgn} \hat{Z}_{\pi} - \Theta \right| \\ &\leq \mathbb{E} \left| \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \left( \text{sgn} \hat{Z}_{\pi} - \text{sgn}(\Theta - \pi) \right) \right| + \mathbb{E} \left| \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{sgn}(\Theta - \pi) - \Theta \right| \\ &\leq \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{MAE}(\text{sgn} \hat{Z}_{\pi}, \text{sgn}(\Theta - \pi)) + \frac{1}{H}, \end{aligned} \tag{41}$$

where the last line comes from the triangle inequality and the inequality

$$\left| \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{sgn}(\Theta(\omega) - \pi) - \Theta(\omega) \right| \leq \frac{1}{H}, \quad \text{for all } \omega \in [d_1] \times \cdots \times [d_K].$$

Write  $n = |\Omega|$ . Now it suffices to bound the first term in (41). For any given  $t \geq t_n = \frac{d_{\max} r K \log n}{n}$ , define the event

$$A = \left\{ \text{MAE}(\text{sgn} \hat{Z}_{\pi}, \text{sgn}(\Theta - \pi)) \lesssim t^{\alpha/(2+\alpha)} + \frac{t}{\rho^2(\pi, \mathcal{N})} \text{ for all } \pi \in \mathcal{H} \right\}.$$

We shall prove that under the event  $A$ ,

$$\frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)) \lesssim t^{\alpha/(\alpha+2)} + \frac{1+|\mathcal{N}|}{H} + Ht. \quad (42)$$

Theorem 2 implies that the sign estimation accuracy depends on the closeness of  $\pi \in \mathcal{H}$  to the mass points in  $\mathcal{N}$ . Therefore, we partition the level set  $\pi \in \mathcal{H}$  based on their closeness to  $\mathcal{N}$ . Specifically, Define  $\mathcal{H}_1 \stackrel{\text{def}}{=} \{\pi \in \mathcal{H} : \rho(\pi, \mathcal{N}) < \frac{1}{H}\}$  and  $\mathcal{H}_2 = \mathcal{H} \setminus \mathcal{H}_1$ . Notice  $|\mathcal{H}_1| \leq 2|\mathcal{N}|$ . We expand the left hand side of (42) by

$$\begin{aligned} & \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)) \\ &= \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}_1} \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)) + \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}_2} \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)). \end{aligned} \quad (43)$$

The first term involves only  $2|\mathcal{N}|$  many number of summands thus can be bounded by  $4|\mathcal{N}|/(2H+1)$ . We bound the second term using the explicit forms of  $\rho(\pi, \mathcal{N})$  in the sequence  $\pi \in \mathcal{H}_2$ . Under the event  $A$ , we have

$$\begin{aligned} \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}_2} \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)) &\lesssim \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}_2} t^{\alpha/(\alpha+2)} + \frac{t}{2H+1} \sum_{\pi \in \mathcal{H}_2} \frac{1}{\rho^2(\pi, \mathcal{N})} \\ &\leq t^{\alpha/(\alpha+2)} + \frac{t}{2H+1} \sum_{\pi \in \mathcal{H}_2} \sum_{\pi' \in \mathcal{N}} \frac{1}{|\pi - \pi'|^2} \\ &\leq t^{\alpha/(\alpha+2)} + \frac{t}{2H+1} \sum_{\pi' \in \mathcal{N}} \sum_{\pi \in \mathcal{H}_2} \frac{1}{|\pi - \pi'|^2} \\ &\leq t^{\alpha/(\alpha+2)} + 2CHt, \end{aligned}$$

where the first inequality uses the property of event  $A$ , and the last inequality follows from Lemma 3. Combining the bounds for the two terms in (43) completes the proof for conclusion (42); that is

$$\mathbb{P} \left( \text{MAE}(\hat{\Theta}, \Theta) \lesssim t^{\alpha/(\alpha+2)} + \frac{1+|\mathcal{N}|}{H} + Ht \right) \geq \mathbb{P}(A). \quad (44)$$

Based on the proof of Theorem 2 and union bound over  $\pi \in \mathcal{H}$ , we have, for all  $t \geq t_n$ ,

$$\begin{aligned} \mathbb{P}(A) &\geq 1 - \sum_{\pi \in \mathcal{H}} \mathbb{P} \left( \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)) \gtrsim t^{\alpha/(\alpha+2)} + \frac{t}{\rho(\pi, \mathcal{N})^2} \right) \\ &\gtrsim 1 - (2H+1) \exp(-nt) \gtrsim 1 - \exp(-nt + \log H). \end{aligned} \quad (45)$$

We choose  $t \asymp t_n \log H$  in (45) so that  $\log H$  is negligible compared to  $nt$ . Finally, combining (44) and (45) with the choice of  $t$  yields

$$\text{MAE}(\hat{\Theta}, \Theta) \lesssim \left( \frac{d_{\max} r K \log |\Omega| \log H}{|\Omega|} \right)^{\alpha/(\alpha+2)} + \frac{1+|\mathcal{N}|}{H} + \frac{d_{\max} r K \log |\Omega|}{|\Omega|} H \log H,$$

with at least probability  $1 - \exp(-d_{\max} r K \log |\Omega| \log H) \geq 1 - \exp(-d_{\max} r K \log |\Omega|)$ .

□

**Lemma 3.** Fix  $\pi' \in \mathcal{N}$  and a sequence  $\Pi = \{-1, \dots, -1/H, 0, 1/H, \dots, 1\}$  with  $H \geq 2$ . Then,

$$\sum_{\pi \in \mathcal{H}_2} \frac{1}{|\pi - \pi'|^2} \leq 4H^2.$$

*Proof of Lemma 3.* Notice that all points  $\pi \in \mathcal{H}_2$  satisfy  $|\pi - \pi'| \gtrsim \frac{1}{H}$  for all  $\pi' \in \mathcal{N}$  by definition and the fact that  $\Delta s$  is negligible compared to  $1/H$ . We use this fact to compute the sum

$$\begin{aligned} \sum_{\pi \in \mathcal{H}_2} \frac{1}{|\pi - \pi'|^2} &= \sum_{\frac{h}{H} \in \mathcal{H}_2} \frac{1}{|\frac{h}{H} - \pi'|^2} \\ &\leq 2H^2 \sum_{h=1}^H \frac{1}{h^2} \\ &\leq 2H^2 \left\{ 1 + \int_1^2 \frac{1}{x^2} dx + \int_2^3 \frac{1}{x^2} dx + \dots + \int_{H-1}^H \frac{1}{x^2} dx \right\} \\ &= 2H^2 \left( 1 + \int_1^H \frac{1}{x^2} dx \right) \leq 4H^2, \end{aligned}$$

where the third line uses the monotonicity of  $\frac{1}{x^2}$  for  $x \geq 1$ .  $\square$

### 7.3.5 Proof of Theorem 4

*Proof of Theorem 4.* Write  $\bar{\mathcal{Y}} = \mathcal{Y} - \pi$ ,  $\bar{\Theta} = \Theta - \pi$ , and  $n = |\Omega|$ . Here we consider the estimation

$$\hat{\mathcal{Z}}_\pi = \arg \min_{\text{rank}(\mathcal{Z}) \leq r} \sum_{\omega \in \Omega} |\bar{\mathcal{Y}}(\omega)| \times F(\mathcal{Z}(\omega) \text{sgn}(\bar{\mathcal{Y}}(\omega))) + \lambda \|\mathcal{Z}\|_F^2, \quad (46)$$

where  $\lambda > 0$  is the penalty parameter and  $F$  is a large-margin loss satisfying Assumption 2.

The tensor estimation error (11) directly follows from sign tensor estimation error (10) and the proof of Theorem 3. Therefore, it suffices to prove (10). Our proof uses the same techniques used in the proof of Theorem 2. We summarize only the key difference.

Fix  $\pi \notin \mathcal{N}$ . For notational simplicity, we suppress the subscript  $\pi$  and write  $\hat{\mathcal{Z}}$  in place of  $\hat{\mathcal{Z}}_\pi$ . Denote  $n = |\Omega|$  and  $\rho = \rho(\pi, \mathcal{N})$ . Define  $\ell_{\omega, F}(\mathcal{Z}) = |\bar{\mathcal{Y}}(\omega)| \times F(\mathcal{Z}(\omega) \text{sgn}(\bar{\mathcal{Y}}(\omega)))$  and  $\ell_{\omega, F'}(\mathcal{Z}) = |\bar{\mathcal{Y}}(\omega)| \times F'(\mathcal{Z}(\omega) \text{sgn}(\bar{\mathcal{Y}}(\omega)))$  where  $F'$  is T-truncated version of  $F$  such that  $F'(x) = \min(F(x), T)$  with  $T = \max(2, J^2)$ . We focus on the following two empirical processes induced by function  $f_{\mathcal{Z}} \in \mathcal{F}_{\mathcal{T}}$  where  $\mathcal{T} = \{\mathcal{Z} : \text{rank}(\mathcal{Z}) \leq r\}$ ,

$$\frac{1}{n} \sum_{i=1}^n \underbrace{[\ell_{i, F}(\mathcal{Z}, \bar{\mathcal{Y}}) - \ell_{i, F}(\bar{\Theta}, \bar{\mathcal{Y}})]}_{\stackrel{\text{def}}{=} \Delta_{i, F}(f_{\mathcal{Z}}, \bar{\Theta})}, \quad \text{and} \quad \frac{1}{n} \sum_{i=1}^n \underbrace{[\ell_{i, F'}(\mathcal{Z}, \bar{\mathcal{Y}}) - \ell_{i, F'}(\bar{\Theta}, \bar{\mathcal{Y}})]}_{\stackrel{\text{def}}{=} \Delta_{i, F'}(f_{\mathcal{Z}}, \bar{\Theta})}.$$

Note that there is an one-to-one correspondence between sets  $\mathcal{F}_{\mathcal{T}}$  and  $\mathcal{T}$ .

By definition of  $\hat{\mathcal{Z}}$  in (46), we have

$$\frac{1}{n} \sum_{i=1}^n \Delta_{i, F}(f_{\hat{\mathcal{Z}}}, \mathcal{Z}^{(n)}) \leq \lambda J^2 - \lambda \|\hat{\mathcal{Z}}\|_F^2,$$

where  $\mathcal{Z}^{(n)}$  is a sequence of function in Assumption 2(a). Let  $L_n$  denote the desired convergence rate to seek. Then, we have the following inclusion of probability events,

$$\begin{aligned}
& \left\{ (\omega, \mathcal{Y}_\omega) : \text{Risk}_{F'}(\hat{\mathcal{Z}}) - \text{Risk}_{F'}(\bar{\Theta}) \geq 2L_n \right\} \\
& \subset \left\{ (\omega, \mathcal{Y}_\omega) : \exists \mathcal{Z} \text{ s.t. } \text{rank}(\mathcal{Z}) \leq r, \text{Risk}_{F'}(\mathcal{Z}) - \text{Risk}_{F'}(\bar{\Theta}) \geq 2L_n, \right. \\
& \quad \left. \text{and } -\frac{1}{n} \sum_{i=1}^n \Delta_{i,F}(f_{\mathcal{Z}}, \mathcal{Z}^{(n)}) + \lambda J^2 - \lambda \|\hat{\mathcal{Z}}\|_F^2 \geq 0 \right\} \\
& \stackrel{(*)}{\subset} \left\{ (\omega, \mathcal{Y}_\omega) : \exists \mathcal{Z} \text{ s.t. } \text{rank}(\mathcal{Z}) \leq r, \text{Risk}_{F'}(\mathcal{Z}) - \text{Risk}_{F'}(\bar{\Theta}) \geq 2L_n, \right. \\
& \quad \left. \text{and } -\frac{1}{n} \sum_{i=1}^n \Delta_{i,F'}(f_{\mathcal{Z}}, \mathcal{Z}^{(n)}) + \lambda J^2 - \lambda \|\hat{\mathcal{Z}}\|_F^2 \geq 0 \right\} \\
& \subset \left\{ (\omega, \mathcal{Y}_\omega) : \sup_{\substack{\text{rank}(\mathcal{Z}) \leq r \\ \text{Risk}_{F'}(\mathcal{Z}) - \text{Risk}_{F'}(\bar{\Theta}) \geq 2L_n}} -\frac{1}{n} \sum_{i=1}^n \Delta_{i,F'}(f_{\mathcal{Z}}, \mathcal{Z}^{(n)}) + \lambda J^2 - \lambda \|\hat{\mathcal{Z}}\|_F^2 \geq 0 \right\} \\
& \subset \bigcup_{\ell_1, \ell_2=1}^{\infty} \left\{ (\omega, \mathcal{Y}_\omega) : \sup_{\mathcal{Z} \in A_{\ell_1, \ell_2}} -\frac{1}{n} \sum_{i=1}^n \Delta_{i,F'}(f_{\mathcal{Z}}, \mathcal{Z}^{(n)}) + \lambda J^2 - \lambda \|\hat{\mathcal{Z}}\|_F^2 \geq 0 \right\}, \tag{47}
\end{aligned}$$

where  $(*)$  comes from the fact

$$\ell_{\omega, F'}(\mathcal{Z}, \bar{\mathcal{Y}}) \leq \ell_{\omega, F}(\mathcal{Z}, \bar{\mathcal{Y}}) \text{ for all } \mathcal{Z}, \quad \text{and } \ell_{\omega, F'}(\mathcal{Z}^{(n)}, \bar{\mathcal{Y}}) = \ell_{\omega, F}(\mathcal{Z}^{(n)}, \bar{\mathcal{Y}}),$$

because the truncation constant  $T = \max(2, J^2) \geq \max(2, \sup_n \|\mathcal{Z}^{(n)}\|_F^2)$ . In the last line of (47), we have partitioned  $\{\mathcal{Z} : \text{rank}(\mathcal{Z}) \leq r \text{ and } \text{Risk}_{F'}(\mathcal{Z}) - \text{Risk}_{F'}(\bar{\Theta}) \geq 2L_n\}$  into union of  $A_{\ell_1, \ell_2}$  with

$$\begin{aligned}
A_{\ell_1, \ell_2} = & \left\{ \mathcal{Z} : \text{rank}(\mathcal{Z}) \leq r, (\ell_1 + 1)L_n \leq \text{Risk}_{F'}(\mathcal{Z}) - \text{Risk}_{F'}(\bar{\Theta}) < (\ell_1 + 2)L_n, \right. \\
& \left. \text{and } (\ell_2 - 1)J^2 \leq \|\mathcal{Z}\|_F^2 < \ell_2 J^2 \right\},
\end{aligned}$$

for  $\ell_1, \ell_2 = 1, 2, \dots$

Let  $\Gamma$  denote the target probability for the first line in (47). For each  $A_{\ell_1, \ell_2}$ , we consider the centered empirical process,

$$v_n(f_{\mathcal{Z}}) := -\frac{1}{n} \sum_{i=1}^n \left( \Delta_{i,F'}(f_{\mathcal{Z}}, \mathcal{Z}^{(n)}) - \mathbb{E} \Delta_{i,F'}(f_{\mathcal{Z}}, \mathcal{Z}^{(n)}) \right). \tag{48}$$

Notice that

$$\begin{aligned}
\mathbb{E} \Delta_{i,F'}(f_{\mathcal{Z}}, \mathcal{Z}^{(n)}) &= \text{Risk}_{F'}(\mathcal{Z}) - \text{Risk}_{F'}(\bar{\Theta}) + \text{Risk}_{F'}(\bar{\Theta}) - \text{Risk}_{F'}(\mathcal{Z}^{(n)}) \\
&\geq (\ell_1 + 1)L_n - a_n \\
&\geq \ell_1 L_n,
\end{aligned}$$

where the first inequality is from the fact that  $\mathcal{Z} \in A_{\ell_1, \ell_2}$  and Assumption 2(a), and the last inequality uses the condition that  $a_n \lesssim L_n$ .

Combining (47), (48) and the union bound yields

$$\Gamma \leq \sum_{\ell_1, \ell_2=1}^{\infty} \mathbb{P} \left\{ \sup_{\mathcal{Z} \in A_{\ell_1, \ell_2}} v_n(f_{\mathcal{Z}}) \geq \ell_1 L_n + \lambda(\ell_2 - 2)J^2 =: M(\ell_1, \ell_2) \right\}. \quad (49)$$

Similar to the proof of Lemma 1 and Lemma 2 with  $T$ -truncated hinge loss in Lee et al. (2021), the variance of empirical process is bounded by

$$\begin{aligned} \sup_{\mathcal{Z} \in A_{\ell_1, \ell_2}} \text{Var} \Delta_{i, F'}(f_{\mathcal{Z}}, \bar{\Theta}) &\lesssim \sup_{\mathcal{Z} \in A_{\ell_1, \ell_2}} \left( [\mathbb{E} \Delta_{i, F'}(f_{\mathcal{Z}}, \bar{\Theta})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho} \mathbb{E} \Delta_{i, F'}(f_{\mathcal{Z}}, \bar{\Theta}) \right) + \Delta s \\ &\lesssim M(\ell_1, \ell_2)^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho} M(\ell_1, \ell_2) + \Delta s =: V(\ell_1, \ell_2). \end{aligned}$$

To apply Theorem 5, we choose the pair  $(L_n, \lambda)$  satisfying

$$\sup_{\ell_1, \ell_2 \geq 1} \frac{1}{x} \int_x^{\sqrt{x^{\alpha/(\alpha+1)} + x/\rho + \Delta s}} \sqrt{\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}_{\mathcal{T}}(\ell_2), \|\cdot\|_2)} d\varepsilon \lesssim n^{1/2}, \quad (50)$$

where  $x = \ell_1 L_n + \lambda(\ell_2 - 2)J^2$  and  $\mathcal{F}_{\mathcal{T}}(\ell_2) := \{f_{\mathcal{Z}} : \text{rank}(\mathcal{Z}) \leq r, \|\mathcal{Z}\|_F^2 \leq \ell_2 J^2\}$ . Similar to the proof of Lemma 2, we solve the pair  $(L_n, \lambda)$  satisfying (50) as

$$L_n \asymp t_n^{(\alpha+1)/(\alpha+2)} + \frac{t_n}{\rho}, \quad \text{and} \quad \lambda = \frac{L_n}{2J^2}, \quad (51)$$

where  $t_n = \frac{d_{\max} r K \log n}{n}$ . With the choice (51), we bound the right-hand side of (49) based on Theorem 5,

$$\begin{aligned} \Gamma &\lesssim \sum_{\ell_1, \ell_2=1}^{\infty} \exp \left( -\frac{nM^2(\ell_1, \ell_2)}{V(\ell_1, \ell_2) + 2M(\ell_1, \ell_2)} \right) \\ &\lesssim \sum_{\ell_1, \ell_2=1}^{\infty} \exp(-\rho n M(\ell_1, \ell_2)) \\ &\leq \left( \frac{e^{-n\rho L_n}}{1 - e^{-n\rho L_n}} \right) \left( \frac{e^{n\rho \lambda J^2}}{1 - e^{-n\rho \lambda J^2}} \right) \\ &\lesssim e^{-n\rho L_n} \leq e^{-nt_n}, \end{aligned}$$

where the last line uses the fact that  $2\rho\lambda J^2 = \rho L_n \gtrsim t_n \gtrsim n^{-1}$  from (51). The proof is then completed by (37).  $\square$

#### 7.4 Extension of Theorems 2-3 to unbounded observation with sub-Gaussian noise

Consider the signal plus noise model

$$\mathcal{Y} = \Theta + \mathcal{E},$$

where  $\mathcal{E}$  consists of zero-mean, independent noise entries, and  $\Theta \in \mathcal{P}_{\text{sgn}}(r)$  is an  $\alpha$ -smooth tensor. Theoretical results in Section 4 of the main paper are based on bounded observation  $\|\mathcal{Y}\|_{\infty} \leq 1$ . We extend the results to unbounded observation with the following assumption.

**Assumption 3** (Sub-Gaussian noise).

1. There exists a constant  $\beta > 0$ , independent of tensor dimension, such that  $\|\Theta\|_\infty \leq \beta$ . Without loss of generality, we set  $\beta = 1$ .
2. The noise entries  $\mathcal{E}(\omega)$  are independent zero-mean sub-Gaussian random variables with variance proxy  $\sigma^2 > 0$ ; i.e.,  $\mathbb{P}(|\mathcal{E}(\omega)| \geq B) \leq 2e^{-B^2/2\sigma^2}$  for all  $B > 0$ .

We say that an event  $A$  occurs “with high probability” if  $\mathbb{P}(A)$  tends to 1 as the tensor dimension  $d_{\min} = \min_k d_k \rightarrow \infty$ . The following result show that the sub-Gaussian noise incurs an additional  $\log |\Omega|$  factor compared to the bounded case.

**Theorem 6** (Extension to sub-Gaussian noise). *Consider the same condition of Theorem 2. Suppose that Assumption 3 holds. With high probability over training data  $\mathcal{Y}_\Omega$ , we have:*

(a) (Sign matrix estimation). For all  $\pi \notin \mathcal{N}$ ,

$$\text{MAE}(\text{sgn}(\hat{\mathcal{Z}}_\pi), \text{sgn}(\Theta - \pi)) \lesssim t_d^{\frac{\alpha}{\alpha+2}} + \frac{t_d}{\rho^2(\pi, \mathcal{N})}, \text{ where } t_d := \frac{r\sigma^2 d_{\max} \log d_{\max} \log |\Omega|}{|\Omega|}.$$

(b) For all resolution parameter  $H \in \mathbb{N}_+$ ,

$$\text{MAE}(\hat{\Theta}, \Theta) \lesssim (\sigma^2 t_d \log d_{\max} \log H)^{\alpha/(\alpha+2)} + \frac{1 + |\mathcal{N}|}{H} + H(\sigma^2 t_d \log d_{\max} \log H). \quad (52)$$

In particular, setting  $H \asymp \left( \frac{1 + |\mathcal{N}|}{\sigma^2 t_d \log d_{\max}} \right)^{1/2}$  yields the tightest upper bound in (52).

*Proof of Theorem 6.* By setting  $s = K \log(d_{\max})$  in Lemma 4, we have

$$\mathbb{P}(\|\mathcal{E}\|_\infty \geq \sqrt{4\sigma^2 K \log d_{\max}}) \leq 2d_{\max}^{-K}.$$

We divide the sample space into two exclusive events:

- Event I:  $\|\mathcal{E}\|_\infty \geq \sqrt{4\sigma^2 K \log d_{\max}}$ ;
- Event II:  $\|\mathcal{E}\|_\infty < \sqrt{4\sigma^2 K \log d_{\max}}$ .

Because the Event I occurs with probability tending to zero, we restrict ourselves to the Event II only, by following the proof of Theorem 2. We summarize the key difference compared to Section 7.3. We expand the variance by

$$\begin{aligned} \text{Var} [\ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) - \ell_\omega(\bar{\Theta}, \bar{\mathcal{Y}}_\Omega)] &\leq \mathbb{E} |\ell_\omega(\mathcal{Z}(\omega), \bar{\mathcal{Y}}(\omega)) - \ell_\omega(\bar{\Theta}(\omega), \bar{\mathcal{Y}}(\omega))|^2 \\ &= \mathbb{E} |\bar{\mathcal{Y}}(\omega) - \bar{\Theta}(\omega) + \bar{\Theta}(\omega)|^2 |\text{sgn} \mathcal{Z}(\omega) - \text{sgn} \bar{\Theta}(\omega)| \\ &\leq 2(4\sigma^2 K \log d_{\max} + 2) \mathbb{E} |\text{sgn} \mathcal{Z} - \text{sgn} \bar{\Theta}| \\ &\lesssim (\sigma^2 K \log d_{\max}) \text{MAE}(\text{sgn} \mathcal{Z}, \text{sgn} \bar{\Theta}), \end{aligned} \quad (53)$$

where the third line uses the facts  $\|\bar{\Theta}\|_\infty \leq 2$  and  $\|\bar{\mathcal{Y}} - \bar{\Theta}\|_\infty^2 = \|\mathcal{E}\|_\infty^2 < 4\sigma^2 K \log d_{\max}$  within the Event II; the last line comes from the definition of MAE and the asymptotic  $\sigma^2 \log d_{\max} \gg 1$  provided that  $\sigma > 0$  with  $d_{\max}$  sufficiently large.

Based on (53), the  $\alpha$ -smoothness of  $\Theta$  implies that for all measurable functions  $f_{\mathcal{Z}}$ , we have

$$\text{Var} \Delta_i(f_{\mathcal{Z}}, \bar{\Theta}) \lesssim (\sigma^2 K \log d_{\max}) \left\{ [\mathbb{E} \Delta_i(f_{\mathcal{Z}}, \bar{\Theta})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho} \mathbb{E} \Delta_i(f_{\mathcal{Z}}, \bar{\Theta}) + \Delta s \right\}. \quad (54)$$

Based on the proof of Theorem 2, the empirical process with variance-to-mean relationship (54) gives that

$$\mathbb{P}\left(\text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta}) \geq L_n\right) \lesssim \exp(-nt_n), \quad (55)$$

where the convergence rate  $L_n$  is obtained by the same way in the proof of Lemma 2,

$$L_n \asymp t_n^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho} t_n, \quad \text{with } t_n = \frac{r\sigma^2 d_{\max} \log d_{\max} \log n}{n}, \quad (56)$$

where constants (possibly depending on  $K$ ) have been absorbed into the  $\asymp$  relationship. Combining (55) and (56), we obtain that, with high probability,

$$\text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta}) \lesssim \left( \frac{r\sigma^2 d_{\max} \log d_{\max} \log |\Omega|}{|\Omega|} \right)^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho(\pi, \mathcal{N})} \left( \frac{r\sigma^2 d_{\max} \log d_{\max} \log |\Omega|}{|\Omega|} \right), \quad (57)$$

Therefore, combining (57) and (37) completes the proof. The tensor estimation error follows readily from the proof of Theorem 3 and Theorem 6.  $\square$

**Lemma 4** (sub-Gaussian maximum). *Let  $X_1, \dots, X_n$  be independent sub-Gaussian zero-mean random variables with variance proxy  $\sigma^2$ . Then, for any  $s > 0$ ,*

$$\mathbb{P}\left\{\max_{1 \leq i \leq n} |X_i| \geq \sqrt{2\sigma^2(\log n + s)}\right\} \leq 2e^{-s}.$$

*Proof of Lemma 4.* The conclusion follows from

$$\mathbb{P}\left[\max_{1 \leq i \leq n} |X_i| \geq u\right] \leq \sum_{i=1}^n \mathbb{P}[X_i \geq u] \leq 2ne^{-\frac{u^2}{2\sigma^2}} = 2e^{-s},$$

where we set  $u = \sqrt{2\sigma^2(\log n + s)}$ .  $\square$

## 8 Conclusion

We have developed a tensor estimation method that addresses both low- and high-rankness. Our sign series representation leads to nonparametric guarantees previously impossible. We hope the work opens up new inquiry that allows more researchers to contribute to this field.

## Acknowledgements

This research is supported in part by NSF grant DMS-1915978, DMS-2023239 and Wisconsin Alumni Research Foundation.

## References

Alon, N., Moran, S., and Yehudayoff, A. (2016). Sign rank versus VC dimension. In *Conference on Learning Theory*, pages 47–80.

- Alquier, P., Cottet, V., Lecué, G., et al. (2019). Estimation bounds and sharp oracle inequalities of regularized procedures with lipschitz loss functions. *Annals of Statistics*, 47(4):2117–2144.
- Anandkumar, A., Ge, R., Hsu, D., Kakade, S. M., and Telgarsky, M. (2014). Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15(1):2773–2832.
- Anandkumar, A., Ge, R., and Janzamin, M. (2017). Analyzing tensor power method dynamics in overcomplete regime. *Journal of Machine Learning Research*, 18(1):752–791.
- Balabdaoui, F., Durot, C., and Jankowski, H. (2019). Least squares estimation in the monotone single index model. *Bernoulli*, 25(4B):3276–3310.
- Bartlett, P. L., Jordan, M. I., and McAuliffe, J. D. (2006). Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156.
- Cai, C., Li, G., Poor, H. V., and Chen, Y. (2019). Nonconvex low-rank tensor completion from noisy data. In *Advances in Neural Information Processing Systems*, pages 1863–1874.
- Chan, S. and Airoldi, E. (2014). A consistent histogram estimator for exchangeable graph models. In *International Conference on Machine Learning*, pages 208–216.
- Chi, E. C., Gaines, B. J., Sun, W. W., Zhou, H., and Yang, J. (2020). Provable convex co-clustering of tensors. *Journal of Machine Learning Research*, 21(214):1–58.
- Cohn, H. and Umans, C. (2013). Fast matrix multiplication using coherent configurations. In *Proceedings of the twenty-fourth annual ACM-SIAM symposium on Discrete algorithms*, pages 1074–1087.
- De Lathauwer, L., De Moor, B., and Vandewalle, J. (2000). A multilinear singular value decomposition. *SIAM Journal on Matrix Analysis and Applications*, 21(4):1253–1278.
- De Wolf, R. (2003). Nondeterministic quantum query and communication complexities. *SIAM Journal on Computing*, 32(3):681–699.
- Fan, J. and Udell, M. (2019). Online high rank matrix completion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8690–8698.
- Ganti, R., Rao, N., Balzano, L., Willett, R., and Nowak, R. (2017). On learning high dimensional structured single index models. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, pages 1898–1904.
- Ganti, R. S., Balzano, L., and Willett, R. (2015). Matrix completion under monotonic single index models. In *Advances in Neural Information Processing Systems*, pages 1873–1881.
- Genzel, M. and Stollenwerk, A. (2020). Robust 1-bit compressed sensing via hinge loss minimization. *Information and Inference: A Journal of the IMA*, 9(2):361–422.
- Ghadermarzy, N., Plan, Y., and Yilmaz, O. (2018). Learning tensors from partial binary measurements. *IEEE Transactions on Signal Processing*, 67(1):29–40.
- Ghadermarzy, N., Plan, Y., and Yilmaz, Ö. (2019). Near-optimal sample complexity for convex tensor completion. *Information and Inference: A Journal of the IMA*, 8(3):577–619.



- Globerson, A., Chechik, G., Pereira, F., and Tishby, N. (2007). Euclidean embedding of co-occurrence data. *Journal of Machine Learning Research*, 8:2265–2295.
- Han, R., Willett, R., and Zhang, A. (2020). An optimal statistical and computational framework for generalized tensor estimation. *arXiv preprint arXiv:2002.11255*.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media.
- He, L., Lu, C.-T., Ma, G., Wang, S., Shen, L., Philip, S. Y., and Ragin, A. B. (2017). Kernelized support tensor machines. In *International Conference on Machine Learning*, pages 1442–1451. PMLR.
- Hillar, C. J. and Lim, L.-H. (2013). Most tensor problems are NP-hard. *Journal of the ACM (JACM)*, 60(6):45.
- Hitchcock, F. L. (1927). The expression of a tensor or a polyadic as a sum of products. *Journal of Mathematics and Physics*, 6(1-4):164–189.
- Hong, D., Kolda, T. G., and Duersch, J. A. (2020). Generalized canonical polyadic tensor decomposition. *SIAM Review*, 62(1):133–163.
- Hore, V., Viñuela, A., Buil, A., Knight, J., McCarthy, M. I., Small, K., and Marchini, J. (2016). Tensor decomposition for multiple-tissue gene expression experiments. *Nature genetics*, 48(9):1094.
- Jain, P. and Oh, S. (2014). Provable tensor factorization with missing data. In *Advances in Neural Information Processing Systems*, volume 27, pages 1431–1439.
- Kolda, T. G. and Bader, B. W. (2009). Tensor decompositions and applications. *SIAM Review*, 51(3):455–500.
- Kosorok, M. R. (2007). *Introduction to empirical processes and semiparametric inference*. Springer Science & Business Media.
- Lee, C., Li, L., Zhang, H. H., and Wang, M. (2021). Nonparametric trace regression in high dimensions via sign series representation. *arXiv preprint arXiv:2105.01783*.
- Lee, C. and Wang, M. (2020). Tensor denoising and completion based on ordinal observations. In *International Conference on Machine Learning*, pages 5778–5788.
- Li, Y., Liu, Y., Li, J., Qin, W., Li, K., Yu, C., and Jiang, T. (2009). Brain anatomical network and intelligence. *PLoS Comput Biol*, 5(5):e1000395.
- Montanari, A. and Sun, N. (2018). Spectral algorithms for tensor completion. *Communications on Pure and Applied Mathematics*, 71(11):2381–2425.
- Mu, C., Huang, B., Wright, J., and Goldfarb, D. (2014). Square deal: Lower bounds and improved relaxations for tensor recovery. In *International Conference on Machine Learning*, pages 73–81.
- Ongie, G., Willett, R., Nowak, R. D., and Balzano, L. (2017). Algebraic variety models for high-rank matrix completion. In *International Conference on Machine Learning*, pages 2691–2700.
- Robinson, P. M. (1988). Root-N-consistent semiparametric regression. *Econometrica: Journal of the Econometric Society*, 56(4):931–954.

- Scott, C. (2011). Surrogate losses and regret bounds for cost-sensitive classification with example-dependent costs. In *ICML*.
- Shen, X. and Wong, W. H. (1994). Convergence rate of sieve estimates. *The Annals of Statistics*, 22:580–615.
- Wang, L., Durante, D., Jung, R. E., and Dunson, D. B. (2017). Bayesian network–response regression. *Bioinformatics*, 33(12):1859–1866.
- Wang, M. and Li, L. (2020). Learning from binary multiway data: Probabilistic tensor decomposition and its statistical optimality. *Journal of Machine Learning Research*, 21(154):1–38.
- Wang, M. and Zeng, Y. (2019). Multiway clustering via tensor block models. In *Advances in Neural Information Processing Systems*, pages 713–723.
- Yuan, M. and Zhang, C.-H. (2016). On tensor completion via nuclear norm minimization. *Foundations of Computational Mathematics*, 16(4):1031–1068.
- Zhang, A. and Xia, D. (2018). Tensor SVD: Statistical and computational limits. *IEEE Transactions on Information Theory*, 64(11):7311 – 7338.