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

Botao Hao Department of Statistics Purdue University

### 2017 International Conference on Data Science Dec. 18, 2017

Joint work with Anru Zhang, and Guang Cheng



Order-3 tensor

▲ 御 ▶ → 三 ▶

- (三) → 三三

## Tensor Data Example

#### Color image





ヘロア 人間 アメヨア 人間 アー

æ

fMRI



Botao Hao@Purdue sparse and low-rank tensor recovery

# Motivation: Compressed Image Transmission



Botao Hao@Purdue sparse and low-rank tensor recovery

# Motivation: Interaction Effect Model



source: Contraceptive Method Choice dataset from UCI

<ロ> (四) (四) (三) (三) (三)

# Motivation: Interaction Effect Model



### Sparse and Low-Rank Tensor Recovery

イロト イヨト イヨト イヨト

臣

• Observe  $\{y_i, \mathscr{X}_i\}$  from noisy cubic sketching model,



• Goal: Recover unknown third-order tensor parameter  $\mathcal{T}^*$ .

・ 同 ト ・ ヨ ト ・ ヨ ト

# Key Assumptions on Tensor Parameter

When *T*<sup>\*</sup> ∈ ℝ<sup>p×p×p</sup> is a symmetric tensor...
 CANDECOMP/PARAFAC(CP) low-rank:



2 Sparse components:  $\|\beta_k^*\|_0 \le s$  for  $k \in [K]$ .

- The cubic sketching tensor  $\mathscr{X}_i$  for symmetric case is  $\mathscr{X}_i = x_i \circ x_i \circ x_i$ , where  $\{x_i\}_{i=1}^n$  are Gaussian random vectors.
- $\beta_k^*$  and  $\beta_{k'}^*$  are not orthogonal. Different from eigenvalue decomposition in matrix case.

Key Assumptions on Tensor Parameter

When *T*<sup>\*</sup> ∈ ℝ<sup>p<sub>1</sub>×p<sub>2</sub>×p<sub>3</sub> is a non-symmetric tensor...
 CANDECOMP/PARAFAC(CP) low-rank:
</sup>



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

▲冊 ▶ ▲ 臣 ▶ ▲ 臣 ▶ 二 臣

# Reduced Symmetric Tensor Recovery Model

For symmetric tensor recovery model

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

• 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

$$y_{i} = \langle \sum_{k=1}^{K} \eta_{k}^{*} \boldsymbol{\beta}_{1k}^{*} \circ \boldsymbol{\beta}_{2k}^{*} \circ \boldsymbol{\beta}_{3k}^{*}, \boldsymbol{u}_{i} \circ \boldsymbol{v}_{i} \circ \boldsymbol{w}_{i} \rangle + \epsilon_{i}$$
$$= \sum_{k=1}^{K} \eta_{k}^{*} \underbrace{(\boldsymbol{u}_{i}^{\top} \boldsymbol{\beta}_{1k}^{*})(\boldsymbol{v}_{i}^{\top} \boldsymbol{\beta}_{2k}^{*})(\boldsymbol{w}_{i}^{\top} \boldsymbol{\beta}_{3k}^{*})}_{\text{non-linear}} + \epsilon_{i}$$

• Connect with compressed image transmission model.



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

• Consider Empirical Risk Minimization

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

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

向 ト イヨ ト イヨト

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

• 3 >

Two-stage Implementation

ヘロト 人間 とくほど 人間とう

æ

# Main Algorithm (Symmetric Recovery)



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

$$\mathcal{T}_s := \underbrace{rac{1}{6} \Big[ rac{1}{n} \sum_{i=1}^n y_i oldsymbol{x}_i \circ oldsymbol{x}_i \circ oldsymbol{x}_i - oldsymbol{\mathcal{U}} \Big]}_{i=1}$$

only depends on observations.

where the bias term  $\mathcal{U} = \sum_{j=1}^{p} \left( \boldsymbol{m}_{1} \circ \boldsymbol{e}_{j} \circ \boldsymbol{e}_{j} + \boldsymbol{e}_{j} \circ \boldsymbol{m}_{1} \circ \boldsymbol{e}_{j} + \boldsymbol{e}_{j} \circ \boldsymbol{e}_{j} \circ \boldsymbol{m}_{1} \right), \text{ and }$   $\boldsymbol{m}_{1} = \frac{1}{n} \sum_{i=1}^{n} y_{i} \boldsymbol{x}_{i}. \text{ Here } \{\boldsymbol{e}_{j}\}_{j=1}^{p} \text{ are the basis vectors in } \mathbb{R}^{p}.$ 

回り イヨト イヨト 三日

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

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



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

通 ト イ ヨ ト イ ヨ ト

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

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



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

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

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



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

• • = • • = •

# Initial Step: Sparse Tensor Decomposition

- $\Rightarrow$  Generate L staring points  $\{\beta_l^{\text{start}}\}_{l=1}^L$ .
  - $\Rightarrow$  For each starting point, compute a non-sparse component of moment-based  $T_s$  via symmetric tensor power update:

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

 $\Rightarrow$  Get a sparse solution  $\beta_l^{(t+1)}$  via thresholding or truncation.

 $\Rightarrow \text{ Cluster } L \text{ sets of single component } \{\beta_l^{(T)}, \beta_l^{(T)}, \beta_l^{(T)}\}_{l=1}^L \text{ into } K \text{ clusters to obtain a rank-} K \text{ decomposition } \{\eta_k^{(0)}, \beta_k^{(0)}, \beta_k^{(0)}, \beta_k^{(0)}\}_{k=1}^K.$ 

Different from matrix SVD due to non-orthogonality.

<sup>1</sup>For  $\mathcal{T}_s \in \mathbb{R}^{p \times p \times p}$  and  $\boldsymbol{x} \in \mathbb{R}^p$ , define  $\mathcal{T}_s \times_2 \boldsymbol{x} \times_3 \boldsymbol{x} := \sum_{j,l} \boldsymbol{x}_j \boldsymbol{x}_l [\mathcal{T}]_{:,j,l} = \mathcal{I}_s$ 

# 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

$$\widetilde{\boldsymbol{\beta}}_{k}^{(t+1)} = \boldsymbol{\beta}_{k}^{(t)} - \frac{\mu_{t}}{\phi} \nabla_{\boldsymbol{\beta}_{k}} \mathcal{L}_{1}(\boldsymbol{\eta}_{k}^{(0)}, \boldsymbol{\beta}_{k}^{(t)})$$

where  $\phi = \frac{1}{n} \sum_{i=1}^{n} y_i^2$ ,  $\mu_t$  is the step size.  $\Rightarrow$  Sparsify current update by thresholding  $\beta_k^{(t+1)} = \varphi_{\rho}(\tilde{\beta}_k^{(t+1)})$ .  $\Rightarrow$  Normalize final update  $\beta_k^{(T)} = \frac{\beta_k^{(T)}}{\|\beta_k^{(T)}\|_2}$  and update the weight  $\hat{\eta}_k = \eta_k^{(0)} \times \|\beta_k^{(T)}\|_2^3$ .

Non-asymptotic Analysis

ヘロト 人間 とくほど 人間とう

æ

#### Theorem

Suppose some regularity conditions for the true tensor parameter hold. Assume  $n \ge C_0 s^{3/2} \log p$  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, ..., the factor-wise estimator satisfies



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.
- We conjecture that  $n \gtrsim s^{3/2} \log p$  is the minimum requirement of sample complexity in most tensor problems. This has an essential difference with matrix case, where the optimal sample complexity is  $\mathcal{O}(s \log p)$ .

• When  $t \ge T$  for some enough T, the final estimator is bounded by

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

with high probability.

• Minimax optimal rate!

▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶ □

臣

• Sparse CP decomposition

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

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

$$\max_{k_i \neq k_j \in [K]} |\langle \boldsymbol{\beta}_{k_i}^*, \boldsymbol{\beta}_{k_j}^* \rangle| \le \frac{C}{\sqrt{s}}.$$

### Theorem

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,

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

### Theorem

Consider the class of tensor  $\mathcal{F}_{p,K,s}$  satisfy sparse CP-decomposition and incoherence condition. Suppose we observe n samples  $\{y_i, \mathscr{X}_i\}_{i=1}^n$  from symmetric tensor cubic sketching model, where  $n \geq Cs^{3/2} \log p$  for some large constant C. Then the estimator  $\widehat{\mathscr{T}}$ achieves

$$\inf_{\widetilde{\mathscr{T}}} \sup_{\mathscr{T}\in\mathcal{F}_{p,K,s}} \mathbb{E} \left\| \widetilde{\mathscr{T}} - \mathscr{T} 
ight\|_F^2 symp rac{\sigma^2 rac{Ks\log(p/s)}{n}}{R^*},$$

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

伺下 イヨト イヨト

- Our analysis is non-asymptotic and our estimator is rate-optimal.
- In general, we have a trade-off → 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 the combination of *truncation argument* and  $\psi_{\alpha}$ -norm.

# Numerical Study

symmetric tensor, p = 50, K = 3, s = 0.3, replication = 200.



Botao Hao@Purdue

sparse and low-rank tensor recovery



Botao Hao hao22@purdue.edu Department of Statistics Purdue University

∃ >