

# Smooth tensor estimation

Miaoyan Wang, May 18, 2020

## 1 Model

Key: Cartesian product of piecewise constant function representations. (has full generality...)

Let  $\mathcal{Y} = \llbracket y_{i_1, \dots, i_K} \rrbracket \in \{0, 1\}^{d_1 \times \dots \times d_K}$  be an order- $K$ ,  $(d_1, \dots, d_K)$ -dimensional binary tensor. Let  $\boldsymbol{\xi}^{(k)} = (\xi_1^{(k)}, \dots, \xi_d^{(k)}) \in [0, 1]^d$  be random vectors following (unknown) distributions  $\mathbb{P}^{(k)}$  for all  $k \in [K]$ , and  $\boldsymbol{\xi}^{(k)}$  and  $\boldsymbol{\xi}^{(k')}$  are mutually independent for  $k \neq k' \in [K]$ . Assume that, conditional on  $\{\boldsymbol{\xi}^{(k)}\}$ , the entries of  $\mathcal{Y}$  are independent sub-Gaussian distributed:

$$\mathbb{E}(y_{i_1, \dots, i_K} | \boldsymbol{\xi}) = f(\xi_{1, i_1}, \dots, \xi_{K, i_K}), \quad \text{for all } (i_1, \dots, i_K) \in [d] \times \dots \times [d],$$

where  $f: [0, 1]^K \mapsto [0, 1]$  is an unknown multivariate function belonging to a function class  $f \in \mathcal{F}_\alpha(L)$ . Specifically, the function class is defined as

$$\mathcal{F}_\alpha(L) = \{f: \text{Im}(f) \in [0, 1] \text{ and } \|f\|_{\mathcal{H}_\alpha} \leq L\},$$

where  $\alpha \in (0, 1]$  is the smoothness parameter and  $L > 0$  is the Hölder norm bound for the functions in the class.

Recall that the function Hölder norm  $\|f\|_{\mathcal{H}_\alpha}$  is defined as

$$\|f\|_{\mathcal{H}_\alpha} \stackrel{\text{def}}{=} \max_{|\omega| \leq \lfloor \alpha \rfloor} \sup_{\mathbf{x} \in \mathcal{D}} |\nabla_\omega f(\mathbf{x})| + \max_{|\omega| = \lfloor \alpha \rfloor} \sup_{\mathbf{x} \neq \mathbf{x}' \in \mathcal{D}} \frac{|\nabla_\omega f(\mathbf{x}) - \nabla_\omega f(\mathbf{x}')|}{\|\mathbf{x} - \mathbf{x}'\|_1^{\alpha - \lfloor \alpha \rfloor}},$$

where we have used the short-hand notion

$$\nabla_\omega f(\mathbf{x}) = \frac{\partial^{i_1 + \dots + i_K}}{\partial x_1^{i_1} \dots \partial x_K^{i_K}} f(x_1, \dots, x_K),$$

for multi-indices  $\omega = (i_1, \dots, i_K)$  with  $|\omega| = i_1 + \dots + i_K$ , and  $\mathbf{x} = (x_1, \dots, x_K)$  in the function domain.

## 2 Estimation

Define the objective function

$$F(\mathcal{C}, \{\mathbf{M}_k\}) = \|\mathcal{Y} - \mathcal{C} \times_1 \mathbf{M}_1 \times \dots \times_K \mathbf{M}_K\|_F^2.$$

Denote  $\Theta = \mathcal{C} \times_1 \mathbf{M}_1 \times \cdots \times_K \mathbf{M}_K$  and  $\mathbf{r} = (r_1, \dots, r_K)$ . Then the feasible domain is

$$\mathcal{P}(\mathbf{r}) = \left\{ \Theta \in \mathbb{R}^{d_1 \times \cdots \times d_K} : \Theta = \mathcal{C} \times_1 \mathbf{M}_1 \times \cdots \times_K \mathbf{M}_K, \text{ where } \mathcal{C} \in \mathbb{R}^{r_1 \times \cdots \times r_K} \text{ and } \mathbf{M}_k \in \{0, 1\}^{d_k \times r_k} \text{ are membership matrices for all } k \in [K] \right\}.$$

We propose constrained least-square estimator

$$\hat{\Theta}(\mathbf{r}, M) = \arg \min_{\Theta \in \mathcal{P}(\mathbf{r}), \|\Theta\|_\infty \leq M} F(\Theta | \mathcal{Y}).$$

The function estimator  $\hat{f}: (0, 1]^K \mapsto [0, 1]$  is defined as follows:

$$\hat{f}(x_1, \dots, x_K) \stackrel{\text{def}}{=} \hat{\Theta}(\lceil d_1 x_1 \rceil, \dots, \lceil d_K x_K \rceil), \quad \text{for all } (x_1, \dots, x_K) \in (0, 1]^K.$$

We propose an adaptive smooth (?) estimation,

$$\begin{aligned} \hat{\Theta} &= \arg \min_{\Theta \in \mathcal{P}(\mathbf{r}^*)} F(\Theta), \quad \text{with } \mathbf{r}^* = (r_1^*, \dots, r_K^*), \\ \text{and } r_k^* &= \lceil d_k^{1/(\alpha \wedge 1 + 1)} \rceil \text{ for all } k \in [K]. \end{aligned}$$

**Theorem 2.1.** Consider a function class  $\mathcal{F}_\alpha(R)$  with  $\alpha > 0$  and  $M > 0$ . We have

$$\sup_{f \in \mathcal{F}_\alpha(R)} \sup_{\boldsymbol{\xi}^{(k)} \sim \mathbb{P}^{(k)}, k \in [K]} \frac{1}{d^K} \mathbb{E} \left( \|\hat{\Theta} - f(\xi_{i_1}^{(1)}, \dots, \xi_{i_K}^{(K)})\|_F^2 \right) \leq C \left( d^{-K\alpha/(\alpha+1)} + \frac{\log d}{d^{K-1}} \right),$$

where the constant  $C > 0$  depends only on  $L$ , and the expectation is taken jointly over  $\mathcal{Y}$ ,  $\{\boldsymbol{\xi}^{(k)}\}$  for all  $k \in [K]$ .

Phase transition at  $\alpha = 1$  only for  $K \geq 3$ ??

## 2.1 Non-parametric tensor model

Special cases: low-rank model, additive-multiplicative model (in the literature of non-parametric...)

Let  $\mathcal{Y} = \llbracket Y_\omega \rrbracket$  be a binary tensor, where  $\omega = (i_1, \dots, i_K)$  is a  $K$ -tuple index. We propose the following conditionally-independent tensor model:

$$\begin{aligned} Y_\omega | \boldsymbol{\xi}_\omega &\sim \text{Bernoulli}(\theta_\omega), \\ \theta_\omega &= f(\xi_{i_1}^{(1)}, \dots, \xi_{i_K}^{(K)}), \text{ for all } \omega = (i_1, \dots, i_K) \in [d_1] \times \cdots \times [d_K], \end{aligned}$$

where  $f: [0, 1]^K \mapsto [0, 1]$  is a multivariate function of interest. We use  $\boldsymbol{\xi}_\omega \equiv (\xi_{i_1}^{(1)}, \dots, \xi_{i_K}^{(K)})$  to denote the latent design variables at position  $\omega = (i_1, \dots, i_K)$ . Furthermore, we assume that the collection

of latent variables at the  $k$ -th coordinate,  $(\xi_1^{(k)}, \dots, \xi_{d_K}^{(k)})$ , follow a  $d_k$ -dimensional distribution  $\mathbb{P}^{(k)}$ , and the distributions  $\mathbb{P}^{(k)}$  and  $\mathbb{P}^{(k')}$  are mutually independent for  $k \neq k' \in [K]$ .

$$\begin{aligned} \boldsymbol{\xi} &: [d_1] \times \dots \times [d_K] \mapsto [0, 1]^K \\ \omega \equiv (i_1, \dots, i_K) &\mapsto \boldsymbol{\xi}_\omega \equiv (\xi_{i_1}^{(1)}, \dots, \xi_{i_K}^{(K)}) \sim \mathbb{P}^{(1)} \times \dots \times \mathbb{P}^{(K)}. \end{aligned}$$

$$\begin{aligned} f &: [0, 1]^K \mapsto [0, 1] \\ \boldsymbol{\xi} &\mapsto f(\boldsymbol{\xi}). \end{aligned}$$

**Lemma 1** (Connection between tensor block model and non-parametric tensor model). *Let  $f \in \mathcal{F}_\alpha(L)$  be a target function and  $\boldsymbol{\xi} \sim \mathbb{P}_\xi$  be realized latent variables. Let  $f(\boldsymbol{\xi}_\omega) \in \mathbb{R}$  denote the tensor entry indexed by  $\omega \in [d_1] \times \dots \times [d_K]$ , and let  $f(\boldsymbol{\xi}) = \llbracket f(\boldsymbol{\xi}_\omega) \rrbracket \in \mathbb{R}^{d_1 \times \dots \times d_K}$  denote the non-parametric tensor parameter. Then for each  $\mathbf{r} \in [d_1] \times \dots \times [d_K]$ , there exists a parameter  $\Theta = \llbracket \theta_\omega \rrbracket \in \mathcal{P}(\mathbf{r})$  such that*

$$\text{Loss}(f(\boldsymbol{\xi}), \Theta) \leq L^2 \left( \sum_k \frac{1}{r_k} \right)^{\alpha \wedge 1}.$$

**Remark 1.** The (deterministic) bound holds uniformly over  $\mathbb{P}_\xi$  and  $\mathcal{F}_\alpha(L)$ .

*Proof.* The proof is constructive. We partition the interval  $[0, 1]$  into  $r_k$  equal-sized intervals at each of the  $K$  modes. Denote the grid of intervals  $\mathcal{I}_{s_1, \dots, s_K} = \left( \frac{s_1-1}{r_1}, \frac{s_1}{r_1} \right] \times \dots \times \left( \frac{s_K-1}{r_K}, \frac{s_K}{r_K} \right]$  for all  $(s_1, \dots, s_K) \in [r_1] \times \dots \times [r_K]$ . We use the notation  $\boldsymbol{\xi}_{\omega'} \sim \boldsymbol{\xi}_\omega$  if and only if  $\boldsymbol{\xi}_{\omega'}$  and  $\boldsymbol{\xi}_\omega$  belong to the same grid.

The parameter  $\Theta = \llbracket \theta_\omega \rrbracket$  is constructed as follows. For each  $\omega \in [d_1] \times \dots \times [d_K]$ , we set  $\theta_\omega$  to be the average of  $f(\boldsymbol{\xi}_{\omega'})$  over the index set for which  $\boldsymbol{\xi}_{\omega'}$  in the same grid as  $\boldsymbol{\xi}_\omega$ . Specifically, define

$$\theta_\omega = \frac{1}{c_\omega} \sum_{\omega'} f(\boldsymbol{\xi}_{\omega'}) \mathbb{1}\{\boldsymbol{\xi}_{\omega'} \sim \boldsymbol{\xi}_\omega\},$$

where  $c_\omega = \sum_{\omega'} \mathbb{1}\{\boldsymbol{\xi}_{\omega'} \sim \boldsymbol{\xi}_\omega\}$  is the number of tensor entries in the index set  $\{\omega': \boldsymbol{\xi}_{\omega'} \sim \boldsymbol{\xi}_\omega\}$ . The construction implies that  $\Theta = \llbracket \theta_\omega \rrbracket$  takes constant value in the index set  $\{\omega: \boldsymbol{\xi}_\omega \in \mathcal{I}_{s_1, \dots, s_K}\}$  for any given  $(s_1, \dots, s_K) \in [r_1] \times \dots \times [r_K]$ . Therefore,  $\Theta$  is a block tensor with at most  $r$  blocks on each of the  $K$  modes; that is,  $\Theta \in \mathcal{P}(\mathbf{r})$ .

We will show that the defined  $\Theta$  is close to  $f(\boldsymbol{\xi}_\omega)$  in the square distance. Specifically,

$$|f(\boldsymbol{\xi}_\omega) - \theta_\omega| = \left| f(\boldsymbol{\xi}_\omega) - \frac{1}{c_\omega} \sum_{\omega'} f(\boldsymbol{\xi}_{\omega'}) \mathbb{1}\{\boldsymbol{\xi}_{\omega'} \sim \boldsymbol{\xi}_\omega\} \right|$$

$$\begin{aligned}
&\leq \frac{1}{c_\omega} \sum_{\{\omega': \xi_{\omega'} \sim \xi_\omega\}} |f(\xi_\omega) - f(\xi_{\omega'})| \\
&\leq \frac{1}{c_\omega} \sum_{\{\omega': \xi_{\omega'} \sim \xi_\omega\}} L \|\xi_\omega - \xi_{\omega'}\|^{\alpha \wedge 1} \\
&\leq \frac{1}{c_\omega} \sum_{\{\omega': \xi_{\omega'} \sim \xi_\omega\}} L \left(\frac{K}{r}\right)^{\alpha \wedge 1} \\
&\leq L K r^{-(\alpha \wedge 1)},
\end{aligned}$$

where the third line comes from the Hölder condition for  $f$ , the fourth line comes from the fact that  $\|\xi_\omega - \xi_{\omega'}\| \leq \frac{K}{r}$  for  $\xi_\omega \sim \xi_{\omega'}$ . Summing over all  $\omega \in [d_1] \times \cdots \times [d_K]$  gives the conclusion

$$\frac{1}{\prod_k d_k} \sum_{\omega} (f(\xi_\omega) - \theta_\omega)^2 \leq L^2 K^2 r^{-2(\alpha \wedge 1)}.$$

□

**Theorem 2.2** (MSE for tensor block model). *Let  $\mathcal{Y} \in \{0, 1\}^{d_1 \times \cdots \times d_K}$  be a binary tensor generated from tensor block model, i.e.*

$$\mathcal{Y} | \Theta \sim \text{Bernoulli}(\Theta), \quad \text{where } \Theta \in \mathcal{P}(\mathbf{r}).$$

*Then for constant  $C' > 0$ , there exists a constant  $C > 0$  such that*

$$\frac{1}{\prod_k d_k} \|\Theta - \hat{\Theta}\|_F^2 \leq C \left( \prod_k \frac{r_k}{d_k} + \frac{\sum_k d_k \log r_k}{\prod_k d_k} \right),$$

*with probability at least  $1 - \exp(-C' \sum_k d_k \log r_k)$ , uniformly over  $\Theta \in \mathcal{P}(\mathbf{r})$ .*

Suppose  $r \asymp d^\delta$  for some  $\delta \in [0, 1]$ .

$$\text{MSE}(\Theta, \hat{\Theta}) \asymp \begin{cases} d^{-K}, & \delta = 0 \text{ and } r = 1, \\ d^{-(K-1)}, & \delta = 0 \text{ and } r \geq 2, \\ d^{-(K-1) \log d}, & \delta \in (0, \frac{1}{K}], \\ d^{-K(1-\delta)}, & \delta \in (\frac{1}{K}, 1]. \end{cases}$$

**Theorem 2.3.** *Assume  $r_k = O(r)$  and  $d_k = O(d)$ . Then*

$$\text{MSE}(f(\xi), \hat{\Theta}) \leq r^K + d \log r + L^2 K^2 r^{-2(\alpha \wedge 1)}.$$

*Proof.* Taking  $r = d^\delta$  where  $\delta = \frac{K}{2(\alpha \wedge 1) + K}$  gives

$$\text{MSE}(f(\boldsymbol{\xi}), \hat{\Theta}) \leq \begin{cases} d^{-2\alpha/(2\alpha+1)} + o(1), & K = 1 \\ d^{-2\alpha/(\alpha+1)} + \frac{\log d}{d}, & K = 2 \\ d^{-2\alpha K/(2\alpha+K)} + d^{-2K/(K+2)}, & K \geq 3, \end{cases}$$

Only non-parametric rate appear for  $K \geq 3$ . □

**Remark 2.** The distribution of  $\mathbb{P}_\xi$  is required to cover  $[0, 1]$ . – in order to prove the lower bound..

We restrict to i.i.d. uniform distribution over  $[0, 1]$ . (Erdos Roni)

**Theorem 2.4.** *There exists a constant  $C > 0$  only depending on  $L, \alpha$ , such that*

$$\inf_{\hat{\Theta}} \sup_{f \in \mathcal{F}_\alpha(L)} \sup_{\mathbb{P}_\xi \in \mathcal{P}} \mathbb{E}\{\text{MSE}(\hat{\Theta}, f(\boldsymbol{\xi}))\} \begin{cases} d^{-2\alpha K/(2\alpha+K)}, & 0 < \alpha < 1, \\ d^{-1} \log d, & \alpha \geq 1 \text{ and } K = 2, \\ d^{-2K/(K+2)}, & \alpha \geq 1 \text{ and } K \geq 3. \end{cases}$$

### 3 Assumptions on function families

We consider two families of functions.

**Definition 1** ( $\mathcal{F}(R)$ , piecewise constant functions with at most  $r$  marginal pieces). A function  $f: [0, 1]^K \mapsto [0, 1]$  is called a multivariate  $R$ -step function, if there exists a set of mappings  $\{\phi_k: [0, 1] \mapsto [R]\}_{k \in [K]}$  and an order- $K$  tensor  $\mathcal{C} \in \mathbb{R}^{R \times \dots \times R}$  such that

$$f(x_1, \dots, x_K) = \mathcal{C}(\phi_1(x_1), \dots, \phi_K(x_K)), \quad \text{for all } (x_1, \dots, x_K) \in [0, 1]^K.$$

Denote  $\mathbf{x} = (x_1, \dots, x_K)^T$  and  $\phi(\mathbf{x}) = (\phi_1(x_1), \dots, \phi_K(x_K))^T$ . Then  $f$  can be equivalently written as

$$f(\mathbf{x}) = \sum_{\mathbf{r} \in [R]^K} \mathcal{C}(\mathbf{r}) \mathbb{1}\{\phi(\mathbf{x}) = \mathbf{r}\}, \quad \text{for all } \mathbf{x} \in [0, 1]^K.$$

**Remark 3.** The number of constant pieces need not to be equal along each of the  $K$  modes. We use  $R$  to denote the upper bound for the number of constant pieces over the  $K$  modes.

**Definition 2** ( $\mathcal{F}(\alpha, L)$ , Hölder smooth functions). Let  $\alpha \in (0, 1]$  and  $L > 0$ . A function  $f: [0, 1]^K \mapsto [0, 1]$  is called an  $\alpha$ -Hölder smooth function if

$$|f(\mathbf{x}) - f(\mathbf{x}')| \leq L \|\mathbf{x} - \mathbf{x}'\|_1^\alpha, \quad \text{for all } \mathbf{x}, \mathbf{x}' \in [0, 1]^K.$$

The constant  $\alpha \in (0, 1]$  is called the Hölder smoothness parameter and  $L > 0$  is the Hölder constant.

**Definition 3** (Measure-preserving bijection). Let  $\tau: \mathcal{X} \mapsto \mathcal{X}$  be a bijection and  $U$  be a random variable taking values in  $\mathcal{X}$ . Then,  $\tau$  is called a measure-preserving bijection with respect to  $U$  if for all  $A \subset \mathcal{X}$ .

$$\mathbb{P}(U \in A) = \mathbb{P}(\tau(U) \in A).$$

In particular,  $\tau(U)$  and  $U$  are identically distributed.

Equivalent class ??

$$\{f': f'(\xi') \sim f(\xi) \text{ whenever } \xi' \sim \xi\}$$

**Definition 4** (Weakly isomorphism). Two functions  $f, f': [0, 1]^K \mapsto [0, 1]$  are called weakly isomorphic if there exists a set of measure-preserving bijections  $\{\tau_k: [0, 1] \mapsto [0, 1]\}_{k \in [K]}$  such that

$$f(\xi_{1,1}, \dots, \xi_{K,i_K}) \stackrel{a.s.}{=} f'(\tau_1(\xi_{1,i_1}), \dots, \tau_K(\xi_{K,i_K})), \quad \text{for all } (i_1, \dots, i_K) \in [d_1] \times \dots \times [d_K],$$

where  $\{\xi_{k,i_k}\}$  are i.i.d.  $U[0, 1]$  for all  $i_k \in [d_k]$  and all  $k \in [K]$ . We write  $f \sim f'$  to denote weakly isomorphism. The weakly isomorphism relationship defines a quotient space in  $\mathcal{F}(R)$  or  $\mathcal{F}(\alpha, L)$ .

**Definition 5** (Index permutation). Two tensors  $\Theta, \Theta' \in \mathbb{R}^{d_1 \times \dots \times d_K}$  are called equivalent if there exist a set of index permutations  $\{\sigma_k: [d_k] \mapsto [d_k]\}_{k \in [K]}$ , such that

$$\Theta(i_1, \dots, i_K) = \Theta'(\sigma_1(i_1), \dots, \sigma_K(i_K)), \quad \text{for all } (i_1, \dots, i_K) \in [d_1] \times \dots \times [d_K].$$

Two forms of loss are considered.

**Definition 6** (Integrated loss). Let  $f: [0, 1]^K \mapsto [0, 1]$  be the function of interest. We define the integrated loss

$$\text{Loss}(f, \hat{f}) \stackrel{\text{def}}{=} \inf_{f_{\text{iso}} \sim f} \int_{\mathbf{x} \in [0, 1]^K} |f'(\mathbf{x}) - f_{\text{iso}}(\mathbf{x})|^2 d\mathbf{x},$$

where  $f_{\text{iso}} \sim f$  denotes all (?need to be in the specified function space?) functions that are isomorphic with  $f$ .

**Definition 7** (Discrete loss). Let  $\Theta, \hat{\Theta} \in \mathbb{R}^{d_1 \times \dots \times d_K}$  be two tensors. We define the discrete loss as

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

**Remark 4.** We use the notion  $\text{Loss}(\cdot, \cdot)$  to denote either the discrete loss (for tensors) or the integrated loss (for functions). The meaning should be clear given the contexts.

**Definition 8** (Operations between  $f$  and  $\Theta$ ). Let  $f: [0, 1]^K \mapsto [0, 1]$  be a  $K$ -variate function. Then the  $f$ -induced probability tensor  $\Theta$  is defined as

$$\Theta(i_1, \dots, i_K) = f(\xi_{1,i_1}, \dots, \xi_{K,i_K}), \quad \text{for all } (i_1, \dots, i_K) \in [d_1] \times \dots \times [d_K],$$

where  $\{\xi_{k,i_k}\}$  are i.i.d. Uniform $[0,1]$  for all  $i_k \in [d_k]$  and all  $k \in [K]$ .

Conversely, let  $\Theta \in \mathbb{R}^{d_1 \times \dots \times d_K}$  be an order- $K$  tensor. Then the  $\Theta$ -induced function is defined as

$$\begin{aligned} f(x_1, \dots, x_K) &= \Theta(\lceil d_1 x_1 \rceil, \dots, \lceil d_K x_K \rceil) \\ &= \sum_{(i_1, \dots, i_K)} \Theta(i_1, \dots, i_K) \mathbb{1}\{(i_k - 1) < x_k d_k \leq i_k, \text{ for all } k \in [K]\}, \end{aligned}$$

for all  $(x_1, \dots, x_K) \in [0, 1]^K$ .

**Remark 5.** The two operations are not inverse.  $f \Rightarrow \Theta(f, \xi) \Rightarrow f' \stackrel{\text{def}}{=} f(\Theta(f, \xi))$ , but  $f \neq f'$  (not even isomorphic). Similarly,  $\Theta \Rightarrow f(\Theta, \xi) \Rightarrow \Theta' \stackrel{\text{ref}}{=} \Theta(f(\Theta, \xi))$ , but  $\Theta' \neq \Theta$ . The inverse relationship holds only if  $\xi_{i_k, k} = \frac{i_k}{d_k}$  for all  $i_k \in [d_k]$  and all  $k \in [K]$ .

**Proposition 1** (Connection between  $f$  and  $\Theta$ ). The following properties hold:

1. [From functions to tensors] Let  $\Theta$  denote the induced tensor from  $f$ . Then,  $f \in \mathcal{F}(R) \Rightarrow \Theta \in \mathcal{P}(R, 1)$ . Furthermore,  $\Theta' \sim \Theta \Rightarrow$  there exists  $f'$  such that  $\Theta'$  is the induced tensor from  $f'$ .
2. [From tensors to functions] Let  $f$  denote the induced function from  $\Theta$ . Then  $\Theta \in \mathcal{P}(R, 1) \Rightarrow f \in \mathcal{F}(R)$ .
3. [From tensor pairs to function pairs] Let  $f, f'$  denote the induced functions from  $\Theta$  and  $\Theta'$ , respectively. Then,  $f \sim f' \Leftrightarrow \Theta \sim \Theta'$ . Furthermore,

$$\text{Loss}(f, f') \leq \text{Loss}(\Theta, \Theta').$$

**Problem 1** (Non-parametric estimation with unobserved designs). Let  $f: [0, 1]^K \mapsto [-1, 1]$  be a target function of interest. For each  $k \in [K]$ , we draw a random i.i.d. sample of  $d$  points  $\{x_i^{(k)}\}_{i \in [d]}$  uniformly from  $[0, 1]$ . Write  $\omega = (i_1, \dots, i_K) \in [d]^K$ ,  $\mathbf{x}_\omega = (x_{i_1}^{(1)}, \dots, x_{i_K}^{(K)})^T \in \mathbb{R}^K$ , and  $f_\omega = f(\mathbf{x}_\omega) \in \mathbb{R}$ . The goal is to estimation  $f$  given the set  $\{(\omega, f_\omega)\}$  but without knowing  $\{\mathbf{x}_\omega\}$ !

We estimate  $f$  by using the null design; that is, we partition  $[0, 1]^K$  into  $d^K$  equal-sized grids. Specifically, we assume  $\check{x}_i^{(k)} = \frac{i}{d}$  for all  $i \in [d]$  and  $k \in [K]$ , and write  $\check{\mathbf{x}}_\omega = (\check{x}_{i_1}^{(1)}, \dots, \check{x}_{i_K}^{(K)})^T \in \mathbb{R}^K$ . Then the estimation is based on  $\{(\check{\mathbf{x}}_\omega, f_\omega)\}$ .

Q1: existing results for known  $\{\mathbf{x}_\omega\}$ ? Q2: quantify the difference between  $\{\mathbf{x}_\omega\}$  vs.  $\{\check{\mathbf{x}}_\omega\}$ .

The following results give the estimation error when  $f$  belongs to  $\mathcal{F}(R, 1)$  or  $\mathcal{F}(\alpha, L)$ .

**Proposition 2** (Agnostic error). Consider the chain  $f \Rightarrow \Theta(f, \xi) \Rightarrow \hat{f} \stackrel{\text{def}}{=} f(\Theta(f, \xi))$ .

- If  $f \in \mathcal{F}(R, 1)$ , then

$$\mathbb{E} \left[ \text{Loss}(f, \hat{f}) \right] \leq \frac{CKR}{\sqrt{d}},$$

where the expectation is over  $\xi \sim \mathbb{P}_\xi$ .

- If  $f \in \mathcal{F}(\alpha, L)$ , then

$$\mathbb{E} \left[ \text{Loss}(f, \hat{f}) \right] \leq \frac{CK^2L^2}{d^\alpha},$$

where the expectation is over  $\boldsymbol{\xi} \sim \mathbb{P}_{\boldsymbol{\xi}}$ .

**Remark 6.** Does it matter how to define  $\hat{f}$ ? say kernel, or smooth version?

*Proof.* Note that

$$\Theta(i_1, \dots, i_K) = f(\xi_{i_1}^{(1)}, \dots, \xi_{i_K}^{(K)}), \quad \text{for all } (i_1, \dots, i_K) \in [d]^K.$$

For each  $k \in [K]$ , we sort  $\{\xi_i^{(k)}\}_{i \in [d]}$  in an increasing order and write  $0 \leq \xi_{(1)}^{(k)} \leq \dots \leq \xi_{(d)}^{(k)} \leq 1$ . Define

$$\Theta^*(i_1, \dots, i_K) = f(\xi_{(i_1)}^{(1)}, \dots, \xi_{(i_K)}^{(K)}), \quad \text{for all } (i_1, \dots, i_K) \in [d]^K.$$

Since  $\Theta^* \sim \Theta$ , they induce weakly isomorphic functions. Now consider the function induced by  $\Theta^*$ .

$$\begin{aligned} \hat{f}(\mathbf{x}) &= \sum_{(i_1, \dots, i_K)} \Theta^*(i_1, \dots, i_K) \mathbb{1} \left\{ x_k \in \left( \frac{i_k - 1}{d}, \frac{i_k}{d} \right] \right\} \\ &= \sum_{(i_1, \dots, i_K)} f(\xi_{(i_1)}^{(1)}, \dots, \xi_{(i_K)}^{(K)}) \mathbb{1} \left\{ x_k \in \left( \frac{i_k - 1}{d}, \frac{i_k}{d} \right], \text{ for all } k \in [K] \right\}. \end{aligned}$$

We aim to evaluate the integral over the grid  $I_i = (\frac{i-1}{d}, \frac{i}{d}]$  for  $i \in [d]$ .

$$\text{Loss}(f, \hat{f}) = \sum_{(i_1, \dots, i_K)} \left[ \int_{\mathbf{x} \in I_{i_1} \times \dots \times I_{i_K}} |f(\mathbf{x}) - \hat{f}(\mathbf{x})|^2 d\mathbf{x} \right] \quad (1)$$

We now evaluate the integral over the region  $I_\omega = I_{i_1} \times \dots \times I_{i_K}$ . Define  $\tilde{\mathbf{x}} = (\frac{i_1-1}{d}, \frac{i_1}{d}] \times \dots \times (\frac{i_K-1}{d}, \frac{i_K}{d}] \in I_\omega$ . For any  $\mathbf{x} \in I_\omega$ , we have

$$|f(\mathbf{x}) - \hat{f}(\mathbf{x})| \leq |f(\mathbf{x}) - f(\tilde{\mathbf{x}})| + \left| \hat{f}(\mathbf{x}) - f(\tilde{\mathbf{x}}) \right|.$$

Note that both  $\tilde{\mathbf{x}}, \mathbf{x} \in I_\omega$

$$|f(\mathbf{x}) - f(\tilde{\mathbf{x}})| \leq L|\mathbf{x} - \tilde{\mathbf{x}}|^\alpha \leq LKd^{-\alpha}. \quad (2)$$

Furthermore,

$$\begin{aligned} |\hat{f}(\mathbf{x}) - f(\tilde{\mathbf{x}})| &\leq |f(\xi_{(i_1)}^{(1)}, \dots, \xi_{(i_K)}^{(K)}) - f(\frac{i_1}{d}, \dots, \frac{i_K}{d})| \\ &\leq L \max_k \left| \xi_{(i_k)}^{(k)} - \frac{i_k}{d} \right|^\alpha. \end{aligned} \quad (3)$$

Plugging (2) and (3) to (1), we obtain



$$\begin{aligned} \int_{\mathbf{x} \in [0,1]^K} |f(\mathbf{x}) - \hat{f}(\mathbf{x})|^2 d\mathbf{x} &\leq L^2 K^2 \left\{ d^{-\alpha} + \max_{k, i_k} \left| \xi_{(i_k)}^{(k)} - \frac{i_k}{d} \right|^\alpha \right\}^2 \\ &\leq L^2 K^2 \left\{ d^{-2\alpha} + \max_{k, i_k} \left| \xi_{(i_k)}^{(k)} - \frac{i_k}{d} \right|^{2\alpha} + 2d^{-\alpha} \max_k \left| \xi_{(i_k)}^{(k)} - \frac{i_k}{d} \right|^\alpha \right\}. \end{aligned}$$

By Jensen's inequality,  $f(x) = x^\alpha$  is concave for  $\alpha \in [0, 1)$ , so hence  $\mathbb{E}[f(x)] \leq f(\mathbb{E}(x))$ :

$$\mathbb{E} \left| \xi_{(i_k)}^{(k)} - \frac{i_k}{d} \right|^{2\alpha} \leq \left[ \text{Var} \left( \xi_{(i_k)}^{(k)} \right) \right]^\alpha \leq C d^{-\alpha}.$$

The last line comes from the fact that the order statistics of the uniform distribution belongs to Beta distribution,  $\xi_{(i)} \sim \text{Beta}(i, d+1-i)$  for all  $i \in [d]$ . Then  $\text{Var}(\xi_{(i)}) \leq \sqrt{\mathbb{E}(X^2)}$ , we have

$$\mathbb{E} \left| \xi_{(i_k)}^{(k)} - \frac{i_k}{d} \right|^\alpha \leq C d^{-\alpha/2}.$$

Therefore,

$$\text{Loss}(f, \hat{f}) \leq L^2 K^2 \left( d^{-2\alpha} + C d^{-\alpha} + 2C d^{-3\alpha/2} \right) \leq C L^2 K^2 d^{-\alpha}.$$

□

*Proof.* By definition,  $f \in \mathcal{F}(R, 1)$  implies there exist a set of mappings  $\phi_k: [0, 1] \mapsto [R]$  such that

$$f(\mathbf{x}) = \sum_{\mathbf{r} \in [R]^K} \mathcal{C}(\mathbf{r}) \mathbb{1} \{ \phi_k(x_k) = r_k, \text{ for all } k \in [K] \}, \text{ for all } \mathbf{x} \in [0, 1]^K.$$

Without loss of generality, assume  $\phi_k^{-1}(r) = (\lambda_{r-1}^{(k)}, \lambda_r^{(k)}]$ , where  $0 = \lambda_1^{(k)} \leq \dots \leq \lambda_{R-1}^{(k)} \leq \lambda_R^{(k)} = 1$  are a sequence of cut-off points on  $[0, 1]$ , and the interval length preserves the Lebesgue measure  $|\{x_k \in \phi_k^{-1}(r)\}| = |\lambda_r^{(k)} - \lambda_{r-1}^{(k)}|$  for all  $r \in [R]$ . Under this assumption,  $f$  has an isomorphic form

$$f(\mathbf{x}) = \sum_{\mathbf{r} \in [R]^K} \mathcal{C}(\mathbf{r}) \mathbb{1} \left\{ x_k \in (\lambda_{r_k-1}^{(k)}, \lambda_{r_k}^{(k)}], \text{ for all } k \in [K] \right\}, \text{ for all } \mathbf{x} \in [0, 1]^K. \quad (4)$$

Now consider  $f \Rightarrow \Theta$ . By definition, for all  $(i_1, \dots, i_K) \in [d_1] \times \dots \times [d_K]$ ,

$$\Theta(i_1, \dots, i_K) = \sum_{\mathbf{r} \in [R]^K} \mathcal{C}(\mathbf{r}) \mathbb{1} \left\{ \xi_{i_k}^{(k)} \in (\lambda_{r_k-1}^{(k)}, \lambda_{r_k}^{(k)}], \text{ for all } k \in [K] \right\},$$

where  $\{\xi_{i_k}^{(k)}\}$  are i.i.d.  $\text{Unif}[0, 1]$  for all  $i_k \in [d]$  and  $k \in [K]$ .

From  $\Theta \Rightarrow \hat{f}$ , we have

$$\begin{aligned}
\hat{f}(\mathbf{x}) &= \sum_{(i_1, \dots, i_K)} \Theta(i_1, \dots, i_K) \mathbb{1} \left\{ x_k \in \left( \frac{i_k - 1}{d}, \frac{i_k}{d} \right], \text{ for all } k \in [K] \right\} \\
&= \sum_{\mathbf{r} \in [R]^K} \mathcal{C}(\mathbf{r}) \sum_{(i_1, \dots, i_K)} \mathbb{1} \left\{ \xi_{i_k}^{(k)} \in \left( \lambda_{r_k-1}^{(k)}, \lambda_{r_k}^{(k)} \right] \text{ and } x_k \in \left( \frac{i_k - 1}{d}, \frac{i_k}{d} \right], \text{ for all } k \in [K] \right\} \\
&= \sum_{\mathbf{r}} \mathcal{C}(\mathbf{r}) \prod_{k \in [K]} \left( \sum_{i_k \in [d]} \mathbb{1} \left\{ \xi_{i_k}^{(k)} \in \left( \lambda_{r_k-1}^{(k)}, \lambda_{r_k}^{(k)} \right] \text{ and } x_k \in \left( \frac{i_k - 1}{d}, \frac{i_k}{d} \right] \right\} \right)
\end{aligned}$$

For each  $k \in [K]$ , define the empirical cumulative proportion of categories among the set  $\{\xi_1^{(k)}, \dots, \xi_d^{(k)}\}$

$$\hat{\lambda}_r^{(k)} = \frac{1}{d} \sum_{i \in [d]} \mathbb{1} \left\{ \xi_i^{(k)} \leq \lambda_r^{(k)} \right\}, \text{ for all } r \in [R].$$

The function  $\hat{f}$  can be equivalently written as follows:

$$\begin{aligned}
\hat{f}(\mathbf{x}) &= \sum_{\mathbf{r}} \mathcal{C}(\mathbf{r}) \prod_{k \in [K]} \mathbb{1} \left\{ x_k \in (\hat{\lambda}_{r_k-1}^{(k)}, \hat{\lambda}_{r_k}^{(k)}] \right\} \\
&= \sum_{\mathbf{r}} \mathcal{C}(\mathbf{r}) \mathbb{1} \{ x_k \in (\hat{\lambda}_{r_k-1}^{(k)}, \hat{\lambda}_{r_k}^{(k)}] \text{ for all } k \in [K] \}.
\end{aligned} \tag{5}$$

Applying Lemma 3 to (4) and (5), we conclude

$$\mathbb{E}[\text{loss}(f, \hat{f})] \leq \frac{8KR}{\sqrt{d}}.$$

□

We partition the set  $[0, 1]^K$  in two ways:

1. Partition based on  $f$ .

$$[0, 1]^K = \cup \{ \phi^{-1}(\mathbf{r}) : \mathbf{r} \in [R]^K \}, \text{ where } \phi^{-1}(\mathbf{r}) \stackrel{\text{def}}{=} \{ \mathbf{x} \in [0, 1]^K : \phi(\mathbf{x}) = \mathbf{r} \}.$$

2. Partition based on  $\hat{f}$ .

$$[0, 1]^K = \cup \{ h^{-1}(\mathbf{r}) : \mathbf{r} \in [R]^K \}, \text{ where } h^{-1}(\mathbf{r}) \stackrel{\text{def}}{=} \{ \mathbf{x} \in [0, 1]^K : h(\mathbf{x}, \mathbf{r}) = 1 \}.$$

We define a one-to-one mapping between  $\phi^{-1}(\mathbf{r})$  and  $h^{-1}(\mathbf{r})$ . The mapping is measure-preserving in the sense that  $|p_{\mathbf{r}}|$  (??)

We will evaluate the difference between  $\phi^{-1}(\mathbf{r})$  and  $h^{-1}(\mathbf{r})$ . The  $p_{\mathbf{r}} := |\phi^{-1}(\mathbf{r})| \in [0, 1]$  denote the Lebesgue measure of the set  $\phi^{-1}(\mathbf{r})$ . Note that  $h^{-1}(\mathbf{r})$  is a random set where the randomness

comes from  $\xi_\omega$ . In particular,

$$\begin{aligned}
\hat{p}_{\mathbf{r}} &\stackrel{\text{def}}{=} |h^{-1}(\mathbf{r})| = |\{\mathbf{x} : \sum_{\omega} \mathbb{1}\{\phi(\xi_\omega) = \mathbf{r}\} \mathbb{1}\{\mathbf{x} \in I_\omega\}\}| \\
&= \frac{1}{d^K} \sum_{\omega} \mathbb{1}\{\phi(\xi_\omega) = \mathbf{r}\} \\
&= \frac{1}{d^K} \sum_{(i_1, \dots, i_K)} \left( \prod_k \mathbb{1}\{\phi_k(\xi_{k, i_k}) = r_k\} \right) \\
&= \frac{1}{d^K} \prod_k \left( \sum_{i_k \in [d]} \mathbb{1}\{\phi_k(\xi_{k, i_k}) = r_k\} \right) \\
&= \frac{1}{d^K} \prod_k \text{Bin}(d, \lambda_k(r_k))
\end{aligned}$$

where the event  $\mathbb{1}\{\phi_k(\xi_{k, i_k}) = r_k\}$  are i.i.d. Bernoulli random variables with success probability  $\lambda_k(r_k) = |\{x_k \in [0, 1] : \phi_k(x_k) = r_k\}| \in [0, 1]$  for all  $i_k \in [d]$  and  $k \in [K]$ . Therefore,

$$\hat{p}_{\mathbf{r}} \sim \frac{1}{d^K} \prod_k \text{Bin}(d, \lambda_k(r_k)) \quad \text{and} \quad p_{\mathbf{r}} = \prod_k \lambda_k(r_k).$$

which implies

$$\mathbb{E}(|\hat{p}_{\mathbf{r}} - p_{\mathbf{r}}|^2) \leq \frac{K p_{\mathbf{r}}}{d}.$$

Now, we evaluate the loss:

$$\int_{\mathbf{x} \in [0, 1]^K} |f(\mathbf{x}) - f'(\tau(\mathbf{x}))|^2 d\mathbf{x} \leq \sum_{\mathbf{r} \in [R]^K} \int_{\mathbf{x} \in \phi^{-1}(\mathbf{r})} |f(\mathbf{x}) - f'(\tau(\mathbf{x}))|^2 d\mathbf{x}$$

*Proof.*

$$\begin{aligned}
\text{Loss}(f, f') &= \inf_{f'_{\text{iso}} \sim f' \sim f} \int_{\mathbf{x} \in [0, 1]^K} |f'_{\text{iso}}(\mathbf{x}) - f(\mathbf{x})|^2 d\mathbf{x} \\
&\leq \inf_{\tau: \text{measure-preserving map}} \int_{\mathbf{x} \in [0, 1]^K} |f'(\tau(\mathbf{x})) - f(\mathbf{x})|^2 d\mathbf{x}.
\end{aligned}$$

By assumption,  $f$  is piecewise constant over the grid,  $\mathcal{G} = \cup_i I_i$ , where the intervals  $I_i \subset [0, 1]^K$  are disjoint for all  $i \in [R^K]$ . Note that

$$\int_{\mathbf{x} \in I_i} |f(\mathbf{x}) - f'(\mathbf{x})|^2 d\mathbf{x} = 4 \int_{\mathbf{x} \in I_i} \mathbb{1}\{f(\mathbf{x}) \neq f'(\mathbf{x})\} d\mathbf{x} = 4\mu\{\mathbf{x} \in I_i : f(\mathbf{x}) \neq f'(\mathbf{x})\},$$

where  $\mu\{\cdot\}$  denotes the Lebesgue measure. We aim to find two isomorphisms to upper bound the Lebesgue measure. Specifically, let  $c = f(\mathbf{x})$  denote the constant in the interval  $I_i$ . We want to

find the region for which  $\{\mathbf{x} \in I_i : f'(\tau(\mathbf{x})) \neq c\}$  where  $\tau$  is a measure-preserving map. Recall the definition  $f \Rightarrow \Theta(f, \xi) \Rightarrow f'$ .

First, we consider  $f \Rightarrow \Theta(f, \xi)$ . The random tensor  $\Theta = \Theta(f, \xi)$  is expressed as

$$\Theta(i_1, \dots, i_K) = f(\xi_{1,i_1}, \dots, \xi_{K,i_K}), \text{ for all } [i_1, \dots, i_K] \in [d_1] \times \dots \times [d_K].$$

For each fixed  $k \in [K]$ , we sort the elements in  $\{\xi_{k,i_k} : i_k \in [d_k]\}$  from smallest to largest. With a little abuse of notation, we denote  $\xi_{k,1} \leq \dots \leq \xi_{k,d_k}$  for all  $k \in [K]$ . The sorted tensor  $\Theta^*$  is expressed as

$$\Theta^*(i_1, \dots, i_K) = \Theta(\sigma_1(i_1), \dots, \sigma_K(i_K)), \text{ for all } [i_1, \dots, i_K] \in [d_1] \times \dots \times [d_K],$$

where  $\sigma_k : [d_k] \mapsto [d_k]$  denotes the permutation that sorts  $\{\xi_{k,i_k}\}_{i_k \in [d_k]}$  in increasing order. Note that  $\Theta^* \sim \Theta$ . By property (i), there exists  $f^* \sim f$  such that  $\Theta^*$  is induced by  $f^*$ . Therefore, it suffices to evaluate the loss between functions  $f^*$  and  $f'$ . Furthermore, by property (ii), without loss of generality, we define  $f'$  the function induced by  $\Theta^*$ . We aim to investigate the value  $f'(\mathbf{x})$  for  $\mathbf{x} \in I$ .

Then, we consider  $\Theta^* \Rightarrow f'$ . By definition, for all  $(x_1, \dots, x_K) \in [0, 1]^K$ ,

$$\begin{aligned} f'(\mathbf{x}) &= \sum_{(i_1, \dots, i_K)} \Theta^*(i_1, \dots, i_K) \mathbb{1} \left\{ \mathbf{x} \in \left( \frac{i_1-1}{d_1}, \frac{i_1}{d_1} \right] \times \dots \times \left( \frac{i_K-1}{d_K}, \frac{i_K}{d_K} \right] \right\}, \\ &= \sum_{(i_1, \dots, i_K)} f(\xi_{1,i_1}, \dots, \xi_{K,i_K}) \mathbb{1} \left\{ \mathbf{x} \in \left( \frac{i_1-1}{d_1}, \frac{i_1}{d_1} \right] \times \dots \times \left( \frac{i_K-1}{d_K}, \frac{i_K}{d_K} \right] \right\} \\ &= \sum_{\mathbf{r}} C(\mathbf{r}) \mathbb{1} \{ \phi(\mathbf{x}) = \mathbf{r} \} \end{aligned}$$

Recall that  $\Theta \in \mathcal{P}(R, 1) \dots$ . Suppose  $\xi_1 < \dots < \xi_d$ .  $\xi_k \in R$ ,  $\xi_{k+1} \in R, \dots$  then  $d_1 x_1$

The only places that  $f$  and  $f'$  differ are

$$|I - \hat{I}|, \quad \text{where } d_{\text{total}} \hat{I} \sim \text{Bernoulli}(d_{\text{total}}, I).$$

Therefore

$$\mathbb{E}(|I - \hat{I}|) \leq \sqrt{\frac{I}{d_{\text{total}}}}.$$

Summing over  $I_i$  over  $i \in [R^K]$  and note that  $\sum_i I_i = 1$ , we have

$$|\{\mathbf{x} : f(\mathbf{x}) \neq f'(\mathbf{x})\}|_{\lambda} \leq \sum_{i \in [R^K]} \mathbb{E}(|I_i - \hat{I}_i|) \leq \sum_{i \in [R^K]} \frac{\sqrt{I_i}}{\sqrt{d^K}} \leq \left( \frac{R}{d} \right)^{K/2}.$$

□

## 4 Main results

**Theorem 4.1** (loss in  $\mathcal{F}(R)$  family). *Let  $f \in \mathcal{F}(R)$  and  $\mathcal{Y} \sim \mathbb{P}(f, \boldsymbol{\xi})$  generated from the tensor nonparametric model. Let  $\hat{f}$  be the least-square estimator with pre-specified block size  $R$ . Then,*

$$\mathbb{E} \left[ \text{Loss}(f, \hat{f}) \right] \leq C \left( \frac{R^K}{d^K} + \frac{K \log R}{d^{K-1}} \right) + \frac{8KR}{\sqrt{d}},$$

where the expectation is taken jointly over  $\mathcal{Y}$  and  $\boldsymbol{\xi}$  (??) (uniformly over  $f \in \mathcal{F}(R)$ ??). In particular, let  $\Theta = f(\boldsymbol{\xi})$  and  $\hat{\Theta} = \hat{f}(\boldsymbol{\xi})$ . Then,

$$\mathbb{E} \left[ \text{Loss}(\Theta, \hat{\Theta}) \right] \leq C \left( \frac{R^K}{d^K} + \frac{\log R}{d^{K-1}} \right), \quad (6)$$

where again the expectation is taken jointly over  $\mathcal{Y}$  and  $\boldsymbol{\xi}$  (? or conditinoally).

*Proof.* We first proof (6). Recall that  $\hat{\Theta} = \arg \min_{\Theta \in \mathcal{F}(R, M)} F(\Theta)$  and  $\omega_{\Theta \in \mathcal{F}(R, M)} \nabla^2 F(\Theta) = \frac{1}{2}$ . By Taylor expansion of  $F(\Theta)$  around  $\hat{\Theta}$ , we have

$$\begin{aligned} \|\hat{\Theta} - \Theta^{\text{true}}\|_F &\leq 2 \left\langle \Theta^{\text{true}} - \mathcal{Y}, \frac{\hat{\Theta} - \Theta^{\text{true}}}{\|\hat{\Theta} - \Theta^{\text{true}}\|_F} \right\rangle \\ &\leq 2 \max_{\Theta \in \mathcal{P}(2R, 1)} \langle \mathcal{E}, \Theta \rangle. \end{aligned}$$

By union bound, for any  $t > 0$

$$\begin{aligned} \mathbb{P} \left( \max_{\Theta^{\text{true}} \in \mathcal{P}(R, M)} \|\hat{\Theta} - \Theta^{\text{true}}\|_F \geq t \right) &\leq \sup_{\boldsymbol{\xi}} \mathbb{P} \left( \max_{\Theta \in \mathcal{P}(2R, 1)} \langle \mathcal{E}, \Theta \rangle \geq \frac{t}{2} \middle| \boldsymbol{\xi} \right) \\ &\leq \sup_{\boldsymbol{\xi}} C |\mathcal{P}(2R, 1)| \exp \left( -\frac{t^2}{4\sigma^2} \right) \\ &\leq C \exp \left( -\frac{t^2}{\sigma^2} + Kd \log R + R^K \right). \end{aligned}$$

We take  $t = \sigma \sqrt{Kd \log R + R^K}$  and obtain that,

$$\max_{\Theta^{\text{true}} \in \mathcal{P}(R, M)} \|\hat{\Theta} - \Theta^{\text{true}}\|_F^2 \leq \sigma^2 (Kd \log R + R^K),$$

with probability at least  $1 - \exp(-Kd \log R - R^K)$ , (or equivalently, uniformly over  $f \in \mathcal{F}(R, M)$  and  $\boldsymbol{\xi} \in \mathbb{P}_{\boldsymbol{\xi}}$ ). Therefore,

$$\mathbb{E}(\text{Loss}(\hat{\Theta}, \Theta^{\text{true}})) \leq \sigma \left( \frac{r^K}{d^K} + \frac{K \log r}{d^{K-1}} \right),$$

where the expectation with respect to  $\mathcal{Y}$  (and  $\boldsymbol{\xi}$ ). □

**Lemma 2.**

$$\text{Loss}(f, \hat{f}) \leq \text{Loss}(\Theta, \hat{\Theta}) + \frac{2KR}{\sqrt{d}}.$$

*Proof.* Define  $f \Rightarrow \Theta(f, \xi) \Rightarrow \tilde{f}$  (a random function measure-able w.r.t.  $\xi$ ).

$$\begin{aligned} \mathbb{E} [\text{Loss}(f, \hat{f})] &\leq \mathbb{E} [\text{Loss}(\tilde{f}, \hat{f})] + \mathbb{E} [\text{Loss}(\tilde{f}, f)] \\ &\leq \mathbb{E} [\text{Loss}(\Theta, \hat{\Theta})] + \mathbb{E} [\text{Loss}(\tilde{f}, f)] \\ &\leq \sigma^2 \left( \frac{K \log R}{d^{K-1}} + \frac{R}{d^K} \right) + \frac{2KR}{\sqrt{d}}, \end{aligned}$$

where the expectation is over  $\mathcal{E}$  and  $\xi$  jointly.  $\square$

**Remark 7.** Dense region: agonistic error dominates. Sparse region: estimation error dominates... intuition?

**Problem 2** (Non-parametric estimation with unobserved designs). Let  $f: [0, 1]^K \mapsto [-1, 1]$  be a target function of interest. For each  $k \in [K]$ , we draw a random i.i.d. sample of  $d$  points  $\{x_i^{(k)}\}_{i \in [d]}$  uniformly from  $[0, 1]$ . Write  $\omega = (i_1, \dots, i_K) \in [d]^K$  and  $\mathbf{x}_\omega = (x_{i_1}^{(1)}, \dots, x_{i_K}^{(K)})^T \in \mathbb{R}^K$ . The goal is to estimation  $f$  given the set  $\{(\omega, f(\mathbf{x}_\omega))\}$  but without knowing  $\{\mathbf{x}_\omega\}$ !

The following lemma gives the estimation error for stepwise constant functions  $f \in \mathcal{F}(R, 1)$ .

**Lemma 3** (Step function approximation with unobserved designs). *Let  $f: [0, 1]^K \mapsto [-1, 1]$  be an  $R$ -step function:*

$$f(\mathbf{x}) = \sum_{\mathbf{r} \in [R]^K} c_{\mathbf{r}} \mathbb{1} \left\{ \mathbf{x} \in \underbrace{(\lambda_{r_1-1}^{(1)}, \lambda_{r_1}^{(1)}) \times \dots \times (\lambda_{r_K-1}^{(K)}, \lambda_{r_K}^{(K)})}_{=: I_{\mathbf{r}}} \right\}, \quad \text{for all } \mathbf{x} \in [0, 1]^K,$$

where  $\mathbf{r} = (r_1, \dots, r_K) \in [R]^K$  denotes multi-index,  $\{c_{\mathbf{r}} \in [-1, 1]\}$  are a set of real numbers, and  $0 = \lambda_1^{(k)} \leq \lambda_2^{(k)} \leq \dots \leq \lambda_{R-1}^{(k)} \leq \lambda_R^{(k)} = 1$  are a sequence of cutoff points over  $[0, 1]$ , for  $k \in [K]$ . Let  $p_r^{(k)} = \lambda_r^{(k)} - \lambda_{r-1}^{(k)}$  denote the length of the  $r$ -th interval at the mode  $k$ , where  $r \in [R]$  and  $k \in [K]$ .

For each  $k \in [K]$ , draw a random i.i.d. sample of  $d$  points  $\{N_i^{(k)}\}_{i \in [d]}$  from a categorical distribution with parameter  $(p_1^{(k)}, \dots, p_R^{(k)})$ . Let

$$\hat{\lambda}_r^{(k)} = \frac{1}{d} \sum_i \mathbb{1}\{N_i^{(k)} \leq r\}, \quad \text{for all } r \in [R]$$

denote the empirical cumulative proportions based on the  $d$  trials. Consider the estimator  $\hat{f}$ :

$$\hat{f}(\mathbf{x}) = \sum_{\mathbf{r} \in [R]^K} c_{\mathbf{r}} \mathbb{1} \left\{ \mathbf{x} \in \underbrace{(\hat{\lambda}_{r_1-1}^{(1)}, \hat{\lambda}_{r_1}^{(1)}) \times \dots \times (\hat{\lambda}_{r_K-1}^{(K)}, \hat{\lambda}_{r_K}^{(K)})}_{=: \hat{I}_{\mathbf{r}}} \right\}, \quad \text{for all } \mathbf{x} \in [0, 1]^K.$$

Then

$$\mathbb{E} \left( \int_{\mathbf{x} \in [0,1]^K} |f(\mathbf{x}) - \hat{f}(\mathbf{x})|^2 dx \right) \leq \frac{8KR}{\sqrt{N}},$$

where the expectation is taken with respect to  $\{N_i^{(k)}\}$ .

*Proof.* Note that

$$\mathbb{E} \left( \int_{\mathbf{x} \in [0,1]^K} |f(\mathbf{x}) - \hat{f}(\mathbf{x})|^2 dx \right) \leq 4\mathbb{E} \int_{\mathbf{x} \in [0,1]^K} \mathbb{1}\{f(\mathbf{x}) \neq \hat{f}(\mathbf{x})\} d\mathbf{x}.$$

Therefore, it suffices to evaluate the Lebesgue measure of  $\{\mathbf{x}: f(\mathbf{x}) \neq \hat{f}(\mathbf{x})\}$ . Note that  $f$  and  $f'$  are the same over  $I_{\mathbf{r}} \cap \hat{I}_{\mathbf{r}}$  for all  $\mathbf{r}$ . This implies

$$\begin{aligned} \{\mathbf{x}: f(\mathbf{x}) \neq f'(\mathbf{x})\} &\subset \cup_{\mathbf{r}} \{I_{\mathbf{r}} \Delta \hat{I}_{\mathbf{r}}\} \\ &\subset \cup_k \cup_r \left\{ \mathbf{x}: x_k \in (\lambda_{r-1}^{(k)}, \lambda_r^{(k)}) \Delta (\hat{\lambda}_{r-1}^{(k)}, \hat{\lambda}_r^{(k)}) \right\}, \end{aligned}$$

where the operation  $\Delta$  denotes the symmetric difference between two sets, and the second line comes from the property for Cartesian product of intervals. Then,

$$\begin{aligned} |\{\mathbf{x}: f(\mathbf{x}) \neq f'(\mathbf{x})\}| &\leq \sum_{k \in [K]} \sum_{r \in [R]} \left| (\lambda_{r-1}^{(k)}, \lambda_r^{(k)}) \Delta (\hat{\lambda}_{r-1}^{(k)}, \hat{\lambda}_r^{(k)}) \right| \\ &\leq 2 \sum_{k \in [K]} \sum_{r \in [R]} |\lambda_r^{(k)} - \hat{\lambda}_r^{(k)}| \end{aligned} \tag{7}$$

Note that by definition,  $d\hat{\lambda}_r^{(k)} = \sum_i \mathbb{1}\{N_i^{(k)} \leq r\}$  follows from Binomial distribution with parameters  $(d, \lambda_r^{(k)})$ . Therefore,

$$\mathbb{E} |\lambda_r^{(k)} - \hat{\lambda}_r^{(k)}|^2 = \frac{\lambda_r^{(k)}(1 - \lambda_r^{(k)})}{d} \leq \frac{1}{d}. \tag{8}$$

Plugging (8) into (7), we obtain

$$\begin{aligned} \left| \{\mathbf{x}: f(\mathbf{x}) \neq \hat{f}(\mathbf{x})\} \right|^2 &\leq 4KR \sum_{k \in [K]} \sum_{r \in [R]} |\lambda_r^{(k)} - \hat{\lambda}_r^{(k)}|^2 \\ &\leq \frac{4K^2 R^2}{d} \end{aligned}$$

Henceforth,

$$\mathbb{E} \left| \{\mathbf{x}: f(\mathbf{x}) \neq \hat{f}(\mathbf{x})\} \right| \leq \frac{2KR}{\sqrt{d}}.$$

□

## 5 Minimax lower bound

**Theorem 5.1** (Minimax). *For stochastic tensor model, there exists a constant  $C > 0$  such that*

$$\inf_{\hat{\Theta}} \sup_{\Theta \in \mathcal{P}(R)} \mathbb{P} \left\{ \|\hat{\Theta} - \Theta\|_F > C \frac{\sigma^2}{p} \left( \frac{R^K + Kd \log R}{d^K} \right) \right\} > 0.2.$$

*Proof. Nonparametric rate.* We construct fixed  $\mathbf{M}^{(k)} \in [0, 1]^{d \times R}$  for all  $k \in [K]$  as follows. For each  $k \in [K]$ ,  $\mathbf{M}: [d] \mapsto [R]$  partition the set  $[d]$  into  $R$  equal-sized clusters, and we denote the mapping rule  $\mathbf{M}(i) = \lceil \frac{iR}{d} \rceil$  for  $i \in [d]$ . For any binary tensor  $\mathcal{C} \in \{0, 1\}^{r_1 \times \dots \times r_K}$ , we define the core tensor

$$\check{\mathcal{C}} = c \sqrt{\frac{\sigma^2 r^k}{p d^k}} \mathcal{C}.$$

We identify the tensor in  $\{0, 1\}^{r_1 \times \dots \times r_K}$  by vectors in  $\{0, 1\}^{r^K}$ . By Lemma .., there exists some set  $\Omega$  such that  $|\Omega| \geq \exp\left(\frac{r^K}{4}\right)$  and  $H(\mathcal{C}, \mathcal{C}') \geq \frac{r^K}{4}$  for any  $\mathcal{C} \neq \mathcal{C}' \in \Omega$ . We consider the subspace in the original tensor induced by \$ \square\$

**Definition 9** (Clustering function). Let  $M: [d] \mapsto [R]$  denotes a clustering function, where we use  $M(i) \in [R]$  denote the cluster label to which the  $i$  was assigned. We use  $M^{-1}(r) = \{d: M(d) = r\} \subset [d]$  to denote the set of indices that was assigned to cluster  $r$ . For notational convenience, we use  $M \in [R]^d$  or  $M \in \{0, 1\}^{d \times R}$  exchangeable to denote the collection of all clustering functions that map  $[d]$  to  $[R]$ .

**Lemma 4.** *There exists a subset  $\Omega \in [R]^d$ , such that  $|\Omega| \geq \exp(d \log R/2)$  and*

$$H(\omega, \omega') \stackrel{\text{def}}{=} |\{d: \omega(d) \neq \omega'(d)\}| \geq \frac{d}{4}, \quad \text{for all } \omega \neq \omega' \in \Omega.$$

*Proof.* Define

$$\Omega = \left\{ \omega \in [R]^d: |\omega^{-1}(i)| = \frac{d}{R} \text{ for } i \in [R] \right\},$$

that is,  $\Omega$  is the collection of clustering functions that generates the equal-sized clusters. Given any  $\omega \in \Omega$ , define its  $\varepsilon$ -neighborhood by

$$B(\omega, \varepsilon) = \{\omega' \in \Omega: H(\omega, \omega') \leq \varepsilon\}.$$

The packing number of  $\Omega$

$$\mathcal{M}(\varepsilon, \Omega, H) \geq \frac{|\Omega|}{\max_{\omega \in \Omega} |B(\omega, \varepsilon)|}.$$

Taking  $\varepsilon = \frac{d}{4}$  gives

$$|B(\omega, \varepsilon)| \leq \binom{d}{d/4} R^{d/4} \leq (4e)^{d/4} R^{d/4} \leq \exp\left(\frac{1}{4} d \log R\right),$$



where we have used the inequality  $\binom{d}{k} \leq (\frac{ed}{k})^k$  for all  $k \leq d$ . By stirling's formula

$$|\Omega| = \frac{d!}{[(d/R)!]^R} \approx \frac{\sqrt{d}(\frac{d}{e})^d}{\sqrt{\frac{d}{R}}(\frac{d}{Re})^d} \approx \sqrt{R}(R)^d \geq \exp(d \log R + o(d \log R)) \geq \exp\left(\frac{1}{2}d \log R\right).$$

Therefore,

$$\mathcal{M}\left(\frac{d}{4}, \Omega, H\right) \geq \exp\left(\frac{1}{2}d \log R\right).$$

□

In Wang et al, the MSE for the tensor block model has the following asymptotical rate,

$$\text{MSE}(\hat{\Theta}, \Theta) \stackrel{\text{def}}{=} \frac{1}{d^K} \|\Theta - \hat{\Theta}\|_F^2 \asymp \frac{R^K}{d^K} + \frac{K \log R}{d^{K-1}}, \quad \text{as } d \rightarrow \infty \text{ while fixing } R.$$

We investigate the asymptotic behavior as  $R \asymp d^\delta$ , for some  $\delta \in [0, 1]$ .

$$\text{MSE}(\hat{\Theta}, \Theta) = \begin{cases} d^{-K}, & R = 1, \\ d^{-(K-1)}, & \delta = 0 \text{ (i.e. constant } R = \mathcal{O}(1)), \\ d^{-(K-1)} \log d, & \delta \in (0, \frac{1}{K}], \\ d^{-K(1-\delta)}, & \delta \in (\frac{1}{K}, 1]. \end{cases}$$

In particular, in the matrix case when  $K = 2$ , the asymptotic rate is  $\mathcal{O}(d^{-1} \log d)$  whenever  $R \asymp \mathcal{O}(\sqrt{d})$ . This is faster than regular low-rank matrix decomposition but slower than Gao's paper.