

# Sparse and Low-Rank Tensor Recovery via Cubic-Sketching

Guang Cheng  
Department of Statistics  
Purdue University  
[www.science.purdue.edu/bigdata](http://www.science.purdue.edu/bigdata)

CCAM@Purdue Math  
Oct. 27, 2017

Joint work with Botao Hao and Anru Zhang

# Tensor: Multi-dimensional Array



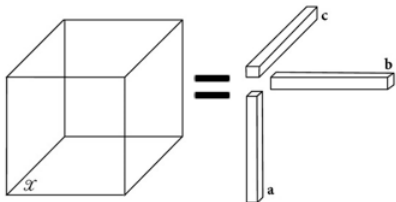
Vector: order-1 tensor

Matrix: order-2 tensor



Order-3 tensor

Rank-one Tensor

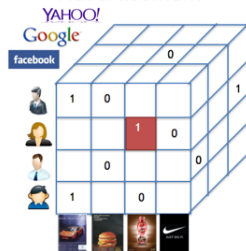


# Tensor Data Example

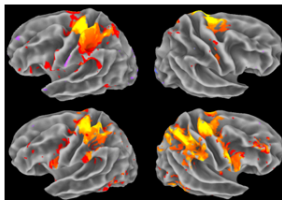
Color image



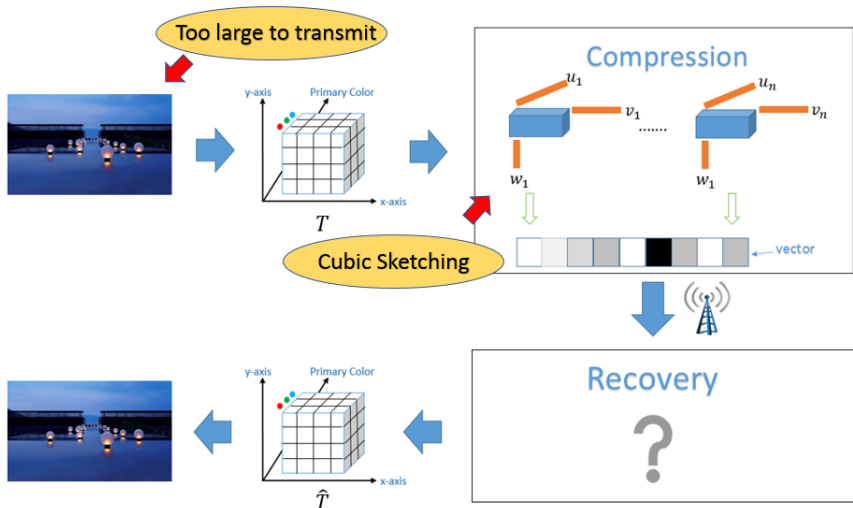
Advertisement



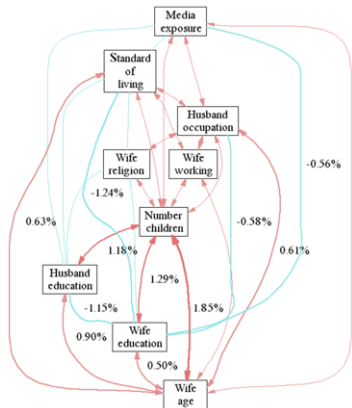
fMRI



# Motivation: Compressed Image Transmission



# Motivation: Interaction Effect Model



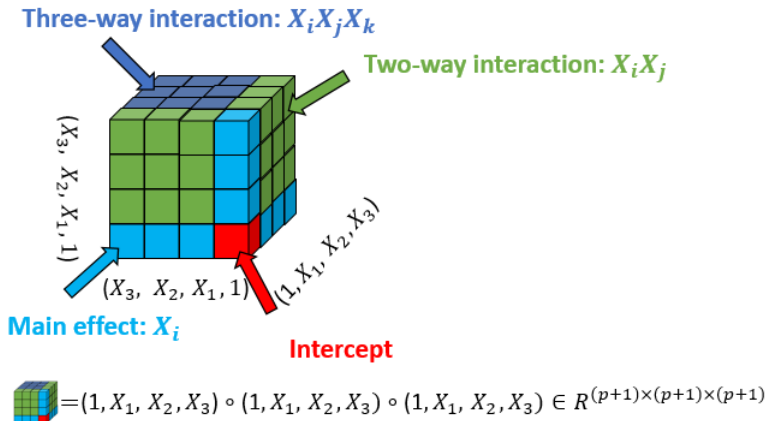
covariate

$$\mathbb{E}(y|\mathbf{X}) = \beta_0 + \sum_{i=1}^p X_i \beta_i + \sum_{i,j=1}^p \gamma_{ij} X_i X_j + \sum_{i,j,k=1}^p \eta_{ijk} X_i X_j X_k$$

Intercept Main effect Pairwise interaction Triple-wise interaction

source: Contraceptive Method Choice dataset from UCI

# Motivation: Interaction Effect Model



# Sparse and Low-Rank Tensor Recovery

# Noisy Cubic Sketching Model

- Observe  $\{y_i, \mathcal{X}_i\}$  from noisy cubic sketching model,

$$\underbrace{y_i}_{\text{scalar}} = \underbrace{\langle \mathcal{T}^*, \mathcal{X}_i \rangle}_{\text{tensor inner product}} + \underbrace{\epsilon_i}_{\text{noise}}, \quad i = 1, \dots, n.$$

- For two tensors  $\mathcal{A} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$  and  $\mathcal{B} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$ , the tensor inner product is defined as

$$\langle \mathcal{A}, \mathcal{B} \rangle = \sum_{ijk} \mathcal{A}_{ijk} \mathcal{B}_{ijk}.$$



# Noisy Cubic Sketching Model

- General model including: tensor regression (Zhou, Li, Zhu (2013)), tensor completion (Yuan and Zhang (2014)).

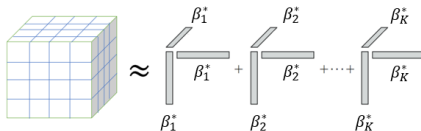


- Goal: Recover unknown third-order tensor parameter  $\mathcal{T}^*$ .
- High-dimensional problem:  $n \ll \dim(\mathcal{T}^*) \approx p^3$ .

# Key Assumptions on Tensor Parameter

- When  $\mathcal{T}^* \in \mathbb{R}^{p \times p \times p}$  is a symmetric tensor...

① CANDECOMP/PARAFAC(CP) low-rank:



$$\mathcal{T}^* = \sum_{k=1}^K \eta_k^* \beta_k^* \circ \beta_k^* \circ \beta_k^*, \text{ with } \|\beta_k^*\|_2 = 1$$

Represented as sum of rank-one tensor, where  $k \ll p$ .

- ② Sparse components:  $\|\beta_k^*\|_0 \leq s$  for  $k \in [K]$ .
- The cubic sketching tensor  $\mathcal{X}_i$  for symmetric case is  $\mathcal{X}_i = \mathbf{x}_i \circ \mathbf{x}_i \circ \mathbf{x}_i$ , where  $\{\mathbf{x}_i\}_{i=1}^n$  are Gaussian random vectors.
- $\beta_k^*$  and  $\beta_{k'}^*$  are **not orthogonal**. Different from singular-value decomposition in **matrix** case.

# Key Assumptions on Tensor Parameter

- When  $\mathcal{T}^* \in \mathbb{R}^{p_1 \times p_2 \times p_3}$  is a non-symmetric tensor...

① CANDECOMP/PARAFAC (CP) low-rank:

$$\mathcal{T}^* = \sum_{k=1}^K \eta_k^* \beta_{1k}^* \circ \beta_{2k}^* \circ \beta_{3k}^*, \text{ with } \|\beta_{1k}^*\|_2 = \|\beta_{2k}^*\|_2 = \|\beta_{3k}^*\|_2 = 1$$

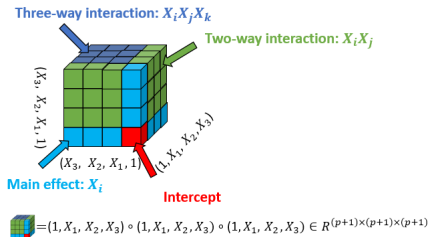
- ② Sparse components:  $\|\beta_{1k}^*\|_0 \leq s_1$ ,  $\|\beta_{2k}^*\|_0 \leq s_2$ ,  $\|\beta_{3k}^*\|_0 \leq s_3$  for  $k \in [K]$ .
- The cubic sketching tensor  $\mathcal{X}_i$  for non-symmetric case is  $\mathcal{X}_i = \mathbf{u}_i \circ \mathbf{v}_i \circ \mathbf{w}_i$ , where  $\{\mathbf{u}_i, \mathbf{v}_i, \mathbf{w}_i\}_{i=1}^n$  are Gaussian random vectors.

# Reduced Symmetric Tensor Recovery Model

- For symmetric tensor recovery model

$$y_i = \left\langle \sum_{k=1}^K \eta_k^* \beta_k^* \circ \beta_k^* \circ \beta_k^*, \mathbf{x}_i \circ \mathbf{x}_i \circ \mathbf{x}_i \right\rangle + \epsilon_i = \sum_{k=1}^K \eta_k^* \underbrace{(\mathbf{x}_i^\top \beta_k^*)^3}_{\text{non-linear}} + \epsilon_i$$

- Connect with *interaction effect model*.



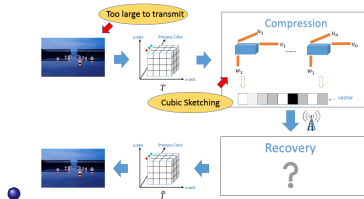
- New Goal: Recover  $\{\eta_k^*, \beta_k^*\}_{k=1}^K$

# Reduced Non-symmetric Tensor Recovery Model

- For non-symmetric tensor recovery model

$$\begin{aligned}
 y_i &= \left\langle \sum_{k=1}^K \eta_k^* \beta_{1k}^* \circ \beta_{2k}^* \circ \beta_{3k}^*, \mathbf{u}_i \circ \mathbf{v}_i \circ \mathbf{w}_i \right\rangle + \epsilon_i \\
 &= \sum_{k=1}^K \eta_k^* \underbrace{(\mathbf{u}_i^\top \beta_{1k}^*)(\mathbf{v}_i^\top \beta_{2k}^*)(\mathbf{w}_i^\top \beta_{3k}^*)}_{\text{non-linear}} + \epsilon_i
 \end{aligned}$$

- Connect with *compressed image transmission model*.



- New Goal: Recover  $\{\eta_k^*, \beta_{1k}^*, \beta_{2k}^*, \beta_{3k}^*\}_{k=1}^K$ .

- Consider Empirical Risk Minimization

$$\widehat{\mathcal{J}} = \operatorname{argmin}_{\{\eta_k, \beta_k\}} \underbrace{\sum_{i=1}^n (y_i - \sum_{k=1}^K \eta_k (\mathbf{x}_i^\top \beta_k)^3)^2}_{\mathcal{L}_1(\eta_k, \beta_k)}$$
$$\widehat{\mathcal{J}} = \operatorname{argmin}_{\{\eta_k, \beta_{ik}\}} \underbrace{\sum_{i=1}^n (y_i - \sum_{k=1}^K \eta_k (\mathbf{u}_i^\top \beta_{1k})(\mathbf{v}_i^\top \beta_{2k})(\mathbf{w}_i^\top \beta_{3k}))^2}_{\mathcal{L}_2(\eta_k, \beta_{ik})}$$

- Difficulties: *Non-convex optimization!* Non-convexity from **cube structure** or **tri-convexity**.

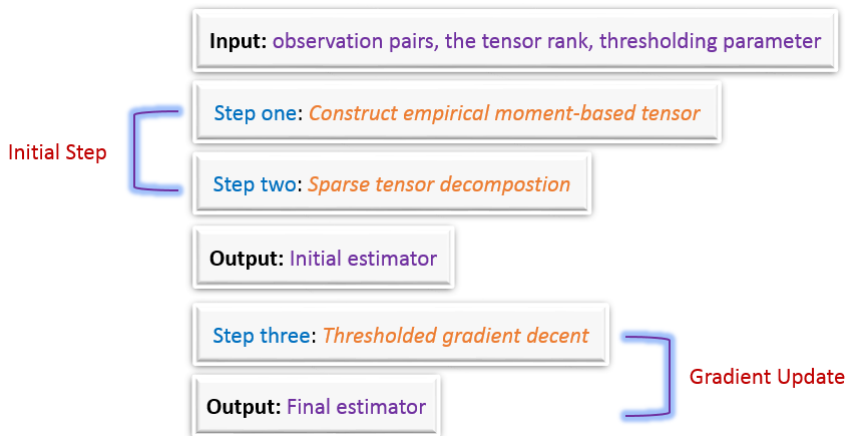
# Our Contributions

- ① Efficient two-stage implementation to non-convex optimization problem.
- ② Non-asymptotic analysis. Provide optimal estimation rate.

## Two-stage Implementation



# Main Algorithm



# Initial Step: Non-symmetric unbiased estimator

- Construct an unbiased empirical moment based tensor  $\mathcal{T}(y_i, \mathcal{X}_i) \in \mathbb{R}^{p_1 \times p_2 \times p_3}$  as following

$$\mathcal{T} := \underbrace{\frac{1}{n} \sum_{i=1}^n y_i \mathbf{u}_i \circ \mathbf{v}_i \circ \mathbf{w}_i}_{\text{only depends on observations.}}$$

- This is used for non-symmetric case and  $\mathbf{u}_i, \mathbf{v}_i, \mathbf{w}_i$  are independent standard Gaussian random vectors.

# Initial Step: Symmetric unbiased estimator

- Construct an unbiased empirical moment based tensor  $\mathcal{T}_s(y_i, \mathcal{X}_i) \in \mathbb{R}^{p \times p \times p}$  as following

$$\mathcal{T}_s := \underbrace{\frac{1}{6} \left[ \frac{1}{n} \sum_{i=1}^n y_i \mathbf{x}_i \circ \mathbf{x}_i \circ \mathbf{x}_i - \mathcal{U} \right]}_{\text{only depends on observations.}}$$

where the bias term

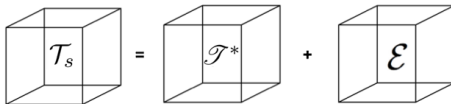
$\mathcal{U} = \sum_{j=1}^p (\mathbf{m}_1 \circ \mathbf{e}_j \circ \mathbf{e}_j + \mathbf{e}_j \circ \mathbf{m}_1 \circ \mathbf{e}_j + \mathbf{e}_j \circ \mathbf{e}_j \circ \mathbf{m}_1)$ , and  $\mathbf{m}_1 = \frac{1}{n} \sum_{i=1}^n y_i \mathbf{x}_i$ . Here  $\{\mathbf{e}_j\}_{j=1}^p$  are the canonical vectors in  $\mathbb{R}^p$ .

- Bias term  $\mathcal{U}$  is due to the correlation among three “identical” Gaussian random vectors.

# Initial Step: Decompose unbiased estimator

- Intuition:  $\mathbb{E}[\mathcal{T}_s] = \mathcal{T}^*$ .

Tensor Denoising Model:  $\mathcal{T}_s = \mathcal{T}^* + \mathcal{E}$

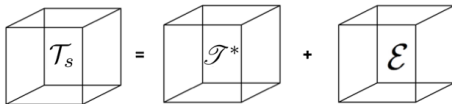


- Observation  $\mathcal{T}_s$ .
- Noise  $\mathcal{E} = \mathcal{T}_s - \mathbb{E}(\mathcal{T}_s)$ : approximation error.
- Decompose  $\mathcal{T}_s$  to obtain  $\{\eta_k^{(0)}, \beta_k^{(0)}\}$  or Decompose  $\mathcal{T}$  to obtain  $\{\eta_k^{(0)}, \beta_{1k}^{(0)}, \beta_{2k}^{(0)}, \beta_{3k}^{(0)}\}$  through **sparse tensor decomposition**. See next slide for details.
- *Far from the optimal estimation, but good enough as a warm start.*

# Initial Step: Decompose unbiased estimator

- Intuition:  $\mathbb{E}[\mathcal{T}_s] = \mathcal{T}^*$ .

Tensor Denoising Model:  $\mathcal{T}_s = \mathcal{T}^* + \mathcal{E}$

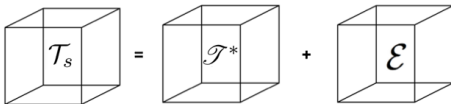


- Observation  $\mathcal{T}_s$ .
- Noise  $\mathcal{E} = \mathcal{T}_s - \mathbb{E}(\mathcal{T}_s)$ : approximation error.
- Decompose  $\mathcal{T}_s$  to obtain  $\{\eta_k^{(0)}, \beta_k^{(0)}\}$  or Decompose  $\mathcal{T}$  to obtain  $\{\eta_k^{(0)}, \beta_{1k}^{(0)}, \beta_{2k}^{(0)}, \beta_{3k}^{(0)}\}$  through **sparse tensor decomposition**. See next slide for details.
- *Far from the optimal estimation, but good enough as a warm start.*

# Initial Step: Decompose unbiased estimator

- Intuition:  $\mathbb{E}[\mathcal{T}_s] = \mathcal{T}^*$ .

Tensor Denoising Model:  $\mathcal{T}_s = \mathcal{T}^* + \mathcal{E}$



- Observation  $\mathcal{T}_s$ .
- Noise  $\mathcal{E} = \mathcal{T}_s - \mathbb{E}(\mathcal{T}_s)$ : approximation error.
- Decompose  $\mathcal{T}_s$  to obtain  $\{\eta_k^{(0)}, \beta_k^{(0)}\}$  or Decompose  $\mathcal{T}$  to obtain  $\{\eta_k^{(0)}, \beta_{1k}^{(0)}, \beta_{2k}^{(0)}, \beta_{3k}^{(0)}\}$  through **sparse tensor decomposition**. See next slide for details.
- *Far from the optimal estimation, but good enough as a warm start.*

# Sparse Tensor Decomposition for Symmetric Tensor

- ⇒ Generate  $L$  starting points  $\{\beta_l^{\text{start}}\}_{l=1}^L$ .
- ⇒ For each starting point, compute a *non-sparse* factor of moment-based  $\mathcal{T}_s$  via symmetric tensor power update:

$$\tilde{\beta}_l^{(t+1)} = \frac{\mathcal{T}_s \times_2 \beta_l^{(t)} \times_3 \beta_l^{(t)}}{\|\mathcal{T}_s \times_2 \beta_l^{(t)} \times_3 \beta_l^{(t)}\|_2},$$

where for  $\mathcal{T}_s \in \mathbb{R}^{p \times p \times p}$  and  $\mathbf{x} \in \mathbb{R}^p$ , define

$$\mathcal{T}_s \times_2 \mathbf{x} \times_3 \mathbf{x} := \sum_{j,l} \mathbf{x}_j \mathbf{x}_l [\mathcal{T}]_{:,j,l}.$$

- ⇒ Get a *sparse solution*  $\beta_l^{(t+1)}$  via thresholding or truncation.
- ⇒ Cluster  $L$  sets of single component  $\{\beta_l^{(T)}, \beta_l^{(T)}, \beta_l^{(T)}\}_{l=1}^L$  into  $K$  clusters to obtain a rank- $K$  decomposition  $\{\beta_k^{(0)}, \beta_k^{(0)}, \beta_k^{(0)}\}_{k=1}^K$ .

*Different from matrix SVD due to non-orthogonality.*

# Sparse Tensor Decomposition for Non-symmetric Tensor

- ⇒ Generate  $L$  starting points  $\{\beta_{1l}^{\text{start}}, \beta_{2l}^{\text{start}}, \beta_{3l}^{\text{start}}\}_{l=1}^L$ .
- ⇒ For each starting point, compute a **non-sparse** factor of moment-based  $\mathcal{T}$  via alternating tensor power update:

$$\tilde{\beta}_{1l}^{(t+1)} = \frac{\mathcal{T}_s \times_2 \beta_{2l}^{(t)} \times_3 \beta_{3l}^{(t)}}{\|\mathcal{T}_s \times_2 \beta_{2l}^{(t)} \times_3 \beta_{3l}^{(t)}\|_2},$$

The updates for  $\tilde{\beta}_{2l}^{(t+1)}$  and  $\tilde{\beta}_{3l}^{(t+1)}$  are similar.

- ⇒ Get a **sparse solution**  $\{\beta_{1l}^{(t+1)}, \beta_{2l}^{(t+1)}, \beta_{3l}^{(t+1)}\}$  via thresholding or truncation.
- ⇒ Cluster  $L$  sets of single component  $\{\beta_{1l}^{(T)}, \beta_{2l}^{(T)}, \beta_{3l}^{(T)}\}_{l=1}^L$  into  $K$  clusters to obtain a rank- $K$  decomposition  $\{\beta_{1k}^{(0)}, \beta_{2k}^{(0)}, \beta_{3k}^{(0)}\}_{k=1}^K$ .



# Gradient Update: Thresholded Gradient Decent

- ⇒ Input initial estimator  $\{\eta_k^{(0)}, \beta_k^{(0)}\}_{k=1}^K$ .  
⇒ In each iteration step, update  $\{\beta_k\}_{k=1}^K$  as

$$\tilde{\beta}_k^{(t+1)} = \beta_k^{(t)} - \frac{\mu_t}{\phi} \nabla_{\beta_k} \mathcal{L}_1(\eta_k^{(0)}, \beta_k^{(t)})$$

where  $\phi = \frac{1}{n} \sum_{i=1}^n y_i^2$ ,  $\mu_t$  is the step size.

- ⇒ Sparsify current update by thresholding  $\beta_k^{(t+1)} = \varphi_\rho(\tilde{\beta}_k^{(t+1)})$ .  
⇒ Normalize final update  $\beta_k^{(T)} = \frac{\beta_k^{(T)}}{\|\beta_k^{(T)}\|_2}$  and update the weight  
 $\hat{\eta}_k = \eta_k^{(0)} \times \|\beta_k^{(T)}\|_2^3$ .

---

<sup>1</sup>Alternating update for non-symmetric tensor recovery.

# Gradient Update: Thresholded Gradient Decent

⇒ Input initial estimator  $\{\eta_k^{(0)}, \beta_{1k}^{(0)}, \beta_{2k}^{(0)}, \beta_{3k}^{(0)}\}_{k=1}^K$ .

⇒ In each iteration step, alternatively update  $\{\beta_{1k}, \beta_{2k}, \beta_{3k}\}_{k=1}^K$  as

$$\tilde{\beta}_{1k}^{(t+1)} = \beta_{1k}^{(t)} - \frac{\mu_t}{\phi} \nabla_{\beta_{1k}} \mathcal{L}_2(\eta_k^{(0)}, \beta_{1k}^{(t)}, \beta_{2k}^{(t)}, \beta_{3k}^{(t)})$$

where  $\phi = \frac{1}{n} \sum_{i=1}^n y_i^2$ ,  $\mu_t$  is the step size. The update for  $\tilde{\beta}_{2k}^{(t+1)}$  and  $\tilde{\beta}_{3k}^{(t+1)}$  is similar.

⇒ Sparsify current update by thresholding  $\beta_{jk}^{(t+1)} = \varphi_\rho(\tilde{\beta}_{jk}^{(t+1)})$  for  $j = 1, 2, 3$ .

⇒ Normalize final update  $\beta_{jk}^{(T)} = \frac{\beta_{jk}^{(T)}}{\|\beta_{jk}^{(T)}\|_2}$  and update the weight

$$\hat{\eta}_k = \eta_k^{(0)} \times \|\beta_{1k}^{(T)}\|_2 \|\beta_{2k}^{(T)}\|_2 \|\beta_{3k}^{(T)}\|_2.$$

---

<sup>1</sup>Alternating update for non-symmetric tensor recovery.

# Non-asymptotic Analysis

# Non-asymptotic Upper Bound

## 定理

Suppose some regularity conditions for the true tensor parameter hold. Assume  $n \geq C_0 s^2 \log p$  for some large constant  $C_0$ . Denote  $Z_k^{(t)} = \sum_{k=1}^K \|\sqrt[3]{\eta_k} \beta_k^{(t)} - \sqrt[3]{\eta_k^*} \beta_k^*\|_2^2$  For **any**  $t = 0, 1, 2, \dots$ , the factor-wise estimator satisfies

$$Z_k^{(t+1)} \leq \underbrace{\kappa^t Z_k^{(t)}}_{\text{computational error}} + \underbrace{\frac{C_1 \eta_{\min}^{*- \frac{4}{3}}}{16} \frac{\sigma^2 s \log p}{n}}_{\text{statistical error}},$$

with high probability, where  $\kappa$  is the contraction parameter between 0 and 1,  $\eta_{\min}^* = \min_k \{\eta_k^*\}$ ,  $\sigma$  is the noise level and  $C_0, C_1$  are some absolute constants.

- Interesting characterization for computational error and statistical error;
- Geometric convergence rate to the truth in the noiseless case and minimax optimal statistical rate shown later;
- The error bound is dominated by **computation error** in the first several iterations and then is dominated by **statistical error**.  
Useful guideline for choosing stopping rule.

- When  $t \geq T$  for some enough  $T$ , the final estimator is bounded by

$$\left\| \mathcal{J}^{(T)} - \mathcal{J}^* \right\|_F^2 \leq \frac{C\sigma^2 K s \log p}{n},$$

with high probability.

- *Minimax optimal rate!*

# Class of Sparse and Low-rank tensor

- Sparse CP decomposition

$$\mathcal{T} = \sum_{k=1}^K \beta_k \circ \beta_k \circ \beta_k, \|\beta_k\|_0 \leq s \text{ for } k \in [K]$$

- Incoherence condition(nearly orthogonal): The true tensor components are incoherent such that

$$\max_{k_i \neq k_j \in [K]} |\langle \beta_{k_i}^*, \beta_{k_j}^* \rangle| \leq \frac{C}{\sqrt{s}}.$$

## 定理

*Consider the class of tensor satisfy sparse CP-decomposition and incoherence condition. Suppose we sample via cubic measurements with i.i.d. standard normal sketches with i.i.d.  $N(0, \sigma^2)$  noise, then we have the following lower bound result for recovery loss for this class of low-rank tensors,*

$$\inf_{\widehat{\mathcal{T}}} \sup_{\mathcal{T} \in \mathcal{F}} \mathbb{E} \left\| \widehat{\mathcal{T}} - \mathcal{T} \right\|_F^2 \geq c \sigma^2 \frac{K s \log(ep/s)}{n}.$$



## 定理

Consider the class of tensor  $\mathcal{F}_{p,K,s}$  satisfy sparse CP-decomposition and incoherence condition. Suppose we observe  $n$  samples  $\{y_i, \mathcal{X}_i\}_{i=1}^n$  from symmetric tensor cubic sketching model, where  $n \geq Cs^2 \log p$  for some large constant  $C$ . Then the estimator  $\widehat{\mathcal{T}}$  achieves

$$\inf_{\widehat{\mathcal{T}}} \sup_{\mathcal{T} \in \mathcal{F}_{p,K,s}} \mathbb{E} \left\| \widehat{\mathcal{T}} - \mathcal{T} \right\|_F^2 \asymp \underbrace{\sigma^2 \frac{Ks \log(p/s)}{n}}_{R^*},$$

when  $\log p \asymp \log p/s$ . Here  $\sigma$  is the noise level.

- Our analysis is **non-asymptotic** and our estimator is **rate-optimal**.
- In general, we have a trade-off  $\rightarrow R^*$  is the outcome of **statistical error** and **optimization error** trade-off.
- Similar argument holds for non-symmetric case. *Different technical tools are used.*
- To overcome the obstacle from high-order Gaussian random variable, we develop novel high-order concentration inequality by using *truncation argument* and  $\psi_\alpha$ -norm.

# Application to Interaction Effect Model

- Given the response  $y \in \mathbb{R}^n$  and covariates  $\mathbf{X} \in \mathbb{R}^{n \times p}$ , the regression model with three-way interactions can be formulated as

$$\begin{aligned} y_l &= \beta_0 + \sum_{i=1}^p X_{li} \beta_i + \sum_{i,j=1}^p \gamma_{ij} X_{li} X_{lj} + \sum_{i,j,k=1}^p \eta_{ijk} X_{li} X_{lj} X_{lk} + \epsilon_l \\ &= \langle \mathcal{B}, \mathbf{X}_l \circ \mathbf{X}_l \circ \mathbf{X}_l \rangle + \epsilon_l \end{aligned}$$

where  $\mathbf{X}_l = (\mathbf{1}, \mathbf{X}_l^\top)^\top \in \mathbb{R}^{p+1}$ .

# Application to Interaction Effect Model

- It is reasonable to assume that  $\mathcal{B}$  possess low-rank and/or sparsity structures in some biometrics studies (Hung, et. 2016).

$$y_l = \left\langle \sum_{k=1}^K \eta_k \beta_k \circ \beta_k \circ \beta_k, \mathbf{X}_l \circ \mathbf{X}_l \circ \mathbf{X}_l \right\rangle + \varepsilon_l$$

- The symmetric tensor recovery model can be treated as a high-order interaction effect model.

# Some Changes for Algorithm

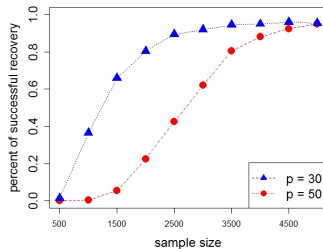
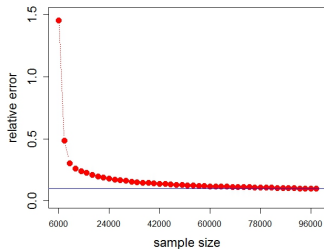
- The unbiased empirical moment based tensor  $\mathcal{A} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$  for interaction effect model is constructed as following:
  - Define three quantities  $\mathbf{a} = \frac{1}{n} \sum_{l=1}^n y_l X_l$ ,  
 $\tilde{\mathcal{A}} = \frac{1}{n} \sum_{l=1}^n y_l X_l \circ X_l \circ X_l$ ,  
 $\bar{\mathcal{A}} = \frac{1}{6}(\tilde{\mathcal{A}} - \sum_{j=1}^p (\mathbf{a} \circ \mathbf{e}_j \circ \mathbf{e}_j + \mathbf{e}_j \circ \mathbf{a} \circ \mathbf{e}_j + \mathbf{e}_j \circ \mathbf{e}_j \circ \mathbf{a}))$ .
  - For  $i, j, k \neq 0$ ,  $\mathcal{A}_{ijk} = \bar{\mathcal{A}}_{ijk}$ .
  - For  $i \neq 0$ ,  $\mathcal{A}_{0,0,i} = \frac{1}{3}\tilde{\mathcal{A}}_{0,0,i} - \frac{1}{6}(\sum_{k=1}^p \tilde{\mathcal{A}}_{k,k,i} - (p+2)a_i)$ . And  
 $\mathcal{A}_{0,0,0} = \frac{1}{2p-2}(\sum_{k=1}^p \tilde{\mathcal{A}}_{0,k,k} - (p+2)\tilde{\mathcal{A}}_{0,0,0})$ .
- The additional intercept term changes the model structure dramatically.

- Recover low-rank and sparse symmetric tensor  $\mathcal{T}^* \in \mathbb{R}^{p \times p \times p}$  from

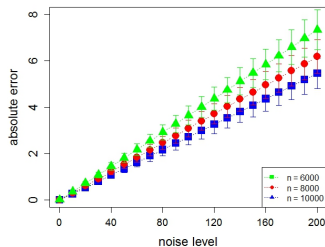
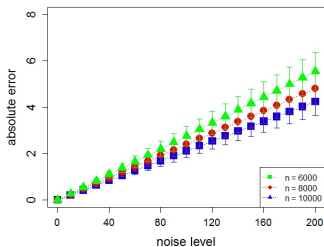
$$y_i = \langle \mathcal{T}^*, \mathcal{X}_i \rangle + \epsilon_i, \quad i = 1, \dots, n.$$

- Proportion of non-zero elements for each factor  $s = 0.3$ .  
Tensor CP-rank  $K = 3$ . Replication = 200.
  - Stopping rule for initialization:  $\|\beta_m^{(l+1)} - \beta_m^{(l)}\|_2 \leq 10^{-6}$ .
  - Stopping rule for gradient update:  $\|\mathbf{B}^{(T+1)} - \mathbf{B}^{(T)}\|_F \leq 10^{-6}$ .
- The dimension, sample size and noise level vary in different scenarios.

- Left panel: relative error for initialization. Right panel: percent of successful recovery with varying sample size. Both are noiseless case with  $p = 30$ .



- Noisy case. Left panel: absolute error for recovering rank-three tensor with varying noise level and sample size. Right panel: absolute error for recovering rank-five tensor with varying noise level and sample size.





Thanks! and Questions?