Communication-Efficient Estimation of

Regularized Smoothed Support Tensor Machine

Zihao Song¹, Lei Wang², Riquan Zhang³, Weihua Zhao¹

1. Nantong University, 2. Nankai University

3. Shanghai University of International Business and Economics

Supplementary Material

S1 Background on t-product

Before presenting the t-product and related theories (Kilmer et al., 2008; Kilmer and Martin, 2011), we introduce some additional notations. For a tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, denote the block circular matrix $\mathrm{bcirc}(\mathcal{A})$ as

$$\overline{m{A}} = \mathrm{bcirc}(m{\mathcal{A}}) := egin{bmatrix} m{\mathcal{A}}^{(1)} & m{\mathcal{A}}^{(I_3)} & \cdots & m{\mathcal{A}}^{(2)} \ m{\mathcal{A}}^{(2)} & m{\mathcal{A}}^{(1)} & \cdots & m{\mathcal{A}}^{(3)} \ & dots & dots & \ddots & dots \ m{\mathcal{A}}^{(I_3)} & m{\mathcal{A}}^{(I_3-1)} & \cdots & m{\mathcal{A}}^{(1)} \ \end{pmatrix} \in \mathbb{R}^{I_1 I_3 \times I_2 I_3}.$$

The block vector-unfolding operator $bvec(\cdot)$ and its inverse operator $bvfold(\cdot)$

are defined as, respectively,

$$\operatorname{bvec}(\mathcal{A}) := [\mathcal{A}^{(1)}; \cdots; \mathcal{A}^{(I_3)}], \ \operatorname{bvfold}(\operatorname{bvec}(\mathcal{A})) := \mathcal{A}.$$

Definition S1.1 (t-product): The t-product between $\mathcal{A} \in \mathbb{R}^{I_1 \times r \times I_3}$ and $\mathcal{B} \in \mathbb{R}^{r \times I_2 \times I_3}$ is defined as

$$\mathbf{\mathcal{A}} * \mathbf{\mathcal{B}} := \text{bvfold}(\text{bcirc}(\mathbf{\mathcal{A}}) \cdot \text{bvec}(\mathbf{\mathcal{B}})) \in \mathbb{R}^{I_1 \times I_2 \times I_3}$$

where * *indicate the t-product.*

To obtain the faster computational efficiency, the discrete Fourier transformation (DFT), connecting t-product and standard matrix product, is introduced. Denote $\boldsymbol{F}_n = [\omega_1, \omega_2, \cdots, \omega_n]$ as the DFT matrix, where

$$\omega_i = [\vartheta^{0 \times (i-1)}; \vartheta^{1 \times (i-1)}; \cdots; \vartheta^{(n-1) \times (i-1)}],$$

 $\vartheta=e^{-\frac{2\pi\varrho}{n}}$ and $\varrho=\sqrt{-1}$. Thus, \boldsymbol{F}_n is a unitary matrix, that is, $\boldsymbol{F}_n^{\mathsf{H}}\boldsymbol{F}_n=\boldsymbol{F}_n^{\mathsf{H}}\boldsymbol{F}_n^{\mathsf{H}}=n\boldsymbol{I}_n$ and $\boldsymbol{F}_n^{-1}=\frac{1}{n}\boldsymbol{F}_n^{\mathsf{H}}$. Then, the block-circulant matrices can be

blockdiagonalized, i.e.,

$$(\boldsymbol{F}_{I_3} \otimes \boldsymbol{I}_{I_1}) \overline{\boldsymbol{A}} (\boldsymbol{F}_{I_3}^{\mathsf{H}} \otimes \boldsymbol{I}_{I_2}) = \widetilde{\boldsymbol{A}} = \begin{bmatrix} \widetilde{\boldsymbol{\mathcal{A}}}^{(1)} & & & & \\ & \widetilde{\boldsymbol{\mathcal{A}}}^{(2)} & & & & \\ & & \ddots & & & \\ & & & \widetilde{\boldsymbol{\mathcal{A}}}^{(I_3)} \end{bmatrix}$$
 (S1.1)

where $\widetilde{\mathcal{A}}$ is obtained by $\widetilde{\mathcal{A}}(i,j,:) = \mathbf{F}_{I_3} \mathcal{A}(i,j,:)$. By the unitarity and (S1.1),

$$\|\mathcal{A}\|_F = \frac{1}{\sqrt{I_3}} \|\widetilde{A}\|_F, \ \langle \mathcal{A}, \mathcal{B} \rangle = \frac{1}{I_3} \langle \widetilde{A}, \widetilde{A} \rangle,$$
 $C = \mathcal{A} * \mathcal{B} \Leftrightarrow \widetilde{C} = \widetilde{\mathcal{A}} \cdot \widetilde{\mathcal{B}}$

Next, we introduce some other definitions on tensor, generalized from the matrix case.

Definition S1.2 (Conjugate transpose, Lu et al. (2018)): The conjugate transpose of $\mathcal{A} \in \mathbb{C}^{I_1 \times I_2 \times I_3}$ is represented as $\mathcal{A}^H \in \mathbb{C}^{I_2 \times I_1 \times I_3}$, achieved by transposing each frontal slice with conjugation and subsequently reversing the sequence of transposed frontal slices from 2 to I_3 .

Definition S1.3 (Orthogonal tensor): A tensor $Q \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ is orthogonal if $Q * Q^H = Q^H * Q = \mathcal{I}$.

Definition S1.4 (F-diagonal tensor): A tensor is considered f-diagonal when

each of its frontal slices represents a diagonal matrix.

Definition S1.5 (*Identity tensor*): A tensor $\mathcal{I} \in \mathbb{R}^{I \times I \times I_3}$ is called as identity tensor if its first frontal slice is the identity matrix and others are all zero.

Definition S1.6 (t-SVD): For a given $A \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, it has the following factorization,

$$\mathcal{A} = \mathcal{U} * \mathcal{S} * \mathcal{V}^{\mathsf{H}},$$

where $\mathcal{U} \in \mathbb{R}^{I_1 \times I_1 \times I_3}$, $\mathcal{V} \in \mathbb{R}^{I_2 \times I_2 \times I_3}$ are orthogonal tensors, and $\mathcal{S} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ is a f-diagonal tensor.

Note that the entries of the first frontal slice in S possess the decreasing property, i.e.,

$$S(1,1,1) \ge S(2,2,1) \ge \cdots \ge S(I_1 \land I_2, I_1 \land I_2, 1).$$
 (S1.2)

The property (S1.2) holds because the inverse DFT gives,

$$\mathcal{S}(i,i,1) = \frac{1}{I_3} \sum_{j=1}^{I_3} \widetilde{\mathcal{S}}(i,i,j), \tag{S1.3}$$

and the $\widetilde{\mathcal{S}}(:,:,j)$ is the singular value matrix of $\widetilde{\mathcal{A}}(:,:,j)$. Therefore, the elements of $\mathcal{S}(:,:,1)$ are regarded as the singular values of \mathcal{A} .

Definition S1.7 (Tensor tubal-rank, Lu et al. (2018)): For $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, its

tubal-rank is the number of nonzero singular value, that is

$$rank_t(\mathbf{A}) := \#\{i, \mathbf{S}(i, i, 1) \neq 0\},\$$

where "#" denotes the cardinality of a set.

It is easy to obtain $\operatorname{rank}_{t}(\mathcal{A}) \leq \operatorname{rank}(\overline{\mathcal{A}}) \leq I_{3} \cdot \operatorname{rank}_{t}(\mathcal{A})$; see Lu et al. (2018) for details.

Definition S1.8 Define $\|A\|:=\|\overline{A}\|_{op}$ as the tensor spectral norm for $A\in\mathbb{R}^{I_1\times I_2\times I_3}$

Definition S1.9 (t-TNN, Lu et al. (2018)): Let $\mathcal{A} = \mathcal{U} * \mathcal{S} * \mathcal{V}^{\mathsf{H}}$ be the t-SVD of $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, and the tubal nuclear norm is defined as

$$\|\mathcal{A}\|_* := \sum_{i=1}^r \mathcal{S}(i, i, 1),$$

where $r = \operatorname{rank}_{t}(\mathcal{A})$.

Note that the tensor spectral norm is the dual norm of t-TNN. It is easy to see that the following properties hold,

$$\|\mathcal{A}\| = \|\overline{A}\|_{op} = \|\widetilde{A}\|_{op}, \|\mathcal{A}\|_{*} = \frac{1}{I_{3}} \|\overline{A}\|_{*} = \frac{1}{I_{3}} \|\widetilde{A}\|_{*}.$$
 (S1.4)

Proposition S1.1 (Lu et al., 2018) Let $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ with $\mathrm{rank_t}(\mathcal{A}) = r$,

and its skinny t-SVD be $\mathcal{U} * \mathcal{S} * \mathcal{V}^{\mathsf{H}}$ where $\mathcal{U} \in \mathbb{R}^{I_1 \times r \times I_3}$, $\mathcal{S} \in \mathbb{R}^{r \times r \times I_3}$ and $\mathcal{V} \in \mathbb{R}^{I_2 \times r \times I_3}$. The subdifferential of t-TNN is $\partial \|\mathcal{A}\|_* = \{\mathcal{U} * \mathcal{V}^{\mathsf{H}} + \mathcal{W} : \mathcal{U}^{\mathsf{H}} * \mathcal{W} = 0, \mathcal{W} * \mathcal{V} = 0, \|\mathcal{W}\| \le 1\}$.

With the above discussion, given $\mathcal{A}, \mathcal{B} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$,

$$\langle \mathcal{A}, \mathcal{B} \rangle \leq \frac{1}{I_3} \langle \widetilde{A}, \widetilde{B} \rangle \leq \frac{1}{I_3} \|\widetilde{A}\|_* \|\widetilde{B}\| = \|\mathcal{A}\|_* \|\mathcal{B}\|.$$

Moreover, let $\mathbf{A} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ with $\mathrm{rank_t}(\mathbf{A}) = r$,

$$\|\mathcal{A}\|_* = \frac{1}{I_3} \|\widetilde{A}\|_* \le \frac{1}{I_3} \sqrt{r'} \|\widetilde{A}\|_F \le \sqrt{r} \|\mathcal{A}\|_F,$$

where r' is the rank of \widetilde{A} and $r' \leq I_3 r$ holds. To prove that t-TNN is decomposable, the additivity of nuclear norm is needed.

Proposition S1.2 (Recht et al., 2010) For any \mathbf{A} and \mathbf{B} with the same dimension, if $\mathbf{A}\mathbf{B}^{\mathsf{H}} = \mathbf{A}^{\mathsf{H}}\mathbf{B} = 0$, then $\|\mathbf{A} + \mathbf{B}\|_* = \|\mathbf{A}\|_* + \|\mathbf{B}\|_*$.

Lemma S1.1 For any \mathcal{A} and \mathcal{B} with the same dimension, if $\mathcal{A}*\mathcal{B}^H = \mathcal{A}^H*\mathcal{B} = 0$, we have $\|\mathcal{A} + \mathcal{B}\|_* = \|\mathcal{A}\|_* + \|\mathcal{B}\|_*$.

Proof of Lemma S1.1: *Note that,*

$$\boldsymbol{\mathcal{A}}\ast\boldsymbol{\mathcal{B}}^{\mathsf{H}}=\boldsymbol{\mathcal{A}}^{\mathsf{H}}\ast\boldsymbol{\mathcal{B}}=0\rightarrow\widetilde{\boldsymbol{A}}\widetilde{\boldsymbol{B}}^{\mathsf{H}}=\widetilde{\boldsymbol{A}}^{\mathsf{H}}\widetilde{\boldsymbol{B}}=0.$$

Thus,

$$\|\mathcal{A} + \mathcal{B}\|_* = \frac{1}{I_3} \|\widetilde{(A + B)}\|_* = \frac{1}{I_3} (\|\widetilde{A}\|_* + \|\widetilde{B}\|_*) = \|\mathcal{A}\|_* + \|\mathcal{B}\|_*.$$

Theorem S1.1 For any $\mathcal{A} \in \mathcal{U}^{\perp}$ and $\mathcal{B} \in \mathcal{B}$, of the same dimension, with tubal rank r, we have $\|\mathcal{A} + \mathcal{B}\|_* = \|\mathcal{A}\|_* + \|\mathcal{B}\|_*$. where

$$egin{aligned} \mathscr{U} &= \Big\{ \mathcal{U} * \mathcal{M} + \mathcal{N} * \mathcal{V}^\mathsf{H} : \mathcal{M} \in \mathbb{R}^{r imes I_2 imes I_3}, \mathcal{N} \in \mathbb{R}^{I_1 imes r imes I_3} \Big\}, \ \mathscr{B} &= \Big\{ \mathcal{U} * \mathcal{U}^\mathsf{H} * \mathcal{Z} + \mathcal{Z} * \mathcal{V} * \mathcal{V}^\mathsf{H} - \mathcal{U} * \mathcal{U}^\mathsf{H} * \mathcal{Z} * \mathcal{V} * \mathcal{V}^\mathsf{H} : \mathcal{Z} \in \mathbb{R}^{I_1 imes I_2 imes I_3} \Big\}. \end{aligned}$$

Proof of Theorem S1.1 Note that, $\mathcal{B} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ has the skinny t-SVD, that is $\mathcal{B} = \mathcal{U} * \mathcal{S} * \mathcal{V}^{\mathsf{H}}$ where $\mathcal{U} \in \mathbb{R}^{I_1 \times r \times I_3}$, $\mathcal{S} \in \mathbb{R}^{r \times r \times I_3}$ and $\mathcal{V} \in \mathbb{R}^{I_2 \times r \times I_3}$. Define $\mathcal{U} = \left\{ \mathcal{U} * \mathcal{M} + \mathcal{N} * \mathcal{V}^{\mathsf{H}} : \mathcal{M} \in \mathbb{R}^{r \times I_2 \times I_3}, \mathcal{N} \in \mathbb{R}^{I_1 \times r \times I_3} \right\}$ and $\mathcal{B} = \left\{ \mathcal{U} * \mathcal{U}^{\mathsf{H}} * \mathcal{Z} + \mathcal{Z} * \mathcal{V} * \mathcal{V}^{\mathsf{H}} - \mathcal{U} * \mathcal{U}^{\mathsf{H}} * \mathcal{Z} * \mathcal{V} * \mathcal{V}^{\mathsf{H}} : \mathcal{Z} \in \mathbb{R}^{I_1 \times I_2 \times I_3} \right\}$, thus, $\mathcal{B}^{\perp} = \left\{ (\mathcal{I} - \mathcal{U} * \mathcal{U}^{\mathsf{H}}) * \mathcal{Z} * (\mathcal{I} - \mathcal{V} * \mathcal{V}^{\mathsf{H}}) : \mathcal{Z} \in \mathbb{R}^{I_1 \times I_2 \times I_3} \right\}$. It is easy to see that $\mathcal{B} \subseteq \mathcal{U}$ and $\mathcal{U}^{\perp} \subseteq \mathcal{B}^{\perp}$. Therefore, $\mathcal{A} * \mathcal{B}^{\mathsf{H}} = 0$ and $\mathcal{A}^{\mathsf{H}} * \mathcal{B} = 0$. Then, by Lemma S1.1, we have $\|\mathcal{A} + \mathcal{B}\|_* = \|\mathcal{A}\|_* + \|\mathcal{B}\|_*$.

S2 Proofs and Lemmas for results in Section 2

The theoretical developments are following the roadmap below: The Lemma S2.1 approximates the difference between empirical risk and population risk via Rademacher complexity. The Lemma S2.2 illustrates that the statistical error $\|\sum_i L'_h(u_i^*)y_i\mathcal{X}_i/n\|$ is bounded by λ through a kind of Talagrand's concentration inequality (Theorem 4.2 of Bartlett et al. (2005)). Furthermore, Theorem 1 investigates the upper bound of estimation error with the aid of Lemma S2.1 and S2.2. Then we study the size of $\mathbb{E}[\|\sum_i \sigma_i \mathcal{X}_i\|]$ and $\mathbb{E}[\|\sum_i L'_h(u_i^*)y_i\mathcal{X}_i\|]$ in **Assumption** 3 through the Lemma S2.3, and the tail probability $P(\|\mathcal{X}\| > M)$ with a proper constant M via the Lemma S2.4, respectively. With the above results, Corollary 1 presents the statistical rate of the RSSTM estimator. Next, we show that the RSSTM estimator enjoys the low rank guarantee of the order O(r). Finally, we investigate the convergence analysis of the proposed algorithm.

To simplify notations, the empirical measure on the n observations is denoted by P_n and the corresponding population counterpart is indicated by P. That is, with any function f on the same probability space as (\mathcal{X}, y) , we have $Pf = \int f dP = \mathbb{E}f$.

Lemma S2.1 Suppose Assumption 1-4 hold, for $\forall M > 0$ with $\mathbb{E}[\|\mathcal{X}\|^2]P(\|\mathcal{X}\| > 0)$

 $M) \leq \lambda/4$, we have

$$\sup_{(\delta, \mathbf{D}) \in \mathcal{T}} \left| (P_n - P) [L_h(y(\beta^* + \delta + \langle \mathbf{B}^* + \mathbf{D}, \mathbf{X} \rangle)) - L_h(y(\beta^* + \langle \mathbf{B}^*, \mathbf{X} \rangle))] \right| \le C \lambda t \sqrt{r},$$

with probability $1 - Ce^{-Cn\lambda/M} - CnP(\|\mathcal{X}\| > M)$, where $\mathscr{T} := \{(\delta, \mathcal{D}) : \|\mathcal{D}^{\perp}\|_* \le 3\|\mathcal{D}^0\|_* + |\delta|, \mathbb{E}[(\delta + \langle \mathcal{X}, \mathcal{D} \rangle)^2] \le t^2\}.$

Proof of Lemma S2.1: For simplicity of notations, let $f(\delta, \mathcal{D}, \mathcal{X}, y) = L_h(y(\beta^* + \delta + \langle \mathcal{B}^* + \mathcal{D}, \mathcal{X} \rangle)) - L_h(y(\beta^* + \langle \mathcal{B}^*, \mathcal{X} \rangle))$, and define $\mathcal{F} = \{f(\delta, \mathcal{D}, \mathcal{X}, y) : (\delta, \mathcal{D}) \in \mathcal{T}\}$. It is trivial that $\operatorname{Var}(f(\delta, \mathcal{D}, \mathcal{X}, y)) \leq C\operatorname{Var}(|y(\delta + \langle \mathcal{X}, \mathcal{D} \rangle)|) \leq Ct^2$ holds for $f \in \mathcal{F}$ by virtue of the Lipschitz continuity of the smoothed function $L_h(\cdot)$. Then the Markov inequality suggests that,

$$P(|(P_n - P)f| < a) \ge 1 - \frac{Ct^2}{na^2}, \ \forall \ a > 0.$$

The following inequality holds via the symmetrization Lemma (Pollard, 1984, Sec. II.3.8),

$$P(\sup_{f \in \mathcal{F}} |(P_n - P)f| > 2a) \le \frac{P(\sup_{f \in \mathcal{F}} |(P_n - P'_n)f| > a)}{1 - Ct^2/(na^2)},$$

where P'_n denotes the empirical measure on the independent copies $\{(y'_i, \mathcal{X}'_i)\}_{i=1}^n$ of $\{(y_i, \mathcal{X}_i)\}_{i=1}^n$.

Let $f_M = fI\{\|\mathcal{X}\| \leq M\}$ for $f \in \mathcal{F}$ where $I\{\cdot\}$ denotes the indicator

function. Note that $\sup_{f\in\mathcal{F}}(P_n-P_n')f_M$ and $\sup_{f\in\mathcal{F}}\sum_i\sigma_i f_M(\delta,\mathcal{D},\mathcal{X}_i,y_i)/n$ are identically distributed. By $\|\mathcal{D}\|_* \leq C(\|\mathcal{D}^0\|_* + |\delta|) \leq Ct\sqrt{r}$ due to the low tubal rank, we have $|f_M| \leq CMt\sqrt{r}$. Then, using $\operatorname{Var}(f) \leq Ct^2$ and the Talagrand's concentration inequality (see Theorem A.2 of Bartlett et al. (2005)), we have

$$P\left(\sup_{f\in\mathcal{F}}\frac{1}{n}\sum_{i}\sigma_{i}f_{M}(\delta,\mathcal{D},\mathcal{X}_{i},y_{i})>\mathbb{E}\left[\sup_{f\in\mathcal{F}}\frac{1}{n}\sum_{i}\sigma_{i}f_{M}(\delta,\mathcal{D},\mathcal{X}_{i},y_{i})\right]+\frac{CMx}{n}t\sqrt{r}\right)$$

$$\leq e^{-x}.$$
(S2.5)

On the other hand, we also have

$$\mathbb{E}(\|\sum_{i} \sigma_{i} y_{i} \mathcal{X}_{i} I\{\|\mathcal{X}_{i}\| > M\}\|) = n \mathbb{E}(\sup_{\|\mathcal{A}\| \leq 1} \sigma y \langle \mathcal{X} I\{\|\mathcal{X}\| > M\}, \mathcal{A}\rangle)$$
$$= n \mathbb{E}(\sup_{\|\mathcal{A}\| \leq 1} \langle \mathcal{X}, \mathcal{A} \rangle I\{\|\mathcal{X}\| > M\}) \leq n \mathbb{E}(\|\mathcal{X}\|^{2}) P(\|\mathcal{X}\| > M) \leq Cn\lambda,$$

which implies

$$\mathbb{E}\left[\sup_{f\in\mathcal{F}}\frac{1}{n}\sum_{i}\sigma_{i}f_{M}(\delta,\boldsymbol{\mathcal{D}},\boldsymbol{\mathcal{X}}_{i},y_{i})\right] \leq \mathbb{E}\left[\sup_{(\delta,\boldsymbol{\mathcal{D}})\in\mathcal{T}}\frac{1}{n}\sum_{i}\sigma_{i}y_{i}\left(\delta+\langle\boldsymbol{\mathcal{X}}_{i}I\{\|\boldsymbol{\mathcal{X}}_{i}\|\leq M\},\boldsymbol{\mathcal{D}}\rangle\right)\right]$$

$$\leq \mathbb{E}\left[\frac{1}{n}\|\sum_{i}\sigma_{i}y_{i}\boldsymbol{\mathcal{X}}_{i}I\{\|\boldsymbol{\mathcal{X}}_{i}\|\leq M\}\|\right] \cdot \left(\sup_{(\delta,\boldsymbol{\mathcal{D}})\in\mathcal{T}}\|\boldsymbol{\mathcal{D}}\|_{*}\right) + \mathbb{E}\left[\sup_{\delta}\frac{1}{n}\sum_{i}\sigma_{i}y_{i}\delta\right]$$

$$\leq \mathbb{E}\left[\frac{1}{n}\|\sum_{i}\sigma_{i}y_{i}\boldsymbol{\mathcal{X}}_{i}\|\right] \cdot \sup_{(\delta,\boldsymbol{\mathcal{D}})\in\mathcal{T}}\|\boldsymbol{\mathcal{D}}\|_{*} + \mathbb{E}\left[\frac{1}{n}\|\sum_{i}\sigma_{i}y_{i}\boldsymbol{\mathcal{X}}_{i}I\{\|\boldsymbol{\mathcal{X}}_{i}\|> M\}\|\right] \cdot \left(\sup_{(\delta,\boldsymbol{\mathcal{D}})\in\mathcal{T}}\|\boldsymbol{\mathcal{D}}\|_{*}\right)$$

$$+ \mathbb{E}\left[\sup_{\delta}\frac{1}{n}\sum_{i}\sigma_{i}y_{i}\delta\right] \leq C\lambda t\sqrt{r}.$$
(S2.6)

Choosing $x \approx n\lambda/M$ and $a = C\lambda t\sqrt{r}$, with (S2.5) and (S2.6), we have

$$P(\sup_{f \in \mathcal{F}} (P_n - P'_n) f_M > a/3) \le e^{-x}.$$

Then

$$P(\sup_{f \in \mathcal{F}} P_n(f - f_M) > a/3) = P(\sup_{f \in \mathcal{F}} P'_n(f - f_M) > a/3)$$

$$\leq P(\max_i ||\mathcal{X}_i|| > M) \leq nP(||\mathcal{X}|| > M).$$

This implies
$$P(\sup_{f \in \mathcal{F}} (P_n - P)f > a) \le Ce^{-x} + CnP(\|\mathcal{X}\| > M)$$
.

Lemma S2.2 For any M>0 satisfying $\mathbb{E}[\|\mathcal{X}\|^2]P(\|\mathcal{X}\|>M)\leq \lambda/4$, with probability $1-e^{-Cn\lambda/M}-nP(\|\mathcal{X}\|>M)$,

$$\lambda \ge 2\|\sum_i L_h'(u_i^*)y_i \mathcal{X}_i/n\|.$$

Proof of Lemma S2.2: For brevity, we indicate $\nu_i = L_h'(u_i^*)y_i$. Using the bound-edness of L_h' and optimality condition, we have $|\nu_i| \le 1$ and $\mathbb{E}[\nu_i \mathcal{X}_i] = 0$. In fact,

$$\|\sum_{i} \nu_{i} \boldsymbol{\mathcal{X}}_{i} / n\| = \sup_{\|\boldsymbol{\mathcal{A}}\|_{*} \leq 1} \frac{1}{n} \sum_{i} \nu_{i} \langle \boldsymbol{\mathcal{X}}_{i}, \boldsymbol{\mathcal{A}} \rangle.$$

Let
$$\mathscr{G}:=\{g:g(\nu,\mathcal{X})=\nu\langle\mathcal{X},\mathcal{A}\rangle,\|\mathcal{A}\|_*\leq 1\}$$
 and $g_M=g\mathrm{I}\{|g|\leq M\}$

 $MI\{g < -M\} + MI\{g > M\}$. Then, we have

$$\|\sum_{i} \nu_{i} \mathcal{X}_{i} / n\| = \sup_{g \in \mathcal{G}} (P_{n} - P)g \le \sup_{g \in \mathcal{G}} (P_{n} - P)g_{M} + |\sup_{g \in \mathcal{G}} (P_{n} - P)(g - g_{M})|.$$

For the first term of the above equation, with probability $1 - e^{-x}$,

$$\sup_{g \in \mathcal{G}} (P_n - P) g_M \leq C \mathbb{E}[\sup_{g \in \mathcal{G}} (P_n - P) g_M] + \frac{CMx}{n} \leq C \mathbb{E}[\|\sum_i \nu_i \mathcal{X}_i / n\|] + \frac{CMx}{n},$$

where the first inequality used the improved version of Talagrand's concentration inequality (see Theorem A.2 of Bartlett et al. (2005)) and the last one used the contraction property of Rademacher complexity for Lipschitz functions (see Theorem 2.2 of Koltchinskii (2011)). Taking $x \approx n\lambda/M$, we have $\sup_{g \in \mathscr{G}} (P_n - P)g_M \leq \lambda/4.$

Thus

$$P(\|\sum_{i} \nu_i \mathcal{X}_i / n\| \le \lambda/4 + |\sup_{g \in \mathcal{G}} (P_n - P)(g - g_M)|) \ge 1 - e^{-Cn\lambda/M}.$$

On the other hand, with probability $1 - nP(||\mathcal{X}|| > M)$,

$$|\sup_{g\in\mathscr{G}}(P_n-P)(g-g_M)|\leq \sup_{g\in\mathscr{G}}|P(g-g_M)|\leq \mathbb{E}[\|\boldsymbol{\mathcal{X}}\|^2]P(\|\boldsymbol{\mathcal{X}}\|>M)\leq \lambda/4.$$

Then,

$$P(\|\sum_{i} L'_{h}(u_{i}^{*})y_{i}\boldsymbol{\mathcal{X}}_{i}/n\| \leq \lambda/2) \geq 1 - e^{-Cn\lambda/M} - nP(\|\boldsymbol{\mathcal{X}}\| > M).$$

Proof of Theorem 1: By definition of the RSSTM estimator in Section 2.2, we have

$$\sum_{i=1}^{n} L_h(y_i(\widehat{\beta} + \langle \widehat{\boldsymbol{\mathcal{B}}}, \boldsymbol{\mathcal{X}}_i \rangle)) + n\lambda \|\widehat{\boldsymbol{\mathcal{B}}}\|_* \le \sum_{i=1}^{n} L_h(y_i(\beta^* + \langle \boldsymbol{\mathcal{B}}^*, \boldsymbol{\mathcal{X}}_i \rangle)) + n\lambda \|\boldsymbol{\mathcal{B}}^*\|_*.$$
(S2.7)

On the other hand, using the convexity of smoothed loss function L_h and the inequality $\langle \mathcal{A}, \mathcal{B} \rangle \leq \|\mathcal{A}\| \|\mathcal{B}\|_*$, Assumption 3 and the result of Lemma S2.2, we have

$$\sum_{i=1}^{n} \{ L_h(y_i(\widehat{\beta} + \langle \widehat{\mathcal{B}}, \mathcal{X}_i \rangle)) - L_h(y_i(\beta^* + \langle \mathcal{B}^*, \mathcal{X}_i \rangle)) \} \ge \sum_{i=1}^{n} L'_h(u_i^*) y_i(\langle \widehat{\mathcal{B}} - \mathcal{B}^*, \mathcal{X}_i \rangle + \widehat{\beta} - \beta^*) \\
\ge - \| \sum_{i=1}^{n} L'_h(u_i^*) y_i \mathcal{X}_i \| \| \widehat{\mathcal{B}} - \mathcal{B}^* \|_* - | \sum_{i=1}^{n} L'_h(u_i^*) y_i | | \widehat{\beta} - \beta | \ge - \frac{n\lambda}{2} (\| \widehat{\mathcal{B}} - \mathcal{B}^* \|_* + | \widehat{\beta} - \beta^* |),$$

where we use that $n\lambda \geq C\sqrt{n\log n} \geq 2|\sum_{i=1}^n L_h'(u_i^*)y_i|$ with high probability by Hoeffding's inequality. Recall (S2.7),

$$\sum_{i=1}^{n} L_h(y_i(\widehat{\beta} + \langle \widehat{\boldsymbol{\mathcal{B}}}, \boldsymbol{\mathcal{X}}_i \rangle)) - \sum_{i=1}^{n} L_h(y_i(\beta^* + \langle \boldsymbol{\mathcal{B}}^*, \boldsymbol{\mathcal{X}}_i \rangle)) \le n\lambda(\|\boldsymbol{\mathcal{B}}^*\|_* - \|\widehat{\boldsymbol{\mathcal{B}}}\|_*),$$

which yields that

$$-\frac{\lambda}{2}(\|\widehat{\boldsymbol{\mathcal{B}}} - \boldsymbol{\mathcal{B}}^*\|_* + |\widehat{\beta} - \beta^*|) \le \lambda(\|\boldsymbol{\mathcal{B}}^*\|_* - \|\widehat{\boldsymbol{\mathcal{B}}}\|_*).$$

Let $\widehat{\mathcal{D}} = \widehat{\mathcal{B}} - \mathcal{B}^*$ and $\widehat{\mathcal{D}}^{\perp}$ be the projection of $\widehat{\mathcal{D}}$ on \mathscr{U}^{\perp} , and $\widehat{\mathcal{D}}^0 = \widehat{\mathcal{D}} - \widehat{\mathcal{D}}^{\perp}$. Also, denote $\widehat{\delta} = \widehat{\beta} - \beta^*$. Thus, by the decomposability of t-TNN, the above inequality implies

$$\begin{split} & -\frac{1}{2}(\|\widehat{\mathcal{D}}^{0}\|_{*} + \|\widehat{\mathcal{D}}^{\perp}\|_{*} + |\widehat{\delta}|) \leq \|\mathcal{B}^{*}\|_{*} - \|\mathcal{B}^{*} + \widehat{\mathcal{D}}^{\perp}\|_{*} + \|\widehat{\mathcal{D}}^{0}\|_{*} \\ & \leq \|\mathcal{B}^{*}\|_{*} - \|\mathcal{B}_{1}^{*} + \widehat{\mathcal{D}}^{\perp}\|_{*} + \|\widehat{\mathcal{D}}^{0}\|_{*} = -\|\widehat{\mathcal{D}}^{\perp}\|_{*} + \|\widehat{\mathcal{D}}^{0}\|_{*}, \end{split}$$

which yields

$$\|\widehat{\boldsymbol{\mathcal{D}}}^{\perp}\|_{*} \leq 3\|\widehat{\boldsymbol{\mathcal{D}}}^{0}\|_{*} + |\widehat{\delta}|.$$

That is, $(\widehat{\delta}, \widehat{\mathcal{D}}) \in \mathscr{H} := \{(\delta, \mathcal{D}) : \|\mathcal{D}^{\perp}\|_{*} \leq 3\|\mathcal{D}^{0}\|_{*} + |\delta|\}$. Then, by Lemma S2.1, we have, with probability $1 - e^{-Cn\lambda/M} - nP(\|\mathcal{X}\| > M)$,

$$\inf_{(\delta, \mathcal{D}) \in \mathscr{T}} \left\{ \frac{1}{n} \sum_{i=1}^{n} L_h(y(\beta^* + \delta + \langle \mathcal{B}^* + \mathcal{D}, \mathcal{X} \rangle)) - \frac{1}{n} \sum_{i=1}^{n} L_h(y(\beta^* + \langle \mathcal{B}^*, \mathcal{X} \rangle)) - \mathbb{E}[L_h(y(\beta^* + \delta + \langle \mathcal{B}^* + \mathcal{D}, \mathcal{X} \rangle))] + \mathbb{E}[L_h(y(\beta^* + \langle \mathcal{B}^*, \mathcal{X} \rangle))] \right\}$$

$$\geq -C\lambda t \sqrt{r}.$$
(S2.8)

Suppose $\mathbb{E}[(\delta + \langle \mathcal{X}, \mathcal{D} \rangle)^2] \geq t^2$ for some t > 0, by (S2.7), then we have

$$\inf_{\substack{(\delta, \mathcal{D}) \in \mathcal{H} \\ \mathbb{E}[(\delta + \langle \boldsymbol{x}, \boldsymbol{\mathcal{D}} \rangle)^2] \ge t^2}} \left\{ \frac{1}{n} \sum_{i=1}^n L_h(y(\beta^* + \delta + \langle \boldsymbol{\mathcal{B}}^* + \boldsymbol{\mathcal{D}}, \boldsymbol{\mathcal{X}} \rangle)) - \frac{1}{n} \sum_{i=1}^n L_h(y(\beta^* + \langle \boldsymbol{\mathcal{B}}^*, \boldsymbol{\mathcal{X}} \rangle)) + \lambda \|\boldsymbol{\mathcal{B}}^* + \boldsymbol{\mathcal{D}}\|_* - \lambda \|\boldsymbol{\mathcal{B}}^*\|_* \right\} < 0,$$

and by the convexity of the objective function,

$$\inf_{\substack{(\delta, \mathcal{D}) \in \mathcal{H} \\ \mathbb{E}[(\delta + \langle \mathcal{X}, \mathcal{D} \rangle)^2] = t^2}} \left\{ \frac{1}{n} \sum_{i=1}^n L_h(y(\beta^* + \delta + \langle \mathcal{B}^* + \mathcal{D}, \mathcal{X} \rangle)) - \frac{1}{n} \sum_{i=1}^n L_h(y(\beta^* + \langle \mathcal{B}^*, \mathcal{X} \rangle)) + \lambda \|\mathcal{B}^* + \mathcal{D}\|_* - \lambda \|\mathcal{B}^*\|_* \right\} < 0.$$
(S2.9)

Thus, for some $(\delta, \mathcal{D}) \in \mathscr{H}$ satisfying $\mathbb{E}[(\delta + \langle \mathcal{X}, \mathcal{D} \rangle)^2] = t^2$,

$$\mathbb{E}[L_{h}(y(\beta^{*} + \delta + \langle \mathcal{B}^{*} + \mathcal{D}, \mathcal{X}\rangle))] - \mathbb{E}[L_{h}(y(\beta^{*} + \langle \mathcal{B}^{*}, \mathcal{X}\rangle))]$$

$$\leq \frac{1}{n} \sum_{i=1}^{n} L_{h}(y(\beta^{*} + \delta + \langle \mathcal{B}^{*} + \mathcal{D}, \mathcal{X}\rangle)) - \frac{1}{n} \sum_{i=1}^{n} L_{h}(y(\beta^{*} + \langle \mathcal{B}^{*}, \mathcal{X}\rangle)) + C\lambda t\sqrt{r}$$

$$\leq \lambda(\|\mathcal{B}^{*}\|_{*} - \|\mathcal{B}^{*} + \mathcal{D}^{\perp}\|_{*} + \|\mathcal{D}^{0}\|_{*}) + C\lambda t\sqrt{r}$$

$$\leq C\lambda(\|\mathcal{D}^{0}\|_{*} + t\sqrt{r}) \leq C\lambda(\sqrt{r}\|\mathcal{D}\|_{F} + t\sqrt{r}) \leq C\lambda t\sqrt{r},$$

where the first two inequalities are based on (S2.8) and (S2.9), the 3rd inequality uses the decomposibility of t-TNN, the penultimate inequality uses the low rank, and the last one applies $\|\mathcal{D}\|_F \leq Ct$ by $\mathbb{E}[(\delta + \langle \mathcal{X}, \mathcal{D} \rangle)^2] = t^2$ and **Assumption** 2, which suggests the first result. Then, under **Assumption** 2, we can immediately

get the second result. By $\|\widehat{\boldsymbol{\mathcal{B}}} - \boldsymbol{\mathcal{B}}^*\|_* \le \|(\widehat{\boldsymbol{\mathcal{B}}} - \boldsymbol{\mathcal{B}}^*)^\perp\|_* + \|(\widehat{\boldsymbol{\mathcal{B}}} - \boldsymbol{\mathcal{B}}^*)^0\|_* \le C(\|(\widehat{\boldsymbol{\mathcal{B}}} - \boldsymbol{\mathcal{B}}^*)^0\|_* + |\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*|) \le C\sqrt{r}(\|\widehat{\boldsymbol{\mathcal{B}}} - \boldsymbol{\mathcal{B}}^*\|_F + |\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*|)$ due to $\widehat{\boldsymbol{\mathcal{B}}} - \boldsymbol{\mathcal{B}}^* \in \mathscr{H}$, we have the last result.

Lemma S2.3 *Under Assumption 2, we have*

$$\mathbb{E}[\|\sum_i L_h'(u_i^*)y_i\boldsymbol{\mathcal{X}}_i\|] \leq C\sqrt{n(I_1I_3 \vee I_2I_3)} \text{ and } \mathbb{E}[\|\sum_i \sigma_i\boldsymbol{\mathcal{X}}_i\|] \leq C\sqrt{n(I_1I_3 \vee I_2I_3)}.$$

Proof of Lemma S2.3: Similar to Lemma S2.2, let $\nu_i = L'_h(u_i^*)y_i$ which is bounded. From Appendix S1, it can be seen that $\|\sum_i \nu_i \boldsymbol{\mathcal{X}}_i\| = \|\sum_i \nu_i \overline{\boldsymbol{\mathcal{X}}}_i\|_{op}$ where $\overline{\boldsymbol{\mathcal{X}}}_i \in \mathbb{R}^{I_1 I_3 \times I_2 I_3}$. It is trivial that $\overline{\boldsymbol{\mathcal{X}}}_i, i = 1, \cdots, n$, are sub-Gaussian matrices via the sub-Gaussianity of $\boldsymbol{x}_i = \text{vec}(\boldsymbol{\mathcal{X}}_i)$. With the boundedness of ν_i , let $\overline{\boldsymbol{x}}_i = \text{vec}(\overline{\boldsymbol{\mathcal{X}}}_i)$ and $\nu_i \overline{\boldsymbol{x}}_i$ is also sub-Gaussian. With the covering of unit spheres \boldsymbol{u}_i and \boldsymbol{v}_j for $\sum_i \nu_i \overline{\boldsymbol{\mathcal{X}}}_i$, defined in the proof of Lemma S2.4, we have, for any t > 0,

$$P(\|\sum_{i} \nu_{i} \overline{\boldsymbol{X}}_{i}\|_{op} > 2a) \leq \sum_{i,j} P(\boldsymbol{u}_{i}^{\mathsf{T}}(\sum_{i} \nu_{i} \overline{\boldsymbol{X}}_{i})\boldsymbol{v}_{j} > a) \leq \sum_{i,j} P((\boldsymbol{v}_{j} \otimes \boldsymbol{u}_{i})^{\mathsf{T}} \operatorname{vec}(\sum_{i} \nu_{i} \overline{\boldsymbol{X}}_{i}))$$

$$\leq \sum_{i,j} e^{-ta} \mathbb{E}[e^{t(\boldsymbol{v}_{j} \otimes \boldsymbol{u}_{i})^{\mathsf{T}} \operatorname{vec}(\sum_{i} \nu_{i} \overline{\boldsymbol{X}}_{i})}] \leq 20^{I_{1}I_{3} + I_{2}I_{3}} e^{-ta + Cnt^{2}} \leq 20^{I_{1}I_{3} + I_{2}I_{3}} e^{-ta + Cnt^{2}}.$$

Taking $t \approx a/n$ and $a = C\sqrt{n((I_1I_3 \vee I_2I_3) - \log \tau)}$ for any $\tau > 0$, we have

$$P(\|\sum_{i} \nu_{i} \overline{X}_{i}\|_{op} > C\sqrt{n((I_{1}I_{3} \vee I_{2}I_{3}) - \log \tau)}) \leq \tau.$$

Denote $t = \log(1/\tau)$ and $a = C\sqrt{n(I_1I_3 \vee I_2I_3)}$, which implies $\|\sum_i \nu_i \overline{X}_i\|_{op} > at^{1/2}$ with probability at least e^{-t} . Thus, we have

$$\mathbb{E}[\|\sum_{i} \nu_{i} \overline{X}_{i}\|_{op}] \leq \sum_{t=1}^{\infty} \mathbb{E}[\|\sum_{i} \nu_{i} \overline{X}_{i}\|_{op} \cdot I\{a(t-1)^{1/2} < \|\sum_{i} \nu_{i} \overline{X}_{i}\|_{op} \leq at^{1/2}\}]$$

$$\leq \sum_{t=1}^{\infty} at^{1/2} P(a(t-1)^{1/2} < \|\sum_{i} \nu_{i} \overline{X}_{i}\|_{op} \leq at^{1/2}) \leq \sum_{t=1}^{\infty} at^{1/2} e^{-t+1} = C\sqrt{n(I_{1}I_{3} \vee I_{2}I_{3})}.$$

We can get $\mathbb{E}[\|\sum_i \sigma_i \mathcal{X}_i\|] \leq C \sqrt{n(I_1 I_3 \vee I_2 I_3)}$ in the same way.

Lemma S2.4 Under Assumption 2, we have $P(\|\mathcal{X}\| > C\sqrt{(I_1I_3 \vee I_2I_3)\log n}) \leq Ce^{-C(I_1I_3 \vee I_2I_3)\log n}$.

Proof of Lemma S2.4: We have $\|X\| = \|\overline{X}\|_{op}$ where $\overline{X} \in \mathbb{R}^{I_1 I_3 \times I_2 I_3}$. Let $\{u_i\}_{i=1}^{M_1}$ and $\{v_j\}_{i=1}^{M_2}$ be the 1/4 covering of the unit sphere in $\mathbb{R}^{I_1 I_3}$ and $\mathbb{R}^{I_2 I_3}$, for \overline{X} , respectively, with $M_1 \leq 20^{I_1 I_3}$ and $M_2 \leq 20^{I_2 I_3}$ referred from Lemma 2.5 in Geer (2000). Thus, for any u, v if $\|u\| = \|v\| = 1$, there exists u_i, v_j in the covering satisfy $\|u - u_i\| \leq 1/4$, $\|v - v_j\| \leq 1/4$, and then

$$\boldsymbol{u}^\mathsf{T}\overline{\boldsymbol{X}}\boldsymbol{v} = \boldsymbol{u}^\mathsf{T}\overline{\boldsymbol{X}}(\boldsymbol{v} - \boldsymbol{v}_j) + (\boldsymbol{u} - \boldsymbol{u}_i)^\mathsf{T}\overline{\boldsymbol{X}}\boldsymbol{v}_j \leq \frac{1}{2}\|\overline{\boldsymbol{X}}\|_{op} + \boldsymbol{u}_i^\mathsf{T}\overline{\boldsymbol{X}}\boldsymbol{v}_j.$$

Then, we have

$$\|\overline{\boldsymbol{X}}\|_{op} = \sup_{\|\boldsymbol{u}\| = \|\boldsymbol{v}\| = 1} \boldsymbol{u}^{\mathsf{T}} \overline{\boldsymbol{X}} \boldsymbol{v} \leq \frac{1}{2} \|\overline{\boldsymbol{X}}\|_{op} + \max_{\boldsymbol{u}_i, \boldsymbol{v}_j} \boldsymbol{u}_i^{\mathsf{T}} \overline{\boldsymbol{X}} \boldsymbol{v}_j,$$

which yields $\|\overline{\boldsymbol{X}}\|_{op} \leq 2 \max_{\boldsymbol{u}_i, \boldsymbol{v}_j} \boldsymbol{u}_i^{\mathsf{T}} \overline{\boldsymbol{X}} \boldsymbol{v}_j$. Thus, for any t > 0, due to the sub-Gaussianity of $\boldsymbol{u}_i^{\mathsf{T}} \overline{\boldsymbol{X}} \boldsymbol{v}_j = (\boldsymbol{v}_j \otimes \boldsymbol{u}_i)^{\mathsf{T}} \text{vec}(\overline{\boldsymbol{X}})$, we have

$$\begin{split} & \mathbb{E}[e^{t\|\overline{\boldsymbol{X}}\|_{op}}] \leq \mathbb{E}\left[e^{2t\max_{\boldsymbol{u}_i,\boldsymbol{v}_j}\boldsymbol{u}_i^{\mathsf{T}}\overline{\boldsymbol{X}}\boldsymbol{v}_j}\right] \leq \mathbb{E}\left[\sum_{i,j}e^{2t\boldsymbol{u}_i^{\mathsf{T}}\overline{\boldsymbol{X}}\boldsymbol{v}_j}\right] \\ & \leq 20^{I_1I_3+I_2I_3}\max_{\boldsymbol{u}_i,\boldsymbol{v}_j}\mathbb{E}[e^{2t\boldsymbol{u}_i^{\mathsf{T}}\overline{\boldsymbol{X}}\boldsymbol{v}_j}] \leq 20^{I_1I_3+I_2I_3}e^{Ct^2}. \end{split}$$

By Markov's inequality,

$$P(\|\overline{\boldsymbol{X}}\|_{op} > s) \le e^{-ts} 20^{I_1 I_3 + I_2 I_3} e^{Ct^2} \le Ce^{-Cs^2}.$$

Taking $s = C\sqrt{(I_1I_3 \vee I_2I_3)\log n}$ completes the proof.

Proof of Theorem 2: By the optimality of $\widehat{\mathcal{B}}$, there exists $\mathcal{Z} \in \partial \|\widehat{\mathcal{B}}\|_*$ such that, $\frac{1}{n} \sum_{i=1}^n L'_h(y_i(\hat{\beta} + \langle \widehat{\mathcal{B}}, \mathcal{X}_i \rangle)) y_i \mathcal{X}_i + \lambda \mathcal{Z} = 0.$ Thus, we have

$$\frac{1}{n} \sum_{i=1}^{n} L_h'(\hat{v}_i) y_i \overline{X}_i + \lambda \overline{Z} = 0,$$
 (S2.10)

where $\hat{v}_i = y_i(\hat{\beta} + \frac{1}{I_3}\langle \overline{\widehat{B}}, \overline{X_i} \rangle)$ and $\overline{Z} \in \partial \|\overline{\widehat{B}}\|_*$. Let $\operatorname{rank}_t(\widehat{\overline{B}}) = \hat{r}$, $\operatorname{rank}(\overline{\widehat{B}}) = \hat{r}' \leq I_3 \hat{r}$ and its SVD given by $\overline{\widehat{B}} = USV^{\mathsf{T}}$, $U \in \mathbb{R}^{I_1 I_3 \times \hat{r}'}$, $V \in \mathbb{R}^{I_2 I_3 \times \hat{r}'}$. Using (2.1) in Koltchinskii et al. (2011), it leads to

$$\partial \|\overline{\widehat{B}}\|_* = \{ UV^\mathsf{T} + U^\perp W(V^\perp)^\mathsf{T} : \|W\|_{op} \le 1 \},$$

where $U^{\perp} \in \mathbb{R}^{I_1 I_3 \times (I_1 I_3 - \hat{r}')}$ and $V^{\perp} \in \mathbb{R}^{I_2 I_3 \times (I_2 I_3 - \hat{r}')}$ are orthogonal matrices whose columns are orthogonal to those in U and V, respectively. Let $u_j, v_j, j = 1, \dots, \hat{r}'$ be the columns of U and V. Pre(post)-multiplying u_j and v_j to (S2.10) yields,

$$\frac{1}{n} \sum_{i=1}^{n} L_h'(\hat{v}_i) y_i \boldsymbol{u}_j^{\mathsf{T}} \overline{\boldsymbol{X}_i} \boldsymbol{v}_j = -\lambda, \ j = 1, \cdots, \hat{r}'.$$
 (S2.11)

According to the optimality condition of $(\beta^*, \mathbf{\mathcal{B}}^*)$, $\mathbb{E}[L_h'(v_i^*)y_i\mathbf{u}_j^{\mathsf{T}}\overline{\mathbf{X}_i}\mathbf{v}_j] = 0$. We write

$$\frac{1}{n} \sum_{i=1}^{n} L'_{h}(\hat{v}_{i}) y_{i} \boldsymbol{u}_{j}^{\mathsf{T}} \overline{\boldsymbol{X}_{i}} \boldsymbol{v}_{j}
= (P_{n} - P)[(L'_{h}(\hat{v}_{i}) - L'_{h}(v_{i}^{*})) y_{i} \boldsymbol{u}_{j}^{\mathsf{T}} \overline{\boldsymbol{X}_{i}} \boldsymbol{v}_{j}] + \mathbb{E}[(L'_{h}(\hat{v}_{i}) - L'_{h}(v_{i}^{*})) y_{i} \boldsymbol{u}_{j}^{\mathsf{T}} \overline{\boldsymbol{X}_{i}} \boldsymbol{v}_{j}]
+ (P_{n} - P)(L'_{h}(v_{i}^{*}) y_{i} \boldsymbol{u}_{j}^{\mathsf{T}} \overline{\boldsymbol{X}_{i}} \boldsymbol{v}_{j}) + \mathbb{E}[L'_{h}(v_{i}^{*}) y_{i} \boldsymbol{u}_{j}^{\mathsf{T}} \overline{\boldsymbol{X}_{i}} \boldsymbol{v}_{j}]
:= g_{1j} + g_{2j} + g_{3j}.$$

Let $\mathbf{g}_k = (g_{kj}, \dots, g_{k\hat{r}'})^\mathsf{T}$ for k = 1, 2, 3 and $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_{\hat{r}'})^\mathsf{T} \in \mathbb{R}^{\hat{r}'}$ be any unit vector and $\boldsymbol{\gamma} = (\boldsymbol{v}_1 \otimes \boldsymbol{u}_1, \dots, \boldsymbol{v}_{\hat{r}'} \otimes \boldsymbol{u}_{\hat{r}'})\boldsymbol{\alpha}$, we have $\|\boldsymbol{\gamma}\|_2^2 = \sum_{j=1}^{\hat{r}'} \alpha_j^2 \|\boldsymbol{v}_j \otimes \boldsymbol{u}_j\|_2^2 = 1$. Denote $\mathscr{D} = \{\overline{\boldsymbol{B}} \in \mathbb{R}^{I_1 I_3 \times I_2 I_3} : \|\overline{\boldsymbol{B}} - \overline{\boldsymbol{B}}^*\|_F \le d_n := C\sqrt{\hat{r}'(I_1 I_3 \vee I_2 I_3)/n}\}$. Due to the Lipschitz continuity of L_h' and the sub-Gaussian property of X, then we have

$$\begin{aligned} \|\mathbf{g}_{2}\|_{2} &= \sup_{\|\alpha\|_{2}=1} \boldsymbol{\alpha}^{\mathsf{T}} \mathbf{g}_{2} \leq \sup_{\|\boldsymbol{\gamma}\|_{2}=1, \overline{\boldsymbol{B}} \in \mathscr{D}} C \mathbb{E}[|\langle \overline{\boldsymbol{X}_{i}}, \overline{\boldsymbol{B}} - \overline{\boldsymbol{B}^{*}} \rangle + (\beta - \beta^{*})||\boldsymbol{\gamma}^{\mathsf{T}} \mathrm{vec}(\overline{\boldsymbol{X}_{i}})|] \\ &\leq \sup_{\|\boldsymbol{\gamma}\|_{2}=1, \overline{\boldsymbol{B}} \in \mathscr{D}} C \sqrt{\mathbb{E}|\langle \overline{\boldsymbol{X}_{i}}, \overline{\boldsymbol{B}} - \overline{\boldsymbol{B}^{*}} \rangle + (\beta - \beta^{*})|^{2}} \sqrt{\mathbb{E}|\boldsymbol{\gamma}^{\mathsf{T}} \mathrm{vec}(\overline{\boldsymbol{X}_{i}})|^{2}} \\ &\leq \sup_{\overline{\boldsymbol{B}} \in \mathscr{D}} C \Big(\|\overline{\boldsymbol{B}} - \overline{\boldsymbol{B}^{*}}\|_{F} + |\beta - \beta^{*}| \Big) \leq C d_{n}. \end{aligned}$$

By similar arguments as Lemma S2.5, with high probability,

$$\|\mathbf{g}_1\|_2 \le C d_n \sqrt{\frac{\hat{r}'(I_1 I_3 \vee I_2 I_3) \log n}{n}} + C d_n \frac{(\hat{r}')^{3/2} (I_1 I_3 \vee I_2 I_3)^2 (\log n)^2}{n}.$$

By similar arguments as Lemma S2.4,

$$\|\mathbf{g}_3\|_2 \le C\sqrt{\frac{(I_1I_3 \lor I_2I_3)\log n}{n}},$$

with high probability. Hence,

$$\lambda \sqrt{\hat{r}'} \leq \|\mathbf{g}_1\|_2 + \|\mathbf{g}_2\|_2 + \|\mathbf{g}_3\|_2$$

$$\leq C d_n \sqrt{\frac{\hat{r}'(I_1 I_3 \vee I_2 I_3) \log n}{n}} + C d_n \frac{(\hat{r}')^{3/2} (I_1 I_3 \vee I_2 I_3)^2 (\log n)^2}{n} + C d_n + C \sqrt{\frac{(I_1 I_3 \vee I_2 I_3) \log n}{n}},$$

which implies that $\hat{r} \leq Cr$ by taking $\lambda = C\sqrt{(I_1I_3 \vee I_2I_3)/n}$.

Lemma S2.5 Assume
$$\mathbb{E}(\langle \mathcal{X}, \mathcal{D} \rangle + \delta)^4 \leq C(\|\mathcal{D}\|_F^4 + |\delta|^4)$$
 for any $\mathcal{D} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$, $\delta \in \mathbb{R}^{I_1 \times I_2 \times I_3}$

 \mathbb{R} and some generic positive constant C, uniformly over $\mathscr{A} = \{(\beta, \mathcal{B}) : |\beta - \beta^*| + \|\mathcal{B} - \mathcal{B}^*\|_F \le d_n/\sqrt{I_3}, \operatorname{rank}_t(\mathcal{B}) \le Cr\},$

$$\sup_{(\beta, \mathcal{B}) \in \mathscr{A}} \left\| \frac{1}{n} \sum_{i=1}^{n} (L'_h(u_i) - L'_h(u_i^*)) y_i \mathcal{X}_i - \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^{n} (L'_h(u_i) - L'_h(u_i^*)) y_i \mathcal{X}_i \right] \right\|$$

$$\leq C d_n \left(\sqrt{\frac{r(I_1 I_3 \vee I_2 I_3) \log n}{n}} + \frac{r^{3/2} (I_1 I_3 \vee I_2 I_3)^2 (\log n)^2}{n} \right).$$

with probability at least $1 - n^{-C}$.

Proof of Lemma S2.5: We consider $\mathscr{A}' = \{(\beta, \overline{B}) : |\beta - \beta^*| + ||\overline{B} - \overline{B}^*||_F \le d_n, \operatorname{rank}_t(\overline{B}) \le Cr'\}$ with $r' \le I_3r$. It can be verified that $(\beta, \overline{B}) \in \mathscr{A}'$ if $(\beta, \mathcal{B}) \in \mathscr{A}$. Let $q_i(\beta, \overline{B}) = ((L'_h(v_i) - L'_h(v_i^*))y_i\overline{X_i} - \mathbb{E}[(L'_h(v_i) - L'_h(v_i^*))y_i\overline{X_i}]$. Denote $\mathcal{N}_{\mathscr{A}'}$ as the $d_n n^{-M}$ -covering of \mathscr{A}' for a sufficiently large M > 0, then

$$\sup_{(\beta,\overline{B})\in\mathscr{A}'} \left\| \frac{1}{n} \sum_{i=1}^{n} q_{i}(\beta,\overline{B}) \right\|_{op}$$

$$\leq \max_{(\beta',\overline{B}')\in\mathscr{N}_{\mathscr{A}'}} \left\| \frac{1}{n} \sum_{i=1}^{n} q_{i}(\beta',\overline{B}') \right\|_{op}$$

$$+ \max_{(\beta',\overline{B}')\in\mathscr{N}_{\mathscr{A}'}} \sup_{|\beta-\beta^{*}|+\|\overline{B}-\overline{B}^{*}\|_{F}\leq d_{n}n^{-M}} \left\| \frac{1}{n} \sum_{i=1}^{n} \{q_{i}(\beta,\overline{B}) - q_{i}(\beta',\overline{B}')\} \right\|_{op}$$

$$:= S_{1} + S_{2}.$$

Next, we separately bound S_1 and S_2 . For any covering of unit spheres $\{u_i\}_{i=1}$, $\{v_j\}_{j=1}$ and the Lipschitz continuity, using Cauchy-Schwartz inequality and assumption

$$\mathbb{E}(\langle \boldsymbol{\mathcal{X}}, \boldsymbol{\mathcal{D}} \rangle + \delta)^4 \leq C(\|\boldsymbol{\mathcal{D}}\|_F^4 + |\delta|^4)$$
, we have

$$\mathbb{E}\left[\left(L_h'(v_i') - L_h'(v_i^*)\right)y_i \boldsymbol{u}_i^{\mathsf{T}} \overline{\boldsymbol{X}_i} \boldsymbol{v}_j\right]^2 \leq C \sqrt{\mathbb{E}\left[\langle \overline{\boldsymbol{X}_i}, \overline{\boldsymbol{B}} - \overline{\boldsymbol{B}}^* \rangle + |\beta' - \beta^*|\right]^4} \sqrt{\mathbb{E}\left[y_i \boldsymbol{u}_i^{\mathsf{T}} \overline{\boldsymbol{X}_i} \boldsymbol{v}_j\right]^4} \\
\leq C(|\beta - \beta^*|^2 + ||\overline{\boldsymbol{B}} - \overline{\boldsymbol{B}}^*||_F^2) \leq C d_n^2.$$

Moreover,

$$\begin{aligned} \left| (L_h'(v_i') - L_h'(v_i^*)) y_i \boldsymbol{u}_i^\mathsf{T} \overline{\boldsymbol{X}}_i \boldsymbol{v}_j \right| &\leq C (\max_{1 \leq i \leq n} \| \overline{\boldsymbol{X}}_i \|_{op}) \| \overline{\boldsymbol{B}} - \overline{\boldsymbol{B}}^* \|_* \\ &\leq C \sqrt{r'} d_n \sqrt{(I_1 I_3 \vee I_2 I_3) \log n}, \end{aligned}$$

where the last result uses that $\max_{1 \le i \le n} \|\overline{X_i}\|_{op} \le C\sqrt{(I_1I_2 \lor I_2I_3)\log n}$ with high probability by similar arguments as that in Lemma S2.4. Hence, by Bernstein's inequality, with probability at least $1 - 2e^{-z}$ for any z > 0, we have

$$\left| \frac{1}{n} \sum_{i=1}^{n} (L'_h(v'_i) - L'_h(v^*_i)) y_i \boldsymbol{u}_i^{\mathsf{T}} \overline{\boldsymbol{X}}_i \boldsymbol{v}_j - \mathbb{E} \left[(L'_h(v'_i) - L'_h(v^*_i)) y_i \boldsymbol{u}_i^{\mathsf{T}} \overline{\boldsymbol{X}}_i \boldsymbol{v}_j \right] \right| \\
\leq C \sqrt{\frac{z d_n^2}{n}} + C \frac{z \sqrt{r'} d_n (I_1 I_3 \vee I_2 I_3)}{n} := \psi(z),$$

By Lemma 5.2 in Vershynin (2011), the covering number satisfies $|\mathcal{N}_{\mathscr{A}'}| \leq (1 +$

 $(2n^M)^{Cr'(I_1I_3+I_2I_3)}$. Then we have

$$P(S_{1} \geq \psi(z)) = P\left(\max_{(\beta', \overline{B}') \in \mathcal{N}_{\mathscr{A}'}} \left\| \frac{1}{n} \sum_{i=1}^{n} q_{i}(\beta', \overline{B}') \right\|_{op} \geq \psi(z)\right)$$

$$\leq Cn^{CMr'(I_{1}I_{3} \vee I_{2}I_{3})} \max_{(\beta', \overline{B}') \in \mathcal{N}_{\mathscr{A}'}} P\left(\max_{\boldsymbol{u}_{i}, \boldsymbol{v}_{j}} \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{u}_{i}^{\mathsf{T}} q_{i}(\beta', \overline{B}') \boldsymbol{v}_{j} \geq \psi(z)\right)$$

$$\leq Cn^{CMr'(I_{1}I_{3} \vee I_{2}I_{3})} 20^{I_{1}I_{3} \vee I_{2}I_{3}} e^{-z}.$$

Taking $z = Cr'(I_1I_3 \vee I_2I_3)\log n$, with probability at least $1 - n^{-C}$,

$$S_1 \le C d_n \sqrt{\frac{r'(I_1 I_3 \vee I_2 I_3) \log n}{n}} + C d_n \frac{(r')^{3/2} (I_1 I_3 \vee I_2 I_3)^2 (\log n)^2}{n}.$$
 (S2.12)

Next we bound S_2 . Denote

$$w(\overline{\boldsymbol{X}_i}) = \sup_{|\beta - \beta^*| + ||\overline{\boldsymbol{B}} - \overline{\boldsymbol{B}}^*||_F \le d_n n^{-M}} [q_i(\beta, \overline{\boldsymbol{B}}) - q_i(\beta', \overline{\boldsymbol{B}}')].$$

Similarly, it leads to

$$\begin{aligned} |\boldsymbol{u}_{i}^{\mathsf{T}}w(\overline{\boldsymbol{X}_{i}})\boldsymbol{v}_{j}| &\leq \left| \sup_{\|\beta-\beta^{*}\|+\|\overline{\boldsymbol{B}}-\overline{\boldsymbol{B}}^{*}\|_{F} \leq d_{n}n^{-M}} \left| \langle \overline{\boldsymbol{X}_{i}}, \overline{\boldsymbol{B}}-\overline{\boldsymbol{B}}^{*} \rangle + (\beta-\beta^{*}) \right| \boldsymbol{u}_{i}^{\mathsf{T}} \overline{\boldsymbol{X}_{i}} \boldsymbol{v}_{j} \right| \\ &\leq C\sqrt{r'} d_{n} n^{-M} (I_{1}I_{2} \vee I_{2}I_{3}) \mathrm{log} n, \end{aligned}$$

which implies

$$S_2 \le C \frac{\sqrt{r'} d_n (I_1 I_3 \vee I_2 I_3) \log n}{n^M},$$
 (S2.13)

with high probability. Therefore, for sufficiently large M, together with (S2.12)

and (S2.13),

$$\sup_{(\beta, \overline{\boldsymbol{B}}) \in \mathscr{A}'} \left\| \frac{1}{n} \sum_{i=1}^{n} q_i(\beta, \overline{\boldsymbol{B}}) \right\|_{om} \le C d_n \sqrt{\frac{r'(I_1 I_3 \vee I_2 I_3) \log n}{n}} + C d_n \frac{(r')^{3/2} (I_1 I_3 \vee I_2 I_3)^2 (\log n)^2}{n},$$

with probability at least $1 - n^{-C}$, which yields

$$\sup_{(\beta, \mathcal{B}) \in \mathscr{A}} \left\| \frac{1}{n} \sum_{i=1}^{n} (L'_h(u_i) - L'_h(u_i^*)) y_i \mathcal{X}_i - \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^{n} (L'_h(u_i) - L'_h(u_i^*)) y_i \mathcal{X}_i \right] \right\|$$

$$\leq C d_n \sqrt{\frac{r(I_1 I_3 \vee I_2 I_3) \log n}{n}} + C d_n \frac{r^{3/2} (I_1 I_3 \vee I_2 I_3)^2 (\log n)^2}{n}.$$

Proof of Lemma 1: We first prove (2.6). By the optimality of \mathcal{B}^{k+1} ,

$$\frac{1}{n} \sum_{i=1}^{n} L_{h}(u_{i}^{k}) + \frac{1}{n} \sum_{i=1}^{n} L'_{h}(u_{i}^{k}) y_{i} \langle \mathbf{\mathcal{B}}^{k+1} - \mathbf{\mathcal{B}}^{k}, \mathbf{\mathcal{X}}_{i} \rangle + \frac{\mu}{2} \|\mathbf{\mathcal{B}}^{k+1} - \mathbf{\mathcal{B}}^{k}\|_{F}^{2} + \lambda \|\mathbf{\mathcal{B}}^{k+1}\|_{*}$$

$$\leq \frac{1}{n} \sum_{i=1}^{n} L_{h}(u_{i}^{k}) + \lambda \|\mathbf{\mathcal{B}}^{k}\|_{*}.$$
(S2.14)

On the other hand, by the descent lemma (Beck and Teboulle, 2009, Lemma 2.1),

we get

$$\frac{1}{n} \sum_{i=1}^{n} L_{h}(y_{i}(\beta^{k} + \langle \mathbf{\mathcal{B}}^{k+1}, \mathbf{\mathcal{X}}_{i} \rangle))$$

$$\leq \frac{1}{n} \sum_{i=1}^{n} L_{h}(u_{i}^{k}) + \frac{1}{n} \sum_{i=1}^{n} L'_{h}(u_{i}^{k}) y_{i} \langle \mathbf{\mathcal{B}}^{k+1} - \mathbf{\mathcal{B}}^{k}, \mathbf{\mathcal{X}}_{i} \rangle + \frac{C_{L_{h}}}{2} \|\mathbf{\mathcal{B}}^{k+1} - \mathbf{\mathcal{B}}^{k}\|_{F}^{2}.$$
(S2.15)

Adding (S2.14) and (S2.15), it can be seen that

$$\frac{1}{n} \sum_{i=1}^{n} L_h(y_i(\beta^k + \langle \mathbf{\mathcal{B}}^{k+1}, \mathbf{\mathcal{X}}_i \rangle)) + \lambda \|\mathbf{\mathcal{B}}^{k+1}\|_* + \frac{\mu - C_{L_h}}{2} \|\mathbf{\mathcal{B}}^{k+1} - \mathbf{\mathcal{B}}^k\|_F^2 \\
\leq \frac{1}{n} \sum_{i=1}^{n} L_h(u_i^k) + \lambda \|\mathbf{\mathcal{B}}^k\|_*,$$

which leads to (2.6). Then we prove (2.7). Note that

$$f(\beta^k, \mathbf{\mathcal{B}}^{k+1}) - f(\beta^{k+1}, \mathbf{\mathcal{B}}^{k+1})$$

$$= \frac{1}{n} \sum_{i=1}^n L_h(y_i(\beta^k + \langle \mathbf{\mathcal{B}}^{k+1}, \mathbf{\mathcal{X}}_i \rangle)) - \frac{1}{n} \sum_{i=1}^n L_h(y_i(\beta^{k+1} + \langle \mathbf{\mathcal{B}}^{k+1}, \mathbf{\mathcal{X}}_i \rangle)).$$

By the iteration equation of β^{k+1} ,

$$-\frac{1}{n}\sum_{i=1}^{n}L'_{h}(y_{i}(\beta^{k}+\langle \boldsymbol{\mathcal{B}}^{k+1},\boldsymbol{\mathcal{X}}_{i}\rangle))y_{i}=\frac{1}{\rho}(\beta^{k+1}-\beta^{k}),$$

which implies, with the descent lemma,

$$f(\beta^k, \mathbf{\mathcal{B}}^{k+1}) - f(\beta^{k+1}, \mathbf{\mathcal{B}}^{k+1}) \ge (\frac{1}{\rho} - \frac{C_{L_h}}{2})|\beta^{k+1} - \beta^k|^2 = \gamma_2|\beta^{k+1} - \beta^k|^2.$$

Proof of Theorem 3: (1). By Lemma 1 and the condition on γ , we have

$$f(\mathcal{W}^k) - f(\mathcal{W}^{k+1}) \ge \gamma \{ \|\mathcal{B}^{k+1} - \mathcal{B}^k\|_F^2 + |\beta^{k+1} - \beta^k|^2 \} = \gamma \|\mathcal{W}^{k+1} - \mathcal{W}^k\|_F^2,$$
(S2.16)

which means $f(\mathbf{W}^k)$ is always monotonically decreasing.

(2). Note that $f(\mathbf{W}^k)$ is bounded. Then, the sequence $\{\mathbf{W}^k\}$ is also bounded via the coerciveness of $f(\cdot)$ (the coerciveness of a function refers to its property of tending toward infinity as the input variables approach infinity, i.e., a function $g(x) \to +\infty$ as $||x|| \to +\infty$ where ||x|| denotes the norm (e.g., Euclidean norm) of x). Summing up (S2.16) for all k > 0, we immediately have

$$f(\mathcal{W}^0) \ge \gamma \sum_{k=0}^{+\infty} \|\mathcal{W}^{k+1} - \mathcal{W}^k\|_F,$$

which implies

$$\sum_{k=0}^{+\infty} \|\boldsymbol{\mathcal{W}}^{k+1} - \boldsymbol{\mathcal{W}}^k\|_F \leq \frac{1}{\gamma} f(\boldsymbol{\mathcal{W}}^0) < +\infty,$$

Hence, we get

$$\lim_{k\to+\infty} (\boldsymbol{\mathcal{W}}^{k+1} - \boldsymbol{\mathcal{W}}^k) = 0.$$

(3). We only prove the accumulation point \mathcal{B}^* is a critical point of $f(\beta, \mathcal{B})$ since the same assertion of β is easy to obtain. By the boundedness of $\{\mathcal{B}^k\}$, we

have a subsequence $\{\mathcal{B}^{k_t}\}$ satisfying $\lim_{t\to+\infty}\mathcal{B}^{k_t}=\mathcal{B}^*$, and $\{\mathcal{B}^{k_t}\}$ is bounded. Due to the optimality condition of \mathcal{B}^{k_t+1} , there exists $\mathcal{Z}^{k_t+1}\in\partial\|\mathcal{B}^{k_t+1}\|_*$ such that

$$\lambda \mathbf{Z}^{k_t+1} + \frac{1}{n} \sum_{i=1}^n L'_h(u_i^{k_t}) y_i \mathbf{X}_i + \mu (\mathbf{B}^{k_t+1} - \mathbf{B}^{k_t}) = 0,$$
 (S2.17)

where $\partial \|\mathbf{\mathcal{B}}\|_*$ indicates the subdifferential of t-TNN. Let $t \to +\infty$ in (S2.17), then there exists $\mathbf{\mathcal{Z}}^* \in \partial \|\mathbf{\mathcal{B}}^*\|_*$ such that

$$0 = \lambda \mathbf{Z}^* + \frac{1}{n} \sum_{i=1}^n L'_h(u_i^*) y_i \mathbf{X}_i \in \partial_{\mathbf{B}} f(\mathbf{W}^*),$$

which completes the proof.

S3 Proofs and Lemmas for results in Section 3

Lemma S3.1 Suppose Assumption 1-5 hold, with probability $1 - n^{-C}$, we have

$$\sup_{(\delta, \mathcal{D}) \in \mathscr{T}} |Q_1(\beta, \mathcal{B}) - Q_1(\beta^*, \mathcal{B}^*) - \delta \nabla_{\beta} Q_1(\beta^*, \mathcal{B}^*) - \langle \mathcal{D}, \nabla_{\mathcal{B}} Q_1(\beta^*, \mathcal{B}^*) \rangle$$

$$- \mathbb{E}[L_h(y_i(\beta + \langle \mathcal{B}, \mathcal{X}_i \rangle))] + \mathbb{E}[L_h(y_i(\beta^* + \langle \mathcal{B}^*, \mathcal{X}_i \rangle))]|$$

$$\leq C \left(\frac{b_n \sqrt{r} t I_1 I_2 I_3^2 \log n}{n} + \frac{(b_n \sqrt{r} t)^{3/2} \sqrt{I_1 I_2 I_3^2 \log n}}{\sqrt{n}}\right).$$

Proof of Lemma S3.1: With the definition of $Q_1(\cdot, \cdot)$, for brevity, let $f(\mathcal{X}_i, \delta, \mathcal{D}) = L_h(u_i) - L_h(u_i^*) - \delta L_h'(u_i^*) y_i - L_h'(u_i^*) y_i \langle \mathcal{D}, \mathcal{X}_i \rangle$. We restrict to the event

 $\max_{1 \leq i \leq n} \| \mathcal{X}_i \| \leq b_n := C \sqrt{(I_1 I_2 \vee I_2 I_3) \log n}$. By a direct computation, we have

$$|f(\mathcal{X}_i, \delta, \mathcal{D})| \le C|\delta + \langle \mathcal{D}, \mathcal{X}_i \rangle| \le Cb_n \sqrt{r}t,$$

and

$$f(\boldsymbol{\mathcal{X}}_i, \delta, \boldsymbol{\mathcal{D}}) = L_h''(u_i^*) y_i^2 \delta^2 + 2L_h''(u_i^*) y_i^2 \delta \langle \boldsymbol{\mathcal{D}}, \boldsymbol{\mathcal{X}}_i \rangle + L_h''(u_i^*) y_i^2 \langle \boldsymbol{\mathcal{D}}, \boldsymbol{\mathcal{X}}_i \rangle^2 + o((\delta + \langle \boldsymbol{\mathcal{D}}, \boldsymbol{\mathcal{X}}_i \rangle)^2).$$

Thus,

$$\mathbb{E}[f^2(\boldsymbol{\mathcal{X}}_i, \delta, \boldsymbol{\mathcal{D}})] \leq Cb_n \sqrt{r}t \cdot \{\mathbb{E}[(\delta + \langle \boldsymbol{\mathcal{D}}, \boldsymbol{\mathcal{X}}_i \rangle)^2] + o((\delta + \langle \boldsymbol{\mathcal{D}}, \boldsymbol{\mathcal{X}}_i \rangle)^2)\} \leq C(b_n \sqrt{r}t)^3.$$

For any fixed (δ, \mathbf{D}) , applying Bernstein's inequality yields

$$P(|(P_n - P)f(\boldsymbol{\mathcal{X}}_i, \delta, \boldsymbol{\mathcal{D}})| > u) \le exp\Big\{ -\frac{Cnu^2}{ub_n\sqrt{r}t + (b_n\sqrt{r}t)^3} \Big\}.$$
 (S3.18)

To make the bound uniformly over $\mathscr{S}:=\{(\delta, \mathcal{D}): |\delta|+\|\mathcal{D}\|_* \leq C\sqrt{r}t\}$, we resort to the covering argument. Suppose $\{(\delta_1, \mathcal{D}_1), \cdots, (\delta_M, \mathcal{D}_M)\}$ is ζ -covering of \mathscr{S} with $M \leq (C/\zeta)^{I_1I_2I_3^2}$. It is obvious that $f(\mathcal{X}_i, \delta, \mathcal{D})$ is Lipschitz continuous. Then, by choosing $\zeta \asymp n^{-C}$ and taking $u = C(\frac{b_n\sqrt{r}tI_1I_2I_3^2\log n}{n} + \frac{b_n\sqrt{r}tI_1I_2I_3^2\log n}{n})$

$$\frac{(b_n\sqrt{r}t)^{3/2}\sqrt{I_1I_2I_3^2\mathrm{log}n}}{\sqrt{n}}\big)$$
 in (S3.18), we have

$$\sup_{(\delta, \mathbf{\mathcal{D}}) \in \mathcal{T}} |(P_n - P)f(\mathbf{\mathcal{X}}_i, \delta, \mathbf{\mathcal{D}})| \le C \left(\frac{b_n \sqrt{r} t I_1 I_2 I_3^2 \log n}{n} + \frac{(b_n \sqrt{r} t)^{3/2} \sqrt{I_1 I_2 I_3^2 \log n}}{\sqrt{n}} \right).$$

with probability at least $1 - n^{-C}$.

Lemma S3.2 Suppose Assumption 1-5 hold, with probability $1 - n^{-C}$, we have

$$\begin{split} &|\nabla_{\beta} \breve{Q}(\beta^*, \mathcal{B}^*)| + \|\nabla_{\mathcal{B}} \breve{Q}(\beta^*, \mathcal{B}^*)\| \\ &\leq C \sqrt{\frac{(I_1 I_3 \vee I_2 I_3) \log N}{N}} + C \Big(\frac{r(I_1 I_3 \vee I_2 I_3)}{n} \sqrt{\log n} + \frac{r^2 (I_1 I_3 \vee I_2 I_3)^{5/2} (\log n)^2}{n^{3/2}}\Big). \end{split}$$

Proof of Lemma S3.2: We only consider $\|\nabla_{\mathcal{B}} \breve{Q}(\beta^*, \mathcal{B}^*)\|$ because the bound of $|\nabla_{\beta} \breve{Q}(\beta^*, \mathcal{B}^*)|$ can be derived in the same way. By the definition of $\breve{Q}(\cdot, \cdot)$, it leads to $\nabla_{\mathcal{B}} \breve{Q}(\beta^*, \mathcal{B}^*) = \nabla_{\mathcal{B}} Q_1(\beta^*, \mathcal{B}^*) - \nabla_{\mathcal{B}} Q_1(\widehat{\beta}, \widehat{\mathcal{B}}) + \nabla_{\mathcal{B}} Q(\widehat{\beta}, \widehat{\mathcal{B}})$, and thus

$$\|\nabla_{\mathcal{B}} \breve{Q}(\beta^*, \mathcal{B}^*)\| \leq \|\nabla_{\mathcal{B}} Q_1(\beta^*, \mathcal{B}^*) - \nabla_{\mathcal{B}} Q_1(\widehat{\beta}, \widehat{\mathcal{B}}) - \nabla_{\mathcal{B}} Q(\beta^*, \mathcal{B}^*) + \nabla_{\mathcal{B}} Q(\widehat{\beta}, \widehat{\mathcal{B}})\| + \|\nabla_{\mathcal{B}} Q(\beta^*, \mathcal{B}^*)\|.$$

By a similar argument as Lemma S2.4,

$$\|\nabla_{\mathcal{B}}Q(\beta^*, \mathcal{B}^*)\| \le C\sqrt{rac{(I_1I_3 \lor I_2I_3)\mathrm{log}N}{N}}.$$

Further, applying Lemma S2.5 twice yields

$$\begin{split} &\|\nabla_{\mathcal{B}}Q_{1}(\beta^{*},\boldsymbol{\mathcal{B}}^{*}) - \nabla_{\boldsymbol{\mathcal{B}}}Q_{1}(\widehat{\beta},\widehat{\boldsymbol{\mathcal{B}}}) - \nabla_{\boldsymbol{\mathcal{B}}}Q(\beta^{*},\boldsymbol{\mathcal{B}}^{*}) + \nabla_{\boldsymbol{\mathcal{B}}}Q(\widehat{\beta},\widehat{\boldsymbol{\mathcal{B}}})\| \\ &\leq C\Big(\frac{r(I_{1}I_{3}\vee I_{2}I_{3})}{n}\sqrt{\log n} + \frac{r^{2}(I_{1}I_{3}\vee I_{2}I_{3})^{5/2}(\log n)^{2}}{n^{3/2}}\Big). \end{split}$$

Integrating the above two bounds completes the proof.

Proof of Theorem 4: By direct algebraic computation,

$$\check{Q}(\beta, \mathcal{B}) - \check{Q}(\beta^*, \mathcal{B}^*) \ge \nabla \check{Q}_{\beta}(\beta^*, \mathcal{B}^*)(\beta - \beta^*) + \langle \check{Q}_{\mathcal{B}}(\beta^*, \mathcal{B}^*), \mathcal{B} - \mathcal{B}^* \rangle.$$

The following proof is similar to that in Theorem 1. We have $(\check{\delta}, \check{\mathcal{D}}) = (\check{\beta} - \beta^*, \check{\mathcal{B}} - \mathcal{B}^*) \in \mathcal{H} := \{(\delta, \mathcal{D}) : \|\mathcal{D}^\perp\|_* \leq 3\|\mathcal{D}^0\|_* + |\delta| \} \text{ if } \lambda \geq 2\|\check{Q}_{\mathcal{B}}(\beta^*, \mathcal{B}^*)\|$ and $\lambda \geq 2|\nabla \check{Q}_{\beta}(\beta^*, \mathcal{B}^*)|$, which has been completed in Lemma S3.2 with λ defined in the statement of the theorem. Next, we will connect $\check{Q}(\beta, \mathcal{B}) - \check{Q}(\beta^*, \mathcal{B}^*)$ to $\mathbb{E}[(\beta - \beta^* + \langle \mathcal{X}, \mathcal{B} - \mathcal{B}^* \rangle)^2]$ to bound the difference between empirical risk and population risk corresponding to the central estimator. Assume $\mathbb{E}[(\check{\delta}, \langle \mathcal{X}, \check{\mathcal{D}} \rangle)^2] \geq (La_n)^2$ with a_n in theorem and a sufficiently large L > 0. As in the proof of Theorem 1, this implies

$$\inf_{\substack{(\delta, \mathcal{D}) \in \mathcal{H} \\ \mathbb{E}[(\delta + \langle \mathcal{X}, \mathcal{D} \rangle)^2] \ge (La_n)^2}} \breve{Q}(\beta^* + \delta, \mathcal{B}^* + \mathcal{D}) - \breve{Q}(\beta^*, \mathcal{B}^*) + \lambda \|\mathcal{B}^* + \mathcal{D}\|_* - \lambda \|\mathcal{B}^*\| < 0.$$

Then, for some (δ, \mathcal{D}) that makes the above expression negative, using $|\delta|$ +

 $\|\mathcal{D}\|_* \leq C\sqrt{r}\|\mathcal{D}\|_F + C|\delta| \leq C\sqrt{r}La_n$ and Lemma S3.1 yields

$$C(La_{n})^{2} \leq \mathbb{E}[L_{h}(y_{i}(\beta + \langle \mathcal{B}, \mathcal{X}_{i} \rangle))] - \mathbb{E}[L_{h}(y_{i}(\beta^{*} + \langle \mathcal{B}^{*}, \mathcal{X}_{i} \rangle))]$$

$$\leq \check{Q}(\beta, \mathcal{B}) - \check{Q}(\beta^{*}, \mathcal{B}^{*}) - \nabla \check{Q}_{\beta}(\beta^{*}, \mathcal{B}^{*})(\beta - \beta^{*}) - \langle \check{Q}_{\mathcal{B}}(\beta^{*}, \mathcal{B}^{*}), \mathcal{B} - \mathcal{B}^{*} \rangle$$

$$+ C\left(\frac{b_{n}\sqrt{r}La_{n}I_{1}I_{2}I_{3}^{2}\log n}{n} + \frac{(b_{n}\sqrt{r}La_{n})^{3/2}\sqrt{I_{1}I_{2}I_{3}^{2}\log n}}{\sqrt{n}}\right)$$

$$\leq \check{Q}(\beta^{*} + \delta, \mathcal{B}^{*} + \mathcal{D}) - \check{Q}(\beta^{*}, \mathcal{B}^{*}) + C\lambda(|\delta| + ||\mathcal{D}||_{*})$$

$$+ C\left(\frac{b_{n}\sqrt{r}La_{n}I_{1}I_{2}I_{3}^{2}\log n}{n} + \frac{(b_{n}\sqrt{r}La_{n})^{3/2}\sqrt{I_{1}I_{2}I_{3}^{2}\log n}}{\sqrt{n}}\right)$$

$$\leq \lambda||\mathcal{B}^{*}||_{*} - \lambda||\mathcal{B}^{*} + \mathcal{D}||_{*} + C\lambda(|\delta| + ||\mathcal{D}||_{*})$$

$$+ C\left(\frac{b_{n}\sqrt{r}La_{n}I_{1}I_{2}I_{3}^{2}\log n}{n} + \frac{(b_{n}\sqrt{r}La_{n})^{3/2}\sqrt{I_{1}I_{2}I_{3}^{2}\log n}}{\sqrt{n}}\right)$$

$$\leq \lambda||\mathcal{D}||_{*} + C\lambda(|\delta| + ||\mathcal{D}||_{*})$$

$$+ C\left(\frac{b_{n}\sqrt{r}La_{n}I_{1}I_{2}I_{3}^{2}\log n}{n} + \frac{(b_{n}\sqrt{r}La_{n})^{3/2}\sqrt{I_{1}I_{2}I_{3}^{2}\log n}}{\sqrt{n}}\right)$$

$$\leq C\lambda\sqrt{r}La_{n} + C\left(\frac{b_{n}\sqrt{r}La_{n}I_{1}I_{2}I_{3}^{2}\log n}{n} + \frac{(b_{n}\sqrt{r}La_{n})^{3/2}\sqrt{I_{1}I_{2}I_{3}^{2}\log n}}{\sqrt{n}}\right)$$

$$\leq CLa_{n}^{2} + CL^{3/2}a_{n}^{2},$$

where the third inequality uses Hölder's inequality and Lemma S3.2. The above displayed leads to a contradiction when L is large enough. Thus we have $\mathbb{E}[(\breve{\delta} + \langle \boldsymbol{\mathcal{X}}, \breve{\boldsymbol{\mathcal{D}}} \rangle)^2] \leq (La_n^2)^2$ and the rest of the statements can be obtained easily as for Theorem 1.

Proof of Theorem 5: There exists $\mathbf{Z} \in \partial \| \mathbf{\breve{B}} \|_*$ by the first order optimality con-

dition. Then, we have

$$\frac{1}{n}\sum_{i=1}^{n}L'_{h}(\check{v}_{i})y_{i}\overline{\boldsymbol{X}_{i}}-\frac{1}{n}\sum_{i=1}^{n}L'_{h}(\hat{v}_{i})y_{i}\overline{\boldsymbol{X}_{i}}+\frac{1}{n}\sum_{i=1}^{N}L'_{h}(\hat{v}_{i})y_{i}\overline{\boldsymbol{X}_{i}}+\lambda\overline{\boldsymbol{Z}}=0.$$

Let $\operatorname{rank}_{t}(\breve{\boldsymbol{\mathcal{B}}}) = \breve{r}$, we have $\operatorname{rank}(\overline{\breve{\boldsymbol{\mathcal{B}}}}) = \breve{r}' \leq I_{3}\breve{r}$ and its SVD given by $\overline{\breve{\boldsymbol{\mathcal{B}}}} = \boldsymbol{U}\boldsymbol{S}\boldsymbol{V}^{\mathsf{T}}$, $\boldsymbol{U} \in \mathbb{R}^{I_{1}I_{3} \times \hat{r}'}$, $\boldsymbol{V} \in \mathbb{R}^{I_{2}I_{3} \times \check{r}'}$. Using (2.1) in Koltchinskii et al. (2011), we have

$$\partial \|\overline{\breve{\boldsymbol{B}}}\|_* = \{ \boldsymbol{U}\boldsymbol{V}^\mathsf{T} + \boldsymbol{U}^\perp \boldsymbol{W} (\boldsymbol{V}^\perp)^\mathsf{T} : \|\boldsymbol{W}\|_{op} \le 1 \},$$

where $U^{\perp} \in \mathbb{R}^{I_1 I_3 \times (I_1 I_3 - \hat{r}')}$ and $V^{\perp} \in \mathbb{R}^{I_2 I_3 \times (I_2 I_3 - \check{r}')}$ are orthogonal matrices and their columns are orthogonal to those in U and V, respectively. Let $u_j, v_j, j = 1, \cdots, \hat{r}'$ be the columns of U and V. Pre(post)-multiplying u_j and v_j to (S2.10) yields,

$$-\lambda = \frac{1}{n} \sum_{i=1}^{n} L'_h(\check{v}_i) y_i \boldsymbol{u}_j^{\mathsf{T}} \overline{\boldsymbol{X}}_i \boldsymbol{v}_j - \frac{1}{n} \sum_{i=1}^{n} L'_h(\hat{v}_i) y_i \boldsymbol{u}_j^{\mathsf{T}} \overline{\boldsymbol{X}}_i \boldsymbol{v}_j + \frac{1}{n} \sum_{i=1}^{N} L'_h(\hat{v}_i) y_i \boldsymbol{u}_j^{\mathsf{T}} \overline{\boldsymbol{X}}_i \boldsymbol{v}_j$$
$$:= g_j, \ j = 1, \cdots, \check{r}'.$$

By the similar argument as that in Theorem 2 and Lemma S2.5, we have $\breve{r} \leq Cr$.

Table S4.1: Prediction errors (%) of different estimators with N=10000 and different values of m.

=	d	m	CSL		Ave		Sub	
r			Gaussian	Epanechnikov	Gaussian	Epanechnikov	Gaussian	Epanechnikov
2	0.5	1	14.97 (0.0034)	13.44 (0.0031)	14.97 (0.0034)	13.44 (0.0031)	14.97 (0.0034)	13.44 (0.0031)
		5	14.46 (0.0172)	12.60 (0.0084)	16.35 (0.0094)	15.98 (0.0077)	18.44 (0.0099)	17.64 (0.0084)
		10	14.79 (0.0195)	13.89 (0.0161)	17.96 (0.0218)	17.43 (0.0148)	21.73 (0.0257)	21.06 (0.0193)
		20	15.14 (0.0406)	14.08 (0.0369)	21.61 (0.0194)	20.01 (0.0161)	27.67 (0.0552)	26.68 (0.0579)
	-1	1	7.53 (0.0032)	6.05 (0.0027)	7.53 (0.0032)	6.05 (0.0027)	7.53 (0.0032)	6.05 (0.0027)
		5	5.04 (0.0074)	4.78 (0.0077)	7.72 (0.0071)	6.49 (0.0079)	7.79 (0.0071)	7.78 (0.0082)
		10	5.50 (0.0154)	5.39 (0.0187)	9.04 (0.0175)	8.16 (0.0183)	10.36 (0.0223)	10.14 (0.0253)
		20	6.32 (0.0224)	6.30 (0.0152)	11.02 (0.0173)	10.96 (0.0194)	15.24 (0.0379)	14.84 (0.0330)
5	-0.5	1	16.83 (0.0083)	14.97 (0.0024)	16.83 (0.0083)	14.97 (0.0024)	16.83 (0.0083)	14.97 (0.0024)
		5	15.12 (0.0120)	13.50 (0.0144)	21.98 (0.0184)	19.19 (0.0179)	23.36 (0.0185)	21.89 (0.0100)
		10	16.67 (0.0203)	14.95 (0.0136)	25.57 (0.0307)	24.45 (0.0237)	28.92 (0.0336)	27.80 (0.0227)
		20	17.48 (0.0453)	15.32 (0.0427)	31.31 (0.0394)	29.63 (0.0141)	37.28 (0.0898)	37.82 (0.0797)
	-1	1	6.94 (0.0027)	6.29 (0.0062)	6.94 (0.0027)	6.29 (0.0062)	6.94 (0.0027)	6.29 (0.0062)
		5	5.94 (0.0061)	5.15 (0.0104)	10.31 (0.0104)	10.24 (0.0109)	11.48 (0.0073)	11.75 (0.0073)
		10	6.54 (0.0109)	6.31 (0.0085)	14.55 (0.0099)	13.92 (0.0197)	15.31 (0.0183)	16.16 (0.0180)
		20	7.70 (0.0114)	7.26 (0.0107)	18.65 (0.0296)	18.48 (0.0145)	21.28 (0.0377)	23.12 (0.0583)

Table S4.2: Prediction errors of different estimators with n=1500 and different values of m.

	d	m	CSL		A	Ave		Sub	
,			Gaussian	Epanechnikov	Gaussian	Epanechnikov	Gaussian	Epanechnikov	
2	-0.5	1	18.46 (0.0137)	20.42 (0.0200)	18.46 (0.0137)	20.42 (0.0200)	18.46 (0.0137)	20.42 (0.0200)	
		5	11.20 (0.0131)	12.61 (0.0153)	13.26 (0.0173)	13.37 (0.0153)	18.29 (0.0154)	20.60 (0.0157)	
		10	9.91 (0.0169)	11.69 (0.0250)	12.02 (0.0134)	12.94 (0.0185)	18.55 (0.0101)	20.38 (0.0206)	
		20	10.24 (0.0224)	11.6 (0.0192)	11.44 (0.0198)	11.90 (0.0210)	18.40 (0.0215)	20.11 (0.0241)	
	-1	1	7.64 (0.0047)	9.54 (0.0093)	7.64 (0.0047)	9.54 (0.0093)	7.64 (0.0047)	9.54 (0.0093)	
		5	3.54 (0.0066)	4.27 (0.0071)	6.34 (0.0073)	6.79 (0.0081)	7.92 (0.0053)	9.40 (0.0114)	
		10	3.38 (0.0101)	3.43 (0.0098)	6.49 (0.0109)	7.02 (0.0122)	7.78 (0.0082)	9.72 (0.0129)	
		20	2.70 (0.0057)	3.8 (0.0085)	5.96 (0.0061)	6.03 (0.0088)	8.12 (0.0069)	9.22 (0.0115)	
5	-0.5	1	25.42 (0.0203)	27.28 (0.0200)	25.42 (0.0203)	27.28 (0.0200)	25.42 (0.0203)	27.28 (0.0200)	
		5	11.23 (0.0205)	11.37 (0.0354)	16.63 (0.0211)	17.60 (0.0382)	25.26 (0.0200)	27.52 (0.0392)	
		10	10.44 (0.0365)	10.85 (0.0242)	15.49 (0.0225)	16.61 (0.0265)	25.10 (0.0210)	27.47 (0.0258)	
		20	9.88 (0.0266)	10.21 (0.0227)	14.95 (0.0237)	15.31 (0.0249)	25.16 (0.0251)	27.24 (0.0218)	
	-1	1	12.15 (0.0125)	14.82 (0.0182)	12.15 (0.0125)	14.82 (0.0182)	12.15 (0.0125)	14.82 (0.0182)	
		5	3.95 (0.0121)	3.56 (0.0094)	9.57 (0.0137)	9.92 (0.0105)	11.86 (0.0091)	14.39 (0.0162)	
		10	3.23 (0.0075)	2.69 (0.0075)	9.07 (0.0089)	9.46 (0.0085)	11.89 (0.0067)	14.73 (0.0165)	
		20	2.72 (0.0072)	2.66 (0.0045)	8.57 (0.0078)	8.71 (0.0096)	12.17 (0.0118)	14.55 (0.0149)	

S4 Some results of simulations

S5 Extensions to order-d tensors

As we mentioned in Section 2, our results can be extended to the order-d tensor. In this section, we briefly state the key ideas by introducing the order-d t-product.

For a given high-order tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_d}$, the block circulant operator

Table S4.3: Prediction errors of different estimators with $n=1500,\,m=10$ and different values of I.

r	d	I	CSL		Ave		Sub	
_	<i>u</i>		Gaussian	Epanechnikov	Gaussian	Epanechnikov	Gaussian	Epanechnikov
2	-0.5	20	9.91 (0.0169)	11.69 (0.0250)	17.32 (0.0106)	19.03 (0.0219)	18.55 (0.0101)	20.38 (0.0206)
		40	11.62 (0.0321)	13.84 (0.0543)	19.28 (0.0152)	21.32 (0.0207)	20.018 (0.0139)	22.89 (0.0237)
		60	12.73 (0.0510)	14.23 (0.0502)	21.61 (0.0193)	23.69 (0.0348)	22.56 (0.0243)	24.20 (0.0317)
		100	16.79 (0.0855)	18.24 (0.0764)	25.90 (0.0257)	27.51 (0.0294)	26.43 (0.0281)	28.41 (0.0318)
		200	20.33 (0.0204)	19.52 (0.0114)	32.21 (0.0943)	33.67 (0.0318)	35.45 (0.1253)	34.48 (0.0254)
	1	20	3.38 (0.0101)	3.43 (0.0098)	7.15 (0.0104)	7.73 (0.0146)	7.78 (0.0082)	9.72 (0.0129)
		40	4.12 (0.0167)	4.57 (0.0100)	9.01 (0.0114)	10.02 (0.0135)	9.38 (0.0110)	11.74 (0.0123)
		60	5.23 (0.0105)	6.23 (0.0143)	11.73 (0.0106)	13.90 (0.0185)	11.94 (0.0089)	14.19 (0.018)
		100	8.43 (0.0321)	8.77 (0.0255)	15.84 (0.0384)	18.10 (0.0419)	16.64 (0.0565)	19.84 (0.0448)
		200	11.92 (0.0434)	11.71 (0.0108)	24.69 (0.0481)	26.81 (0.4326)	26.03 (0.0257)	27.57 (0.0304)
5	-0.5	20	10.44 (0.0365)	10.85 (0.0242)	24.80 (0.0143)	24.95 (0.0164)	25.10 (0.0210)	27.47 (0.0258)
		40	11.72 (0.0247)	11.91 (0.0409)	26.91 (0.0344)	26.71 (0.0272)	27.37 (0.0261)	29.71 (0.0312)
		60	13.01 (0.0142)	13.08 (0.0223)	28.51 (0.0667)	28.94 (0.0331)	29.51 (0.0235)	31.91 (0.0428)
		100	16.63 (0.1101)	16.78 (0.1306)	31.76 (0.0285)	33.08 (0.0271)	33.18 (0.0609)	34.22 (0.0639)
		200	19.09 (0.3309)	18.99 (0.4922)	39.81 (0.0374)	40.64 (0.2074)	41.90 (0.1445)	42.86 (0.1918)
	1	20	3.23 (0.0075)	2.69 (0.0075)	11.25 (0.0077)	14.05 (0.0164)	11.89 (0.0067)	14.73 (0.0165)
		40	4.36 (0.0054)	3.62 (0.0049)	13.22 (0.0162)	15.37 (0.0382)	13.87 (0.0208)	16.28 (0.0186)
		60	5.38 (0.0048)	4.41 (0.0083)	14.53 (0.0203)	16.7 (0.0135)	15.67 (0.0162)	17.48 (0.0504)
		100	9.5 (0.018)	8.94 (0.0168)	19.13 (0.0188)	20.54 (0.0275)	20.14 (0.0285)	22.64 (0.0366)
		200	15.70 (0.022)	15.10 (0.0149)	27.65 (0.0702)	29.48 (0.0615)	30.61 (0.0690)	32.86 (0.0637)

Table S4.4: Rank estimation of different estimators with $n=1500,\,m=10$ and different values of L

	d	I	CSL		Ave		Sub	
r		1	Gaussian	Epanechnikov	Gaussian	Epanechnikov	Gaussian	Epanechnikov
2	-0.5	20	2.30 (0.5408)	2.02 (0.6224)	2.40 (0.6135)	2.20 (0.3448)	2.36 (0.4800)	2.10 (0.3265)
		40	1.96 (0.6351)	1.94 (0.7514)	2.22 (0.4463)	1.76 (0.6853)	2.16 (0.4249)	1.88 (0.2894)
		60	2.02 (0.5996)	2 (0.6531)	1.98 (0.4992)	2.12 (0.2971)	1.96 (0.3873)	2.20 (0.4490)
		100	1.94 (0.8739)	1.98 (0.6792)	2.14 (0.3732)	1.88 (0.3367)	1.88 (0.8462)	2.14 (0.2698)
		200	1.54 (0.6176)	1.34 (0.8820)	1 (0)	1 (0)	1 (0)	1 (0)
	1	20	2.02 (0.7347)	2.02 (0.3873)	2.32 (0.4751)	2.08 (0.3184)	2.08 (0.4016)	2.1 (0.2959)
		40	2.08 (0.6955)	2.06 (0.518)	1.92 (0.4637)	2.12 (0.3491)	2.12 (0.4392)	1.94 (0.5065)
		60	1.96 (0.5310)	2.18 (0.5936)	1.88 (0.5934)	2.22 (0.4715)	1.92 (0.4490)	2.34 (0.3139)
		100	2.08 (0.8506)	1.96 (0.6147)	2.08 (0.3349)	2.14 (0.6947)	2.16 (0.6327)	2.16 (0.3412)
		200	1.63 (0.6174)	1.5 (0.6231)	1 (0)	1 (0)	1 (0)	1 (0)
5	-0.5	20	5 (0.6122)	5.04 (0.4065)	4.78 (0.4632)	4.88 (0.4641)	4.56 (0.3331)	5.14 (0.4494)
		40	4.96 (0.4882)	5.04 (0.6045)	5.12 (0.4673)	4.92 (0.7639)	5.54 (0.5392)	4.84 (0.6106)
		60	4.88 (0.4082)	5.08 (0.4914)	4.98 (0.7614)	5.08 (0.6431)	5.18 (0.2731)	4.92 (0.3269)
		100	4.98 (0.7106)	4.96 (1.0600)	5.06 (0.7616)	5.14 (0.5273)	4.94 (0.6149)	4.86 (0.8555)
		200	4.18 (0.9130)	4.74 (0.9420)	1.96 (0.6427)	1.88 (0.8466)	1.94 (0.1567)	1.92 (0.1982)
	1	20	5.18 (0.3955)	5.02 (0.4282)	5.02 (0.4672)	5.08 (0.6416)	4.56 (0.2922)	4.98 (0.4282)
		40	5 (0.4868)	4.98 (0.6224)	4.88 (0.5271)	5.12 (0.9461)	5.26 (0.3596)	5.06 (0.4392)
		60	5.02 (0.6302)	5.02 (0.5506)	4.96 (0.2597)	4.98 (0.4157)	5.32 (0.6302)	5.18 (0.4812)
		100	5.12 (0.5669)	5.34 (0.3922)	5.12 (0.4634)	4.82 (0.4483)	4.88 (0.0200)	4.32 (0.1016)
		200	2.57 (0.2857)	2.74 (0.4412)	1.88 (0.8196)	1.76 (0.6582)	2 (0)	1.68 (0.7310)

 $\operatorname{circ}(\cdot)$ is defined on the *d*-th index of tensor \mathcal{A} as

$$\operatorname{circ}(\mathcal{A}) := egin{bmatrix} \mathcal{A}_1 & \mathcal{A}_d & \cdots & \mathcal{A}_2 \ & \mathcal{A}_2 & \mathcal{A}_1 & \cdots & \mathcal{A}_3 \ & & & & & \ dots & dots & \ddots & dots \ & \mathcal{A}_d & \mathcal{A}_{d-1} & \cdots & \mathcal{A}_1 \end{bmatrix},$$

which is of size $I_1I_d \times I_2I_d \times \cdots \times I_{d-2}I_d \times I_{d-1}$. A block unfolding operator bvec(·), along with its corresponding inverse operator bvfold(·), functions on the d-th index of a high-order tensor, are defined as

$$\operatorname{bvec}(\boldsymbol{\mathcal{A}}) := [\boldsymbol{\mathcal{A}}_1, \cdots, \boldsymbol{\mathcal{A}}_d]^\mathsf{T}, \ \operatorname{bvfold}(\operatorname{bvec}(\boldsymbol{\mathcal{A}})) := \boldsymbol{\mathcal{A}},$$

where $\operatorname{bvec}(\cdot)$ maps \mathcal{A} to an order-(d-1) tensor of size $I_1I_d \times I_2 \times \cdots \times I_{d-1}$. Thus, the order-d t-product (Martin et al., 2013) is a recursive generalization from the order-3 t-product (Kilmer and Martin, 2011).

Definition S5.1 *Martin et al.* (2013)(order-d t-product): The t-product between $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_d}$ and $\mathcal{B} \in \mathbb{R}^{I_2 \times l \times \cdots \times I_d}$ is defined in a recursive pattern,

$$\boldsymbol{\mathcal{A}}\ast\boldsymbol{\mathcal{B}}:=\mathrm{bvfold}(\mathrm{circ}(\boldsymbol{\mathcal{A}})\ast\mathrm{bvec}(\boldsymbol{\mathcal{B}})),$$

which is an $I_1 \times l \times \cdots \times I_d$ real tensor.

Specifically, at the base level of recursion, we replace the operator $\operatorname{cric}(\mathcal{A})$ with $\overline{A} = \operatorname{bcirc}(\mathcal{A})$ which is a block circulant matrix of size $I_1I_3I_4\cdots I_d \times I_2I_3I_4\cdots I_d$. Similarly, instead of $\operatorname{bvec}(\mathcal{B})$, we adopt $\operatorname{bunfold}(\mathcal{B})$, which is a block matrix of $I_2I_3I_4\cdots I_d \times I$ generated in the form of recursion from \mathcal{B} .

Similarly, by DFT, we have

$$(oldsymbol{P}\otimes oldsymbol{I}_{I_1})\cdot \overline{oldsymbol{A}}\cdot (oldsymbol{P}^{\sf H}\otimes oldsymbol{I}_{I_2}) = \widetilde{oldsymbol{A}} = egin{bmatrix} \widetilde{oldsymbol{\mathcal{A}}}^{(1)} & & & & \ & \widetilde{oldsymbol{\mathcal{A}}}^{(2)} & & & \ & & \ddots & & \ & & \widetilde{oldsymbol{\mathcal{A}}}^{(J)} \end{bmatrix},$$

where $\boldsymbol{P} = \boldsymbol{F}_{I_d} \otimes \boldsymbol{F}_{I_{d-1}} \otimes \cdots \boldsymbol{F}_{I_3}$, $\boldsymbol{P}^{\mathsf{H}} = \boldsymbol{F}_{I_d}^{\mathsf{H}} \otimes \boldsymbol{F}_{I_{d-1}}^{\mathsf{H}} \otimes \cdots \boldsymbol{F}_{I_3}^{\mathsf{H}}$, $J = I_3 \cdots I_d$ and $\widetilde{\boldsymbol{\mathcal{A}}} = \mathrm{fft}(\boldsymbol{\mathcal{A}}, [], i)$ for $i = 3, \cdots, d$ (using the MATLAB command for ease of exposition). By the unitary invariance, it also yields that,

$$\|\mathcal{A}\|_F = \frac{1}{\sqrt{I_3}} \|\widetilde{A}\|_F, \ \langle \mathcal{A}, \mathcal{B} \rangle = \frac{1}{I_3} \langle \widetilde{A}, \widetilde{A} \rangle,$$
 $\mathcal{C} = \mathcal{A} * \mathcal{B} \Leftrightarrow \widetilde{\mathcal{C}} = \widetilde{\mathcal{A}} \cdot \widetilde{\mathcal{B}}.$

The order-d conjugate transpose is relatively complex. Since the denifinition

of the orthogonal tensor, the f-diagonal tensor and so on can be easily generalized from Appendix S1.

Definition S5.2 Qin et al. (2022)(order-d Conjugate transpose): We employ $\mathcal{A}^{\mathsf{H}} \in \mathbb{C}^{I_2 \times I_1 \times \cdots \times I_d}$ to represent the conjugate transpose of $\mathcal{A} \in \mathbb{C}^{I_1 \times I_2 \times \cdots \times I_d}$ where $\widetilde{\mathcal{A}}^{\mathsf{H}}(:,:,i_3,\cdots,i_d) = (\widetilde{\mathcal{A}}(:,:,i_3,\cdots,i_d))^{\mathsf{H}}$ for $i_p \in \{1,\cdots,d_p\}$ with $p \in \{3,\cdots,d\}$.

With the above arguments of order-d t-product, the t-SVD, tubal rank and t-TNN can also be extended to order-d tensors. We refer the readers to Martin et al. (2013); Qin et al. (2022); Song et al. (2025) for more discussions. Our analysis mainly rely on the operator norm and the decomposability of the t-TNN for which we list the key results which can be proved similarly as Appendix S1.

Definition S5.3 Define $\|A\|:=\|\overline{A}\|_{op}$ as the tensor spectral norm for $A \in \mathbb{R}^{I_1 \times I_2 \times \cdots I_d}$.

Theorem S5.1 For any order-d tensor $\mathcal{A} \in \mathcal{U}^{\perp}$ and any $\mathcal{B} \in \mathcal{B}$, of the same dimension, with tubal rank r, we have $\|\mathcal{A} + \mathcal{B}\|_* = \|\mathcal{A}\|_* + \|\mathcal{B}\|_*$, where

$$egin{aligned} \mathscr{U} &= \Big\{ \mathcal{U} * \mathcal{M} + \mathcal{N} * \mathcal{V}^\mathsf{H} : \mathcal{M} \in \mathbb{R}^{r imes I_2 imes \cdots imes I_d}, \mathcal{N} \in \mathbb{R}^{I_1 imes r imes \cdots imes I_d} \Big\}, \ \mathscr{B} &= \Big\{ \mathcal{U} * \mathcal{U}^\mathsf{H} * \mathcal{Z} + \mathcal{Z} * \mathcal{V} * \mathcal{V}^\mathsf{H} - \mathcal{U} * \mathcal{U}^\mathsf{H} * \mathcal{Z} * \mathcal{V} * \mathcal{V}^\mathsf{H} : \mathcal{Z} \in \mathbb{R}^{I_1 imes I_2 imes \cdots imes I_d} \Big\}. \end{aligned}$$

Based on Definition S5.3 and Theorem S5.1, by the same argument as Sections

2 and 3, we can extend the derived convergence rate to the order-d tensor. For example, with the same condition as Corollary 1, we have

$$(\widehat{\beta} - \beta^*)^2 + \|\widehat{\boldsymbol{\mathcal{B}}} - \boldsymbol{\mathcal{B}}^*\|_F^2 \le O_p\left(\frac{r(I_1 I_3 I_4 \cdots I_d \vee I_2 I_3 I_4 \cdots I_d)}{n}\right),$$

for the estimated order-d tensor parameter, where $\widehat{\beta}$ and $\widehat{\mathcal{B}}$ are obtained by the same way as in Section 2. The distributed estimator in Section 3 can be generalized in the same way.

Bibliography

Bartlett, P. L., O. Bousquet, and S. Mendelson (2005). Local rademacher complexities. *Ann. Statist.* 33(4), 1497–1537.

Beck, A. and M. Teboulle (2009, 01). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sciences* 2, 183–202.

Geer, S. v. d. (2000). *Empirical Processes in M-Estimation*. Cambridge: Cambridge University Press.

Kilmer, M. and C. Martin (2011, 08). Factorization strategies for third-order tensors. *Linear Algebra Appl.* 435, 641–658.

Kilmer, M. E., C. D. Martin, and L. Perrone (2008). A third-order generaliza-

- tion of the matrix svd as a product of third-order tensors. *Tufts University*, *Department of Computer Science*, *Tech. Rep. TR-2008-4*.
- Koltchinskii, V. (2011). Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems. New York, NY, USA: Springer.
- Koltchinskii, V., K. Lounici, and A. B. Tsybakov (2011). Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *Ann. Stat.* 39(5), 2302 2329.
- Lu, C., J. Feng, Y. Chen, W. Liu, Z. Lin, and S. Yan (2018, 04). Tensor robust principal component analysis with a new tensor nuclear norm. *IEEE Trans. Pattern Anal. Mach. Intell.* 42(4), 925–938.
- Martin, C. D., R. Shafer, and B. LaRue (2013). An order-p tensor factorization with applications in imaging. *SIAM J. Sci. Comput.* 35(1), A474–A490.
- Pollard, D. (1984). *Convergence of Stochastic Processes*. New York, NY, USA: Springer-Verlag.
- Qin, W., H. Wang, F. Zhang, J. Wang, X. Luo, and T. Huang (2022). Low-rank high-order tensor completion with applications in visual data. *IEEE Trans. Image Process.* 31, 2433–2448.
- Recht, B., M. Fazel, and P. A. Parrilo (2010). Guaranteed minimum-rank so-

lutions of linear matrix equations via nuclear norm minimization. *SIAM REV.* 52(3), 471–501.

Song, Z., Y. Chen, and Z. Weihua (2025). Adaptively robust high-order tensor factorization for low-rank tensor reconstruction. *Pattern Recognit.* 165, 111600.

Vershynin, R. (2011). Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*.