

# Appendix for “Beyond the Signs: Nonparametric Tensor Completion via Sign Series”

The appendix consists of additional theoretical results (Section A), numerical experiments (Section B), and proofs (Section C).

## A Additional results

### A.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.

### A.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 A.1** (Max hypergraphon). 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.* We first prove the results for  $K = 2$ . The full-rankness of  $\Theta$  is verified from elementary row

operations as follows

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

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

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

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 (1),

$$\text{sgn}(\Theta(i_1, \dots, i_K) - \pi) = \begin{cases} -1, & i_k < i^* \text{ for all } k \in [K]; \\ 1, & \text{otherwise,} \end{cases} \quad (2)$$

where  $i^*$  denotes the index that satisfies  $\log(1 + \frac{i^*-1}{d}) < \pi \leq \log(1 + \frac{i^*}{d})$ . The equation (2) 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 A.1 is a special case of the following Proposition.

**Proposition A.1** (Min/Max hypergraphon). *Let  $\mathcal{Z}_{\max} \in \mathbb{R}^{d_1 \times \dots \times d_K}$  denote a tensor with entries*

$$\mathcal{Z}_{\max}(i_1, \dots, i_K) = \max(x_{i_1}^{(1)}, \dots, x_{i_K}^{(K)}), \quad (3)$$

where  $x_{i_k}^{(k)} \in [0, 1]$  are given numbers for all  $i_k \in [d_k]$ . Let  $g: \mathbb{R} \rightarrow \mathbb{R}$  be a continuous function and  $\Theta := g(\mathcal{Z}_{\max})$  be the transformed tensor. For a given  $\pi \in [-1, 1]$ , suppose the function  $g(z) = \pi$  has at most  $r \geq 1$  distinct real roots. 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 (3).

*Proof.* 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 A.1, 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  $\text{sgn}(\Theta - \pi) = -1$ ,

$$\begin{aligned} \{\omega : \Theta(\omega) - \pi < 0\} &= \{\omega : g(\mathcal{Z}_{\max}(\omega)) < \pi\} \\ &= \cup_{s \in \mathcal{I}} \{\omega : \mathcal{Z}_{\max}(\omega) \in (z_s, z_{s+1})\} \\ &= \cup_{s \in \mathcal{I}} \left( \{\omega : x_{i_k}^{(k)} < z_{s+1} \text{ for all } k \in [K]\} \cap \{\omega : x_{i_k}^{(k)} \leq z_s \text{ for all } k \in [K]\}^c \right). \end{aligned} \quad (4)$$

The equation (4) is equivalent to

$$\mathbb{1}(\Theta(i_1, \dots, i_K) < \pi) = \sum_{s \in \mathcal{I}} \left( \prod_k \mathbb{1}(x_{i_k}^{(k)} < z_{s+1}) - \prod_k \mathbb{1}(x_{i_k}^{(k)} \leq z_s) \right), \quad (5)$$

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 (5) implies the low-rank representation of  $\text{sgn}(\Theta - \pi)$ ,

$$\text{sgn}(\Theta - \pi) = 1 - 2 \sum_{s \in \mathcal{I}} \left( \mathbf{a}_{s+1}^{(1)} \otimes \dots \otimes \mathbf{a}_{s+1}^{(K)} - \bar{\mathbf{a}}_s^{(1)} \otimes \dots \otimes \bar{\mathbf{a}}_s^{(K)} \right), \quad (6)$$

where we have denoted the two binary vectors

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

Therefore, by (6) 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 A.2.** *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.* The conclusion directly follows from the definition of tensor rank.  $\square$

**Example A.2** (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.* 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 directly shows that  $\mathbf{M}$  is full rank as follows,

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

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

$$\mathbf{A} = \mathbf{b} \otimes \text{rev}(\mathbf{b}) + \text{rev}(\mathbf{b}) \otimes \mathbf{b}. \quad (7)$$

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

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

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

$$\mathbf{A}(i, j) \geq \mathbf{A}(i', j'), \quad \text{for all } |i - j| \geq |i' - j'|. \quad (8)$$

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 (8). By definition (7),  $\mathbf{A}$  is a rank-2 matrix. Henceforce,  $\text{srnk}(\mathbf{M} - \pi) = \text{srnk}(\mathbf{A} - \pi') \leq 3$ .  $\square$

**Remark A.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 A.3** (Stacked identity matrices). Let  $\mathbf{I}$  be a  $d$ -by- $d$  identity matrix. Then

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

*Proof.* Depending on the value of  $\pi$ , the sign matrix  $\text{sgn}(\mathbf{I} - \pi)$  falls into one of the three cases: 1)  $\text{sgn}(\mathbf{I} - \pi)$  is a matrix of all 1; 2)  $\text{sgn}(\mathbf{I} - \pi)$  is a matrix of all  $-1$ ; 3)  $\text{sgn}(\mathbf{I} - \pi) = 2\mathbf{I} - \mathbf{1}_d \otimes \mathbf{1}_d$ . The former two cases are trivial, so it suffices to show  $\text{srnk}(\mathbf{I} - \pi) \leq 3$  in the third case.

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

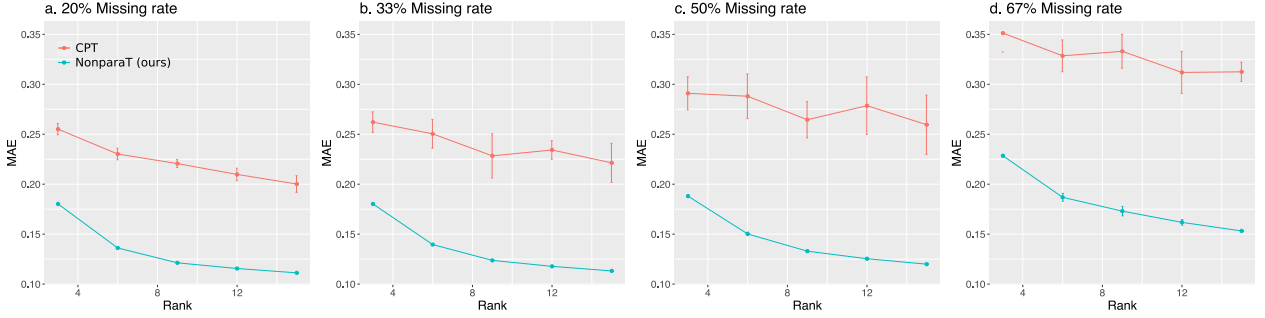
$$\mathbf{A}(i, j) \begin{cases} = 2^{-d}, & i = j, \\ \geq 2^{-d} + 2^{-d-2}, & i \neq j. \end{cases}$$

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

## B Additional data analysis

### B.1 Brain connectivity analysis

Figure S1 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 S1 and Table 2 (in the main paper), 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  $[0, 1]$ .



Supplementary Figure S1: 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.

We next investigate the pattern in the estimated signal tensor. Figure 5 of the main paper 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 our method with  $r = 10$  and  $H = 20$ . Then, we perform a regression analysis of  $\hat{\Theta}(i, j, : ) \in \mathbb{R}^{114}$  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 the role of interhemispheric connectivity in human intelligence.

### B.2 NIPS data analysis

In the main paper we have summarized the MAE in cross-validation experiments for  $r = 6, 9, 12$ . Here we provide additional results for a wider range  $r = 3, 6, \dots, 15$ . Table S1 suggests that further increment of rank appears to have little effect on the performance. In addition, we also perform naive imputation where the missing values are predicted using the sample average. The two tensor methods outperform the naive imputation, implying the necessity of incorporating tensor structure in the analysis.

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	0.32(.001)				

Supplementary Table S1: Prediction accuracy measured in MAE in the NIPS data analysis. The reported MAEs are averaged over five runs of cross-validation, 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  $H = 20$ .

## C Proofs

### C.1 Proofs of Propositions 1-3

*Proof of Proposition 1.* 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  $\text{srnk}(\Theta) = \text{srnk}(g(\Theta)) \leq \text{rank}(g(\Theta))$ .  $\square$

*Proof of Proposition 2.* See Section A.2 for constructive examples.  $\square$

*Proof of Proposition 3.* 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

$$\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) = \mathbb{E}_{\omega \sim \Pi} \underbrace{\mathbb{E}_{\mathcal{Y}(\omega)} \{ |\mathcal{Y}(\omega) - \pi| [ |\mathcal{Z}(\omega) - \text{sgn}(\bar{\mathcal{Y}}(\omega))| - |\bar{\Theta}(\omega) - \text{sgn}(\bar{\mathcal{Y}}(\omega))| ] \}}_{\stackrel{\text{def}}{=} I(\omega)}. \quad (9)$$

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

$$\begin{aligned} I(\omega) &= \mathbb{E}_y [(y - \pi)(\bar{\theta} - z)\mathbf{1}(y \geq \pi) + (\pi - y)(z - \bar{\theta})\mathbf{1}(y < \pi)] \\ &= \mathbb{E}_y [(\bar{\theta} - z)(y - \pi)] \\ &= [\text{sgn}(\theta - \pi) - z](\theta - \pi) \\ &= |\text{sgn}(\theta - \pi) - z||\theta - \pi| \geq 0, \end{aligned} \quad (10)$$

where the third line uses the fact  $\mathbb{E}_y y = \theta$  and  $\bar{\theta} = \text{sgn}(\theta - \pi)$ , and the last line uses the assumption  $z \in \{-1, 1\}$ . The equality (10) is attained when  $z = \text{sgn}(\theta - \pi)$  or  $\theta = \pi$ . Combining (10) with (9), we conclude that, for all  $\mathcal{Z} \in \{-1, 1\}^{d_1 \times \dots \times d_K}$ ,

$$\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) = \mathbb{E}_{\omega \sim \Pi} |\text{sgn}(\Theta(\omega) - \pi) - \mathcal{Z}(\omega)||\Theta(\omega) - \pi| \geq 0, \quad (11)$$

In particular, setting  $\mathcal{Z} = \bar{\Theta} = \text{sgn}(\Theta - \pi)$  in (11) 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.  $\square$

## C.2 Proof of Theorem 1

*Proof of Theorem 1.* Fix  $\pi \in [-1, 1]$ . Based on (11) in Proposition 3 we have

$$\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) = \mathbb{E} [|\text{sgn}\mathcal{Z} - \text{sgn}\bar{\Theta}||\bar{\Theta}|]. \quad (12)$$

The Assumption 1 states that

$$\mathbb{P}(|\bar{\Theta}| \leq t) \leq ct^\alpha, \quad \text{for all } 0 \leq t < \rho(\pi, \mathcal{N}). \quad (13)$$

Without future 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 (12), for all  $0 \leq t < \rho(\pi, \mathcal{N})$ ,

$$\begin{aligned} \text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) &\geq t\mathbb{E} \left( |\text{sgn}\mathcal{Z} - \text{sgn}\hat{\Theta}| \mathbb{1}\{|\hat{\Theta}| > t\} \right) \\ &\geq 2t\mathbb{P}(\text{sgn}\mathcal{Z} \neq \text{sgn}\bar{\Theta} \text{ and } |\bar{\Theta}| > t) \\ &\geq 2t \left\{ \mathbb{P}(\text{sgn}\mathcal{Z} \neq \text{sgn}\bar{\Theta}) - \mathbb{P}(|\bar{\Theta}| \leq t) \right\} \\ &\geq t \left\{ \text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) - 2ct^\alpha \right\}, \end{aligned} \quad (14)$$

where the last line follows from the definition of MAE and (13). We maximize the lower bound (14) with respect to  $t$ , and obtain the optimal  $t_{\text{opt}}$ ,

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

The corresponding lower bound of the inequality (14) becomes

$$\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) \geq \begin{cases} c_1 \rho(\pi, \mathcal{N}) \text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}), & \text{if } \text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) > 2c(1 + \alpha)\rho^\alpha(\pi, \mathcal{N}), \\ c_2 [\text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta})]^{\frac{1+\alpha}{\alpha}}, & \text{if } \text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) \leq 2c(1 + \alpha)\rho^\alpha(\pi, \mathcal{N}), \end{cases}$$

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

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

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

- Case 2:  $\alpha = \infty$ . The inequality (14) now becomes

$$\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta}) \geq t \text{MAE}(\text{sgn}\bar{\Theta}, \text{sgn}\mathcal{Z}), \quad \text{for all } 0 \leq t < \rho(\pi, \mathcal{N}). \quad (15)$$

The conclusion follows by taking  $t = \frac{\rho(\pi, \mathcal{N})}{2}$  in the inequality (15). □

**Remark C.1.** The proof of Theorem 1 shows that, under Assumption 1,

$$\text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}) \lesssim [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})], \quad (16)$$

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

### C.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).

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

$$\begin{aligned} L(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) &= \frac{1}{|\Omega|} \sum_{\omega \in \Omega} \underbrace{|\bar{\mathcal{Y}}(\omega)|}_{\text{weight}} \times \underbrace{|\text{sgn}\mathcal{Z}(\omega) - \text{sgn}\bar{\mathcal{Y}}(\omega)|}_{\text{classification loss}} \\ &= \frac{1}{|\Omega|} \sum_{\omega \in \Omega} \ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}), \end{aligned} \quad (17)$$

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  $(\alpha, \pi)$ -smoothness of  $\Theta$ , we have

$$\text{Var}[\ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}) - \ell_\omega(\bar{\Theta}, \bar{\mathcal{Y}}_\Omega)] \lesssim [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho(\pi, \mathcal{N})} [\text{Risk}(\mathcal{Z}) - \text{Risk}(\bar{\Theta})], \quad (18)$$

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 C.1.* We expand the variance by

$$\begin{aligned} \text{Var}[\ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) - \ell_\omega(\bar{\Theta}, \bar{\mathcal{Y}}_\Omega)] &\lesssim \mathbb{E}|\ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) - \ell_\omega(\bar{\Theta}, \bar{\mathcal{Y}}_\Omega)|^2 \\ &\lesssim \mathbb{E}|\ell_\omega(\mathcal{Z}, \bar{\mathcal{Y}}_\Omega) - \ell_\omega(\bar{\Theta}, \bar{\mathcal{Y}}_\Omega)| \\ &\leq \mathbb{E}|\text{sgn}\mathcal{Z} - \text{sgn}\bar{\Theta}| = \text{MAE}(\text{sgn}\mathcal{Z}, \text{sgn}\bar{\Theta}), \end{aligned} \quad (19)$$

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 (18) then directly follows by applying Remark C.1 to (19).  $\square$

*Proof of Theorem 2.* Fix  $\pi \in [-1, 1]$ . 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

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

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 (17) 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

$$L(\mathcal{Z}, \bar{\mathcal{Y}}_{\Omega}) - L(\bar{\Theta}, \bar{\mathcal{Y}}_{\Omega}) = \frac{1}{n} \sum_{i=1}^n \underbrace{[\ell_i(\mathcal{Z}, \bar{\mathcal{Y}}) - \ell_i(\bar{\Theta}, \bar{\mathcal{Y}})]}_{\stackrel{\text{def}}{=} \Delta_i(f_{\mathcal{Z}}, \bar{\mathcal{Y}})}, \quad (20)$$

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  $\mathcal{T}$ .

Our remaining proof adopts the techniques of Wang et al. (2008, Theorem 3) to bound (20) over the function family  $f_{\mathcal{Z}} \in \mathcal{F}_{\mathcal{T}}$ . We summarize only the key difference here but refer to Wang et al. (2008) for complete proof. Based on Lemma C.1, the  $(\alpha, \pi)$ -smoothness of  $\Theta$  implies

$$\text{Var} \Delta_i(f_{\mathcal{Z}}, \bar{\mathcal{Y}}) \lesssim [\mathbb{E} \Delta_i(f_{\mathcal{Z}}, \bar{\mathcal{Y}})]^{\frac{\alpha}{1+\alpha}} + \frac{1}{\rho} \mathbb{E} \Delta_i(f_{\mathcal{Z}}, \bar{\mathcal{Y}}), \quad \text{for all } f_{\mathcal{Z}} \in \mathcal{F}_{\mathcal{T}}. \quad (21)$$

Applying local iterative techniques in Wang et al. (2008, Theorem 3) to the empirical process (20) with the variance-to-mean relationship (21) gives that

$$\mathbb{P} \left( \text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta}) \geq L_n \right) \lesssim \exp(-nL_n), \quad (22)$$

where the convergence rate  $L_n > 0$  is determined by the solution to the following inequality,

$$\frac{1}{L_n} \int_{L_n}^{\sqrt{L_n^{\alpha/(\alpha+1)} + \frac{L_n}{\rho}}} \sqrt{\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}_{\mathcal{T}}, \|\cdot\|_2)} d\varepsilon \leq C\sqrt{n}, \quad (23)$$

for some constant  $C > 0$ . In particular, the smallest  $L_n$  satisfying (23) yields the best upper bound of the error rate. Here  $\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}_{\mathcal{T}}, \|\cdot\|_2)$  denotes the  $L_2$ -metric,  $\varepsilon$ -bracketing number (c.f. Definition C.1) of family  $\mathcal{F}_{\mathcal{T}}$ .

It remains to solve for the smallest possible  $L_n$  in (23). Based on Lemma C.2, the inequality (23) is satisfied with

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

Therefore, by (22), with very high probability.

$$\text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta}) \leq t_n^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho} t_n.$$

Inserting the above bound into (16) gives

$$\begin{aligned} \text{MAE}(\text{sgn} \hat{\mathcal{Z}}, \text{sgn} \bar{\Theta}) &\lesssim [\text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta})]^{\alpha/(\alpha+1)} + \frac{1}{\rho} [\text{Risk}(\hat{\mathcal{Z}}) - \text{Risk}(\bar{\Theta})] \\ &\lesssim t_n^{\alpha/(\alpha+2)} + \frac{1}{\rho^{\alpha/(\alpha+1)}} t_n^{\alpha/(\alpha+1)} + \frac{1}{\rho} t_n^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho^2} t_n \\ &\leq 4t_n^{\alpha/(\alpha+2)} + \frac{4}{\rho^2} t_n, \end{aligned} \quad (24)$$

where 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 (24) and absorb the term  $K \log K$  into the constant. The conclusion is then proved.  $\square$

**Definition C.1** (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 C.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*

$$L_n \asymp \left( \frac{d_{\max} r K \log K}{n} \right)^{(\alpha+1)/(\alpha+2)} + \frac{1}{\rho(\pi, \mathcal{N})} \left( \frac{d_{\max} r K \log K}{n} \right).$$

Then, the following inequality is satisfied.

$$\frac{1}{L_n} \int_{L_n}^{\sqrt{L_n^{\alpha/(\alpha+1)} + \frac{L_n}{\rho(\pi, \mathcal{N})}}} \sqrt{\mathcal{H}_{[\cdot]}(\varepsilon, \mathcal{F}_{\mathcal{T}}, \|\cdot\|_2)} d\varepsilon \leq C n^{1/2}, \quad (25)$$

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

*Proof of Lemma C.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 (25) gives

$$g(L) = \frac{1}{L} \int_L^{\sqrt{L^{\alpha/(\alpha+1)} + \rho^{-1}L}} \sqrt{d_{\max} r K \log \left( \frac{K}{\varepsilon} \right)} d\varepsilon. \quad (26)$$

By the monotonicity of the integrand in (26), we bound  $g(L)$  by

$$\begin{aligned} g(L) &\leq \frac{\sqrt{d_{\max} r K}}{L} \int_L^{\sqrt{L^{\alpha/(\alpha+1)} + \rho^{-1}L}} \sqrt{\log \left( \frac{K}{\varepsilon} \right)} d\varepsilon \\ &\leq \sqrt{d_{\max} r K (\log K - \log L)} \left( \frac{L^{\alpha/(2\alpha+2)} + \sqrt{\rho^{-1}L}}{L} - 1 \right) \\ &\leq \sqrt{d_{\max} r K \log K} \left( \frac{1}{L^{(\alpha+2)/(2\alpha+2)}} + \frac{1}{\sqrt{\rho L}} \right), \end{aligned} \quad (27)$$

where the second line follows from  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for  $a, b > 0$ . It remains to verify that  $g(L_n) \leq Cn^{1/2}$  for  $L_n$  specified in (25). Plugging  $L_n$  into the last line of (27) gives

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

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

#### C.4 Proof of Theorem 3

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

$$\begin{aligned} \text{MAE}(\hat{\Theta}, \Theta) &= \mathbb{E} \left| \frac{1}{2H+1} \sum_{\pi \in \Pi} \text{sgn} \hat{Z}_\pi - \Theta \right| \\ &\leq \mathbb{E} \left| \frac{1}{2H+1} \sum_{\pi \in \Pi} (\text{sgn} \hat{Z}_\pi - \text{sgn}(\Theta - \pi)) \right| + \mathbb{E} \left| \frac{1}{2H+1} \sum_{\pi \in \Pi} \text{sgn}(\Theta - \pi) - \Theta \right| \\ &\leq \frac{1}{2H+1} \sum_{\pi \in \Pi} \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)) + \frac{1}{H}, \end{aligned} \quad (28)$$

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

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

Write  $n = |\Omega|$ . Now it suffices to bound the first term in (28). We prove that

$$\frac{1}{2H+1} \sum_{\pi \in \Pi} \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)) \lesssim t_n^{\alpha/(\alpha+2)} + \frac{1}{H} + H t_n, \quad \text{with } t_n = \frac{d_{\max} r K \log K}{n}. \quad (29)$$

Theorem 2 implies that the sign estimation accuracy depends on the closeness of  $\pi \in \mathcal{H}$  to the mass points in  $\mathcal{H}$ . Therefore, we partition the level set  $\pi \in \mathcal{H}$  based on their closeness to  $\mathcal{H}$ . Specifically, let  $\mathcal{N}_H \stackrel{\text{def}}{=} \bigcup_{\pi' \in \mathcal{N}} (\pi' - \frac{1}{H}, \pi' + \frac{1}{H})$  denote the set of levels at least  $\frac{1}{H}$ -close to the mass points. We expand (29) by

$$\begin{aligned} &\frac{1}{2H+1} \sum_{\pi \in \Pi} \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)) \\ &= \frac{1}{2H+1} \sum_{\pi \in \Pi \cap \mathcal{W}_H} \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)) + \frac{1}{2H+1} \sum_{\pi \in \Pi \cap \mathcal{W}_H^c} \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)). \end{aligned} \quad (30)$$

By assumption, the first term involves only finite number of summands and thus can be bounded by  $4C/(2H+1)$  where  $C > 0$  is a constant such that  $|\mathcal{N}| \leq C$ . We bound the second term using

the explicit forms of  $\rho(\pi, \mathcal{N})$  in the sequence  $\pi \in \Pi \cap \mathcal{N}_H^c$ . Based on Theorem 2,

$$\begin{aligned}
\frac{1}{2H+1} \sum_{\pi \in \Pi \cap \mathcal{N}_H^c} \text{MAE}(\text{sgn} \hat{Z}_\pi, \text{sgn}(\Theta - \pi)) &\lesssim \frac{1}{2H+1} \sum_{\pi \in \Pi \cap \mathcal{N}_H^c} t_n^{\alpha/(\alpha+2)} + \frac{t_n}{2H+1} \sum_{\pi \in \Pi \cap \mathcal{N}_H^c} \frac{1}{\rho^2(\pi, \mathcal{N})} \\
&\leq t_n^{\alpha/(\alpha+2)} + \frac{t_n}{2H+1} \sum_{\pi \in \Pi \cap \mathcal{N}_H^c} \sum_{\pi' \in \mathcal{N}} \frac{1}{|\pi - \pi'|^2} \\
&\leq t_n^{\alpha/(\alpha+2)} + \frac{t_n}{2H+1} \sum_{\pi' \in \mathcal{N}} \sum_{\pi \in \Pi \cap \mathcal{N}_H^c} \frac{1}{|\pi - \pi'|^2} \\
&\leq t_n^{\alpha/(\alpha+2)} + 2CHt_n,
\end{aligned}$$

where the last inequality follows from the Lemma C.3. Combining the bounds for the two terms in (30) completes the proof for conclusion (29). Finally, plugging (29) into (28) yields

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

The conclusion follows by absorbing  $K \log K$  into the constant term in the statement.  $\square$

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

$$\sum_{\pi \in \Pi \cap \mathcal{N}_H^c} \frac{1}{|\pi - \pi'|^2} \leq 4H^2.$$

*Proof of Lemma C.3.* Notice that all points  $\pi \in \Pi \cap \mathcal{N}_H^c$  satisfy  $|\pi - \pi'| > \frac{1}{H}$  for all  $\pi' \in \mathcal{N}$ . We use this fact to compute the sum

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

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

## References

- Kosorok, M. R. (2007). *Introduction to empirical processes and semiparametric inference*. Springer Science & Business Media.
- Mu, C., B. Huang, J. Wright, and D. Goldfarb (2014). Square deal: Lower bounds and improved relaxations for tensor recovery. In *International Conference on Machine Learning*, pp. 73–81.

- Shen, X. and W. H. Wong (1994). Convergence rate of sieve estimates. *The Annals of Statistics* 22, 580–615.
- Wang, J., X. Shen, and Y. Liu (2008). Probability estimation for large-margin classifiers. *Biometrika* 95(1), 149–167.