

# High-Rank Tensor Estimation Under Smooth Single Index Models

Chanwoo Lee

Miaoyan Wang

Department of Statistics, University of Wisconsin - Madison

CHANWOO.LEE@WISC.EDU

MIAOYAN.WANG@WISC.EDU

## Abstract

We study the semi-parametric single index model with a near low-rank regression coefficient matrix, which extends notion of linear trace regression. Specifically, given matrix covariate  $\mathcal{Y} = f(\mathbf{B}) + \mathcal{E}$  with  $f: \mathbb{R} \rightarrow \mathbb{R}$  an *unknown* analytic function, and  $\mathbf{B}$  is an unknown low-rank coefficient matrix. This model accommodates various types of responses and embraces many important problem setups such as reduced-rank regression, matrix-response regression, 1-bit matrix completion and compressed sensing among others. We quantify the spectral decay of  $\mathbb{E}(\mathcal{Y})$ , develop a higher-order spectral algorithm, and show the distinct behavior between the matrix and tensor version of the problem. We first establish a general theory and then for each specific problem, we derive explicitly the statistical rate of the proposed estimator. They all match the minimax rates in the nonparametric regression up to logarithmic factors. Numerical studies confirm the rates we established and demonstrate the advantage of generalized trace regression over linear trace regression when the response is dichotomous. We also show the benefit of incorporating nuclear norm regularization in dynamic stock return prediction and in image classification.

**Keywords:** Nonparametric regression, matrix completion, Holder smoothness

## 1. Introduction

We use  $f^{(\ell)}(x)$  to denote the  $\ell$ -th derivative of  $f$ , evaluated at  $x$ . We use  $\lesssim$  to denote two sequence up to a universal constant  $C$ .

## 2. Model

**Definition 1 (Analytic function class)** Let  $C > 0$  be a positive constant. The analytic function class  $\mathcal{F}(C)$  on  $[-1, 1]$  is defined as the set of functions  $f: [-1, 1] \rightarrow \mathbb{R}$  whose derivatives satisfy

$$\sup_{x \in [-1, 1]} |f^{(\ell)}(x)| \leq C^{\ell+1} \ell! \quad \text{for all } \ell \in \mathbb{N}_+.$$

Equivalently,  $f$  is infinitely differentiable and its Taylor expansion around any point in its domain converges to the function.

A higher-order tensor  $\mathcal{T}$  can be unfolded into a matrix. We now introduce several quantities that controls the complexity of matrix unfolding. We use  $\text{Mat}(\mathcal{T})$  to denote the matrix unfolding. We use  $\text{Trank}(\cdot)$  and  $\lambda(\mathcal{T})$  the rank and spectral norm of the matrix  $\text{Mat}(\mathcal{T})$ ,

$$\text{Trank}(\mathcal{T}) := \text{rank}(\text{Mat}(\mathcal{T})), \quad \lambda(\mathcal{T}) := \|\text{Mat}(\mathcal{T})\|_{\text{sp}}.$$

The following family consists of  $d$ -dimensional order- $m$  tensor with Tucker rank bounded by  $r$ .

**Definition 2 (Low-rank tensor class)** *The family of  $d$ -dimensional rank- $s$  tensor  $\mathcal{T}(d, s, m)$  is defined as the set of tensors with Tucker rank bounded by  $r$ :*

$$\mathcal{T}(d, s, m) = \{\mathcal{T} \in (\mathbb{R}^d)^{\otimes m} : \text{Trank}(\mathcal{T}) \leq s \text{ and } \lambda(\mathcal{T}) \leq 1\}$$

Equivalently, the tensor in class  $\mathcal{T}(d, s, m)$  admits the rank- $s$  Tucker decomposition:

$$\mathcal{T} = \mathcal{C} \times_1 \mathbf{X} \times \cdots \times_m \mathbf{X}.$$

where  $\mathcal{C} \in \mathbb{R}^{s \times \cdots \times s}$  is core tensor and  $\mathbf{X}$  are factor matrices with orthonormal columns. The condition  $\lambda(\mathcal{T}) \leq 1$  is imposed without loss of generality. The scale between  $\mathcal{T}$  and the  $f$  is undetermined; the tensor  $f(\mathcal{T}) = f'(\mathcal{T}')$  are the same by setting  $\mathcal{T}' = c\mathcal{T}$  and  $f' = f/c$ . We also note that, the rank- $s$  CP tensor is automatically included in  $\mathcal{T}(d, s, m)$ .

Now, we are ready to describe the main model. Let  $\mathcal{Y}$  be the data tensor. We propose the following observation model

$$\mathcal{Y} = f(\mathcal{T}) + \mathcal{E}, \quad \text{for some unknown } \mathcal{T} \in \mathcal{T}(d, s, m) \text{ and } f \in \mathcal{F}(C), \quad (1)$$

where we assume the noise tensor  $\mathcal{E}$  consists of i.i.d. entries with zero-mean and sub-Gaussian parameter  $\sigma^2$ . We call (1) the *single index tensor model*. The name of single index model comes from the observation that

$$\mathbb{E}[\mathcal{Y}(\omega) | \mathcal{X}(\omega)] = f(\langle \mathcal{T}, \mathcal{X}(\omega) \rangle), \quad \text{for all } \omega \in [d]^m,$$

where, for every index  $\omega$ , the predictor  $\mathcal{X}(\omega) \in (\mathbb{R}^d)^{\otimes m}$  is a dummy-variable tensor with 1 at the  $\omega$ -th position, and zero everywhere else. Note that the signal tensor  $f(\mathcal{T})$  is often high rank. Our goal is to address the following two questions:

- What are the *statistical* and *computational* limits for signal estimation in single index model?
- Are there any intrinsic distinctions for matrices  $m = 2$  vs. tensors  $m \geq 3$  for high-rank estimation based on model (1)?

### 3. Identifiability

$f$  and  $\mathcal{T}$  are not identifiable separately. However,  $f(\mathcal{T})$  is identifiable (why??).

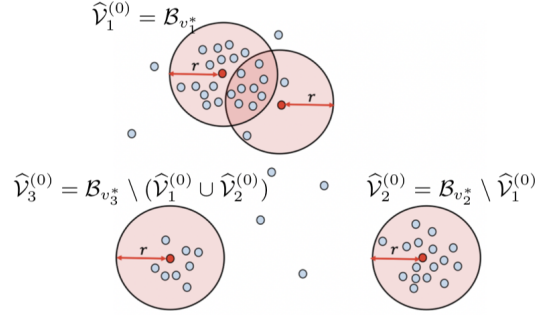
### 4. Smooth single index models are of log-rank

Let  $\mathcal{T} \in \mathcal{T}(d, s, m)$ , and  $\mathcal{T}^{\circ \ell}$  be the tensor of the same size, with entrywise polynomial transformation  $a \rightarrow a^\ell$ . We have the following rank bound.

**Proposition 1 (Monomial tensor)** *For every tensor  $\mathcal{T} \in \mathcal{T}(d, s, m)$  and every natural number  $\ell \in \mathbb{N}_+$*

$$\text{Trank}(\mathcal{T}^{\circ \ell}) \leq \binom{\ell + s - 1}{s - 1}.$$

This proposition shows the polynomial rank growth with respect to  $\ell$ . The exponent depends on the original rank  $s$ , which is assumed small. When  $s = 1$ , then  $\text{Trank}(\mathcal{T}^{\circ \ell}) = 1$  for all  $\ell \in \mathbb{N}_+$ . The bound is nontrivial when  $\ell \leq d$ . In fact, later we will set  $\ell \lesssim \log d$ .



**Proposition 2 (Uniform approximation for single index tensor)** Consider  $\mathcal{T} \in \mathcal{T}(d, s, m)$  and  $f \in \mathcal{F}(C)$ . There exists a set of basis tensors  $\{\mathcal{B}_{\ell, k} : (k, \ell) \in [d] \times \mathbb{N}\}$ , such that, for every number of pieces  $k \in [d]$  and degree  $\ell \in \mathbb{N}$ , we have

$$\|f(\mathcal{T}) - \mathcal{B}_{k, \ell}\|_{\infty} \leq C \left( \frac{k}{sm} \right)^{-(\ell+1)} \quad \text{and} \quad \text{Trank}(\mathcal{B}_{k, \ell}) \leq k^s \ell^s.$$

The following is the key property of single index matrix by taking  $k \gtrsim \Omega(sm)$  and  $\ell = r^{1/s} k^{-1}$  in Proposition 2.

**Theorem 4.1 (Dimension-free approximation error)** Consider a single index tensor  $\Theta = f(\mathcal{T})$ , where  $\mathcal{T} \in \mathcal{T}(d, s, m)$  and  $f \in \mathcal{F}(C)$ . For every rank  $r \in \mathbb{N}_+$ , we have

$$\frac{1}{d^m} \|\Theta - \text{Proj}_r(\Theta)\|_F^2 \lesssim C^2 \exp\left(-\frac{r^{1/s}}{sm}\right).$$

The bound is uniform over  $\mathcal{T} \in \mathcal{T}(d, s, m)$  and  $f \in \mathcal{F}(C)$ .

**Corollary 1 (Nice smooth tensor model is of log rank)** For any fixed  $\varepsilon > 0$ , we have

$$\text{Trank}_{\varepsilon}(\Theta) := \min\{\text{Trank}(\mathcal{A}) : \|\mathcal{A} - \Theta\|_{\infty} \leq \varepsilon\} \lesssim \log^s d.$$

**Corollary 2 (Block tensor model with diverging number of blocks)** Let  $\Theta$  be a block tensor with  $s$  blocks on each mode. The  $\varepsilon$ -rank of  $\Theta$  is no more than  $r = \log d$ .

When  $f(x) = x$ ,  $\text{Trank}(\Theta) = o(1)$ . When oscillating  $f(x) = \sin(x)$ , ....? When monotonic  $f(x) = \frac{\exp(-x)}{1+\exp(-x)}$  at bounded regime  $[a, b]$ ;  $f$  is quantile transformation??

Why do we define domain on  $[-1, 1]$ ??

Marginal quantile transformed tensor estimation. Low rank. Order-preserving Transformed. Analogy of distance matrix???

Connection to sign-representable tensor. Require  $f$  to be monotonic. Here smooth  $f$  suffices...

## 5. Single index model estimation for matrices

### 5.1. Spectral method and matrix lasso are both nearly minimax

We now consider two estimations. The first one is convex estimator

$$\hat{\Theta} = \arg \min_{\Theta} \left\{ \|\mathbf{Y} - \Theta\|_F^2 + \sqrt{d}\sigma \|\Theta\|_* \right\},$$

and the other is non-convex estimator  $r = sm \log^s d$ . The closed-form solution is

$$\hat{\Theta} = \sum_{i=1}^d (\lambda_i(\mathbf{Y}) - \sqrt{d}\sigma)_+ \mathbf{u}_i^T \mathbf{v}_j.$$

### 5.2. Nonconvex

**Theorem 5.1 (Polynomial estimation of single index matrix)** *Let  $\Theta = f(\mathbf{T})$  with  $\mathbf{T} \in \mathcal{T}(d, s, 2)$  and  $f \in \mathcal{F}(C)$ .*

$$\|\hat{\Theta} - \Theta\|_F \lesssim d^{-1} \log^s(d).$$

**Theorem 5.2 (Minimax for matrices)** *Assume  $s \asymp \mathcal{O}(1)$ . Denote the class  $\mathcal{F} \circ \mathcal{T} = \{\Theta : \Theta = f(\mathbf{T}) \text{ for some } (\mathbf{T}, f) \in \mathcal{T}(d, s, 2) \times \mathcal{F}(C)\}$ . Then, when we have*

$$\inf_{\Theta} \sup_{\Theta \in \mathcal{F} \circ \mathcal{T}} \mathbb{P} \left( \mathcal{R}(\Theta, \hat{\Theta}) \gtrsim d^{-1} \right) \geq 0.2.$$

**Proof** For every  $\alpha_1 \leq \alpha_2 \in (0, \infty)$

$$\mathcal{P}(d, \alpha_1, s, L) \supset \mathcal{P}(d, \alpha_2, s, L) \supset \cdots \mathcal{P}(d, \infty, s, L) \supset \{\Theta : \text{Trank}(\Theta) \leq s\}.$$

■

Open question: Can we show sharper lower bound for  $\alpha \neq \infty$ ? Graphon???

## 6. Tensor estimation

### 6.1. Non-convex double spectral algorithm

**Theorem 6.1 (Polynomial algorithm)** *low-rank approximation is substantially better.*

$$\mathcal{R}(\hat{\Theta}, \Theta) \leq \begin{cases} r^m + d^{m/2}r + d^m r^{-2\alpha} & r^{-2\alpha} \leq d^{-(m-1)}, \\ \min(d^{-\frac{m}{2} + \frac{(m-\alpha-1)m}{2\alpha}}, 1) & \alpha \leq \frac{(m-1)^2}{m}, r \asymp \min(d^{(m-1)/2\alpha}, d) \\ d^{-\frac{m}{2} + \frac{m-1}{2\alpha}} & \alpha \geq \frac{(m-1)^2}{m}, r \asymp d^{\frac{m-1}{2\alpha}} \leq d^{\frac{m}{2(m-1)}} \\ d^{-\frac{m}{2}} \log d & \alpha = \infty, r \asymp \log d \end{cases} \quad (2)$$

Block approximation (bottleneck: truncate at  $\alpha = 1$ ):

$$\mathcal{R}(\hat{\Theta}, \Theta) \leq \begin{cases} d^{-\frac{2\alpha m}{m+2\alpha}} & \alpha \leq 1 \\ d^{-\frac{2m}{m+2}} & \alpha \geq 1 \end{cases} \quad (3)$$

**Theorem 6.2 (Minimax optimality for tensors)** *Order- $m$  with known design (bottleneck: truncate at  $d^{-(m-1)}$  with unknown design)*

$$\mathcal{R}(\hat{\Theta}, \Theta) \leq \begin{cases} d^{-\frac{2\alpha m}{m+2\alpha}} & \alpha \leq \frac{m(m-1)}{2} \\ d^{-(m-1)} \log d & \alpha \geq \frac{m(m-1)}{2} \end{cases} \quad (4)$$

## 6.2. Statistical and Computational Trade-offs

### Acknowledgments

Funding here.

Can we apply tensor algorithm to matrices??

## 7. Nice tensor model has $\log d$ rank

### Appendix A. Proof of Proposition 1

**Proof** By definition  $\mathcal{T} \in \mathcal{T}(d, s, m)$ ,  $\text{Trank}(\text{Mat}(\mathcal{T})) \leq r$ . Therefore, there exists matrix SVD such that

$$\text{Mat}(\mathcal{T}) = \sum_{i \in [s]} \lambda_i \mathbf{a}_i \otimes \mathbf{b}_i,$$

where  $\mathbf{a}_i \in \mathbb{R}^d$ ,  $\mathbf{b}_i \in \mathbb{R}^{d^{m-1}}$ , and  $\lambda_1 \geq \dots \geq \lambda_s \geq 0$ . By definition

$$\text{Mat}(\mathcal{T}^{\circ \ell}) = [\text{Mat}(\mathcal{T})]^{\circ \ell} = \left( \sum_{i \in [s]} \lambda_i \mathbf{a}_i \otimes \mathbf{b}_i \right)^{\circ \ell} \quad (5)$$

$$= \sum_{\substack{\kappa_1 + \dots + \kappa_s = \ell, \\ (\kappa_1, \dots, \kappa_s) \in \mathbb{N}_+^s}} \lambda_1^{\kappa_1} \dots \lambda_s^{\kappa_s} (\mathbf{a}_1^{\circ \kappa_1} \circ \dots \circ \mathbf{a}_s^{\circ \kappa_s}) \otimes (\mathbf{b}_1^{\circ \kappa_1} \circ \dots \circ \mathbf{b}_s^{\circ \kappa_s}). \quad (6)$$

Here  $(\mathbf{a}_1^{\circ \kappa_1} \circ \dots \circ \mathbf{a}_s^{\circ \kappa_s}) \in \mathbb{R}^d$  and  $(\mathbf{b}_1^{\circ \kappa_1} \circ \dots \circ \mathbf{b}_s^{\circ \kappa_s}) \in \mathbb{R}^{d^{m-1}}$ . Now notice that by counting argument,

$$\#\{(\kappa_1, \dots, \kappa_s) \in \mathbb{N}_+^s : \kappa_1 + \dots + \kappa_s = \ell\} = \binom{\ell + s - 1}{s - 1}.$$

Therefore, the summation (5) consists of no more than  $\binom{\ell + s - 1}{s - 1}$  rank-1 terms. We conclude that

$$\text{Trank}(\mathcal{T}^{\circ \ell}) = \text{rank}(\text{Mat}(\mathcal{T}^{\circ \ell})) \leq \binom{\ell + s - 1}{s - 1}.$$

■

### Appendix B. Proof of Lemma 2

**Lemma 1** *we have*

$$\|f(\mathcal{T}) - \mathcal{B}_{k, \ell}\|_{\infty} \leq \frac{\max_{|\eta| \leq 1} |f^{(\ell+1)}(\eta)|}{(\ell + 1)!} \|(\mathcal{T} - \mathcal{O})^{\circ \ell+1}\|_{\infty} \leq \left( \frac{k}{sm} \right)^{-(\ell+1)}.$$

Note that the covering number of  $s$ -dimensional bounded set  $\mathcal{X}$  is  $\mathcal{N}(1/k, \mathcal{X}, \|\cdot\|_\infty) \leq k^s$ . Let  $\mathcal{E}_k$  denote the corresponding covering set. Then,  $\mathcal{E}_k$  satisfies

1.  $|\mathcal{E}_k| \leq k^s$ ;
2. For every  $\Delta \in \mathcal{E}_k$ , we have

$$\max_{\mathbf{x}_i, \mathbf{x}_j \in \Delta} \|\mathbf{x}_i - \mathbf{x}_j\|_\infty \lesssim \frac{1}{k}.$$

3.  $\Delta \cap \Delta' = \emptyset$  for all  $\Delta \neq \Delta' \in \mathcal{E}_k$ .

We label the center of the covering set by  $\{\mathbf{o}_1, \dots, \mathbf{o}_{k^s}\} \subset \{\mathbf{x}_1, \dots, \mathbf{x}_d\}$ . We use  $z: [d] \rightarrow [k^s]$  to denote the membership of rows of  $\mathbf{X}$

$$z(i) = \arg \min_{j \in [k^s]} \|\mathbf{x}_i - \mathbf{o}_j\|_\infty, \quad \text{such that} \quad \|\mathbf{x}_i - \mathbf{o}_{z(i)}\|_\infty \leq \frac{1}{k}.$$

We define a block matrix  $\mathcal{O} \in (\mathbb{R}^d)^{\otimes m}$  with entries

$$\mathcal{O}(i_1, \dots, i_m) = \mathcal{C} \times_1 \mathbf{o}_{z(i_1)} \times_2 \cdots \times_m \mathbf{o}_{z(i_m)}$$

Therefore

$$\|\mathcal{T} - \mathcal{O}\|_\infty \leq \max_{(i_1, \dots, i_m) \in [d]^m} |\mathcal{C} \times_1 \mathbf{x}_{i_1} \times_2 \cdots \times_m \mathbf{x}_{i_m} - \mathcal{C} \times_1 \mathbf{o}_{z(i_1)} \times_2 \cdots \times_m \mathbf{o}_{z(i_m)}| \lesssim mk^{-1}s^{1/2}.$$

Define

$$\mathcal{B}_{k,\ell} := f^{(1)}(\mathcal{O}) \circ (\mathcal{T} - \mathcal{O}) + \frac{f^{(2)}(\mathcal{O})}{2} \circ (\mathcal{T} - \mathcal{O})^{\circ 2} + \cdots \frac{f^{(\ell)}(\mathcal{O})}{\ell!} \circ (\mathcal{T} - \mathcal{O})^{\circ \ell}.$$

The membership partition  $[d]^m$  into  $k^{sm}$  blocks. We use  $[d]^m = \cup_{n=1}^{k^{sm}} \Delta_n$  to denote these blocks. Within each block, the tensor  $\mathcal{O}$  takes the same value,

$$\mathcal{B}_{k,\ell}(\omega) = \underbrace{(a_{0,\Delta} + a_{1,\Delta}\mathcal{T} + a_{2,\Delta}\mathcal{T}^{\circ 2} + \cdots a_{\ell,\Delta}\mathcal{T}^{\circ \ell})}_{:=\mathcal{I}} \mathbb{1}\{\omega \in \Delta\}. \quad (7)$$

Because there are  $k^s$  blocks along mode 1, we conclude that

$$\text{Trank}(\mathcal{B}_{k,\ell}) \leq k^s \sum_{n=1}^{\ell} \binom{n+s-1}{s-1} \leq k^s(\ell+s-1)^s.$$

### Appendix C. My Proof of Proposition 1

**Lemma 2 (Polynomial and Unfolding)** Consider a function  $f: (\mathbb{R}^s)^{\otimes m} \rightarrow \mathbb{R}$ . Suppose that, under  $(\mathbb{R}^s)^{\otimes m} \cong \mathbb{R}^{sm}$ , the function  $f$  is a degree- $\ell$  polynomial of  $(sm)$ -variables; i.e.

$$g(\mathbf{x}_1, \dots, \mathbf{x}_m) = \sum_{\kappa_1 + \dots + \kappa_{sm} \leq \ell} c(\kappa_1, \dots, \kappa_{sm}) x_{11}^{\kappa_1} \cdots x_{sm}^{\kappa_{sm}},$$

where  $x_{ij}$  denotes the  $i$ -th element in  $\mathbf{x}_j$ , for  $(i, j) \in [s] \times [m]$ . Then the dimensional- $d$  order- $m$  tensor  $\mathbf{G}$  generated by  $f$  and  $(\mathbf{a}_i)_{i \in [d]}$  admits

$$\mathbf{G}(i_1, \dots, i_m) = g(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_m}), \quad \text{for all } (i_1, \dots, i_m) \in [d]^m.$$

Then

$$\text{Trank}(\mathbf{G}) \leq \sum_{i=0}^{\ell} \binom{i+s-1}{s-1}.$$

**Proof**

$$\text{Unfold}(\mathbf{G}) = \begin{bmatrix} \Phi(\mathbf{a}_1) \\ \vdots \\ \Phi(\mathbf{a}_d) \end{bmatrix} \mathbf{B}^T = \begin{bmatrix} 1 & a_{11} & \cdots & a_{s1} & a_{11}^2 & a_{11}a_{s1} & \cdots & a_{11}a_{21}a_{s1}^{\ell-2} \cdots & a_{s1}^{\ell} \\ \vdots & & & & & & & & \\ 1 & a_{1d} & \cdots & a_{sd} & a_{1d}^2 & a_{1d}a_{1d} & \cdots & a_{1d}a_{2d}a_{sd}^{\ell-2} \cdots & a_{sd}^{\ell} \end{bmatrix} \mathbf{B}^T$$

■

**Lemma 3 (Uniform approximation)** Consider  $\Theta \in \mathcal{P}(d, s, L)$ . There exists a sequence of tensor  $\{\mathcal{A}_{\ell, k} : (k, \ell) \in [d] \times \mathbb{N}\}$  such that, for every number of pieces  $k \in [d]$  and degree  $\ell \in \mathbb{N}$ , we have

$$\|\Theta - \mathcal{A}_{k, \ell}\|_{\infty} \leq Lk^{-\ell} \ell^{ms}, \quad \text{and} \quad \text{Trank}(\mathcal{A}_{k, \ell}) \leq k^s \ell^s.$$

**Proof** Based on the definition of infinitely smooth function

$$|\Theta - \mathcal{A}_{k, \ell}| \leq |f(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_m}) - g(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_m})| \tag{8}$$

$$\leq L \left(\frac{k}{a}\right)^{-\ell} \binom{\ell + ms - 1}{ms - 1} \tag{9}$$

$$\leq L \left(\frac{k}{a}\right)^{-\ell} \ell^{ms} \tag{10}$$

■

**Corollary 3 (Nice tensor model is of  $\log d$  rank)** For any fixed  $\varepsilon > 0$ , we have

$$\text{Trank}_{\varepsilon}(\Theta) := \min\{\text{Trank}(\mathcal{A}) : \|\mathcal{A} - \Theta\|_{\infty} \leq \varepsilon\} \lesssim \log d.$$

Where does  $d$  come from??

**Theorem C.1 (Dimension-free approximation)** *Let  $\Theta \in \mathcal{P}(d, s, L)$ . For every rank  $r \in \mathbb{N}_+$ ,*

$$\frac{1}{d^m} \|\Theta, \text{Proj}_r(\Theta)\|_F^2 \lesssim L^2 \exp(-c_0 r^{1/s}).$$

**Proof** Let  $k = \Omega(2a)$  and  $\ell = r^{1/s} k^{-1}$ . Then

$$\|\Theta - \mathcal{A}_{k,\ell}\|_\infty \leq L 2^{-\ell} \ell^{ms} \lesssim L 2^{-\ell} \asymp L \exp(-r^{1/s})$$

■

## Appendix D. My Proof of Theorem 2

This is a complete version of a proof sketched in the main text.

**Lemma 4 (Tensor polynomial approximation)** *Suppose that the function  $f: [0, 1]^m \rightarrow \mathbb{R}$  generating the signal tensor  $\Theta$  is  $\alpha$ -Hölder smooth with  $\alpha \in (0, \infty)$ . Then, for any  $k \geq d$ , there exists a  $k$ -membership function  $z: [d] \rightarrow [k]$  and a polynomial function such that*

$$\Theta(i_1, \dots, i_K) - P_{[\alpha]}(\omega - \mathcal{S}(z(i_1), \dots, z(i_K))) \leq \frac{m^2}{k^\alpha}$$

,

$$M(\mathbf{x} - \mathbf{x}_0) = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 - x_0 & x_2 - x_0 & \cdots & x_d - x_0 \\ \vdots & \vdots & \vdots & \\ (x_1 - x_0)^{[\alpha]} & (x_2 - x_0)^{[\alpha]} & \cdots & (x_d - x_0)^{[\alpha]} \end{bmatrix}$$



Let  $z: [d] \rightarrow [k]$  denote the clustering function. For notational convenience, we use the shorthand  $\mathbf{z}(i_1, \dots, i_m) := (z(i_1), \dots, z(i_m))^T \in [k]^m$  to denote the block membership by applying  $z$  to each of the  $m$  modes. We make the convention that blockwise constant tensor is of degree 1 (not zero, for notational convenience).

1. blockwise degree-1 (constant) tensor

$$\mathcal{L}(0, k) = \left\{ \mathcal{B} \in (\mathbb{R}^d)^{\otimes m} : \mathcal{B}(i_1, \dots, i_m) = \mathcal{C}(z(i_1), \dots, z(i_m)) \text{ for some tensor } \mathcal{C} \in (\mathbb{R}^k)^{\otimes m} \right\} \quad (11)$$

$$= \left\{ \mathcal{B} \in (\mathbb{R}^d)^{\otimes m} : \Theta(\omega) = \sum_{\Delta \in [k]^m} c_{\Delta} \mathbb{1}\{\mathbf{z}(\omega) = \Delta\} \right\} \quad (12)$$

$$\cong \{ \mathcal{C} \in \mathbb{R}^{k^m}, z: [d] \rightarrow [k] \} \quad (13)$$

where, for every cluster  $\Delta \in [k]^m$ , the coefficient  $c_{\Delta} \in \mathbb{R}$  represents the mean within the cluster. We use  $\mathcal{C} = \{c_{\Delta} : \Delta \in [k]^m\}$  to collect all unknown coefficients.

2. blockwise degree-2 linear tensor

$$\mathcal{L}(1, k) = \left\{ \mathcal{B} \in (\mathbb{R}^d)^{\otimes m} : \mathcal{B}(\omega) = \sum_{\Delta \in [k]^m} [c_{\Delta} + \langle \beta_{\Delta}, \omega \rangle] \mathbb{1}\{\mathbf{z}(\omega) = \Delta\} \right\} \quad (14)$$

$$\cong \{ \mathcal{C} \in \mathbb{R}^{(1+m) \times k^m}, z: [d] \rightarrow [k] \} \quad (15)$$

where, for every cluster  $\Delta \in [k]^m$ , the coefficient  $(c_{\Delta}, \beta_{\Delta}) \in \mathbb{R} \times \mathbb{R}^d$  the mean and slope within the cluster. We use  $\mathcal{C} = \{(c_{\Delta}, \beta_{\Delta}) : \Delta \in [k]^m\} \cong \mathbb{R}^{(1+m) \times k^m}$  to collect all unknown coefficients.

3. blockwise degree- $(\ell + 1)$  polynomial tensor

$$\mathcal{L}(\ell + 1, k) = \left\{ \mathcal{B} \in (\mathbb{R}^d)^{\otimes m} : \mathcal{B}(\omega) = \sum_{\Delta \in [k]^m} \text{Poly}_{\ell, \Delta}(\omega) \mathbb{1}\{\mathbf{z}(\omega) = \Delta\} \right\} \quad (16)$$

$$\cong \{ \mathcal{C} \in \mathbb{R}^{(\ell+m) \times k^m}, z: [d] \rightarrow [k] \} \quad (17)$$

where, for every cluster  $\Delta \in [k]^m$ , the coefficient  $(c_{\Delta}, \beta_{\Delta}) \in \mathbb{R} \times \mathbb{R}^d$  the mean and slope within the cluster. Note that the polynomial function  $\text{Poly}_{\ell, \Delta}(\cdot)$  has at most  $(\ell + m)^{\ell}$  unknown coefficients. We use  $\mathcal{C} \subset \mathbb{R}^{(\ell+m) \times k^m}$  to collect the unknown coefficients in the blockwise degree- $(\ell + 1)$  polynomial tensor family.

Model:

$$\Theta(i_1, \dots, i_m) = f\left(\frac{\pi(i_1)}{d}, \dots, \frac{\pi(i_m)}{d}\right), \quad \text{for all } i_1 < \dots < i_d. \quad (18)$$

**Theorem D.1** (*Blockwise polynomial tensor approximation*) Suppose the function  $f : [0, 1]^m \rightarrow \mathbb{R}$  generating the signal tensor  $\Theta$  as in model (18) is  $\alpha$ -Hölder smooth with  $\alpha \in (0, \infty)$ . Then, for every block size  $k \leq d$  and degree  $\ell \in \mathbb{N}_+$ , we have the approximation error

$$\inf_{\mathcal{B} \in \mathcal{L}(\ell, k)} \frac{1}{d^m} \|\Theta - \mathcal{B}\|_F^2 \lesssim \frac{m^2}{k^{2 \min(\alpha, \ell)}}.$$

We propose the least-square estimate

$$\hat{\Theta}^{\text{LSE}} = \arg \min_{\Theta \in \mathcal{L}(\ell, k)} \|\mathcal{Y} - \Theta\|_F.$$

Rigorously, the least-square estimate  $\hat{\Theta}^{\text{LSE}}$  is a family of estimates with tuning parameters  $(\ell, k)$ ; we suppress the dependence on  $(\ell, k)$  when no confusion arises.

**Theorem D.2 (Least-square estimator)** Let  $\hat{\Theta}^{\text{LSE}}$  denote the least-square estimate with degree  $\ell^* = \lfloor \min(\alpha, \frac{m(m-1)}{2}) \rfloor$  with block size  $k^* = \lfloor d^{\frac{m}{m+2\ell^*}} \rfloor$ . Then,  $\hat{\Theta}^{\text{LSE}}$  obeys the error bound

$$\begin{aligned} \frac{1}{d^m} \|\hat{\Theta}^{\text{LSE}} - \Theta\|_F^2 &\lesssim \inf_{(\ell, k) \in \mathbb{N}_+ \times [d]} \left\{ \frac{1}{k^{2 \min(\alpha, \ell)}} + \frac{k^m \ell^m}{d^m} + \frac{\log d}{d^{m-1}} \right\} \\ &\asymp \begin{cases} d^{-\frac{2m\alpha}{m+2\alpha}} & \text{when } \alpha < m(m-1)/2, \\ d^{-(m-1)} \log d & \text{when } \alpha \geq m(m-1)/2. \end{cases} \end{aligned}$$

**Theorem D.3 (Polynomial estimator)** Suppose that the signal tensor  $\Theta$  is generated in model (18) with  $f \in \mathcal{H}(\alpha) \cap \mathcal{M}(\beta)$ . Let  $\hat{\Theta}^{\text{BC}}$  be the estimator in with degree  $\ell^* = \lfloor \min(\alpha, \frac{m(m-1)}{2}) \rfloor$  and block size  $k^* = \lfloor d^{\frac{m}{m+2\ell^*}} \rfloor$ . Then the estimator  $\hat{\Theta}^{\text{BC}}$  satisfies

$$\frac{1}{d^m} \|\hat{\Theta}^{\text{BC}} - \Theta\|_F^2 \lesssim d^{-\beta(m-1)} + \begin{cases} d^{-\frac{2m\alpha}{m+2\alpha}} & \text{when } \alpha < m(m-1)/2, \\ d^{-(m-1)} \log d & \text{when } \alpha \geq m(m-1)/2. \end{cases}$$

with very high probability.

**Remark 1** For order-3 tensor with sufficiently smooth function  $\alpha \geq 3$ , the optimal choice of degree and block sizes is  $(\ell^*, k^*) = (3, d^{1/3})$ .

**Theorem D.4 (Minimax lower bound)** For any given  $\alpha \in (0, \infty)$ , the problem of estimating  $\alpha$ -smooth tensors obeys the minimax lower bound

$$\mathbb{P} \left( \inf_{\hat{\Theta}} \max_{\Theta \in \mathcal{P}(\alpha)} \|\Theta - \hat{\Theta}\|_F^2 \geq d^{-\frac{2m\alpha}{m+2\alpha}} + d^{-(m-1)} \log d \right) > 0.8.$$