

Tensor completion and tensor estimation

Andrea Montanari and Nike Sun

Stanford University and UC Berkeley

June 25, 2017

When did I first meet David?

2005: MSRI?

2007? 'There are some interesting mathematical things happening'

When did I first meet David?

2005: MSRI?

2007? ‘There are some interesting mathematical things happening’

When did I first meet David?

2005: MSRI?

2007? 'There are some interesting mathematical things happening'

Tensor estimation: General question

Unknown tensor

$$\mathbf{X} \in \mathbb{R}^{d_1} \otimes \cdots \otimes \mathbb{R}^{d_k}$$

$$\mathbf{X} = (X_{i_1, \dots, i_k})_{i_1 \leq d_1, \dots, i_k \leq d_k}$$

Estimate \mathbf{X} from noisy/incomplete observations

To simplify notations...

- ▶ $d_1 = d_2 = \dots = d_k \equiv d$

- ▶ Tensors will be symmetric, e.g.:

$$X_{i_1, i_2, i_3} = X_{i_2, i_1, i_3} = X_{i_1, i_3, i_2} = \dots$$

- ▶ Results generalize.

Two concrete models:

Model #1: Spiked tensors

$$\begin{aligned} Y &= X + W \\ &= \lambda v_0^{\otimes k} + W \end{aligned}$$

Signal: $v_0 \in S^{d-1} \equiv \{x \in \mathbb{R}^d : \|x\|_2 = 1\}$.

Noise: $(W_{i_1, i_2, \dots, i_k})_{i_1 < i_2 < \dots < i_k} \sim iid N(0, 1/n)$

SNR: λ

Given Y , estimate v_0

[Montanari, Richard, 2015]

Model #1: Spiked tensors

$$\begin{aligned} Y &= X + W \\ &= \lambda v_0^{\otimes k} + W \end{aligned}$$

Signal: $v_0 \in S^{d-1} \equiv \{x \in \mathbb{R}^d : \|x\|_2 = 1\}$.

Noise: $(W_{i_1, i_2, \dots, i_k})_{i_1 < i_2 < \dots < i_k} \sim iid N(0, 1/n)$

SNR: λ

Given Y , estimate v_0

[Montanari, Richard, 2015]

Model #1: Spiked tensors

$k = 1$ (sequence model):

$$y_i = \lambda v_{0,i} + w_i$$

$k = 2$ (spiked matrix model):

$$Y_{ij} = \lambda v_{0,i} v_{0,j} + W_{ij}$$

$k = 3$ (spiked matrix model):

$$Y_{ijl} = \lambda v_{0,i} v_{0,i} v_{0,l} + W_{ijl}$$

Model #2: Tensor completion

$$\mathbf{X} = \sum_{\ell=1}^r \mathbf{v}_\ell^{\otimes k}$$

Observed entries $E \subseteq [d]^k \equiv \{1, \dots, d\}^k$

$$\Pi_E(\mathbf{X})_{i_1, \dots, i_k} = \begin{cases} X_{i_1, \dots, i_k} & \text{if } (i_1, \dots, i_k) \in E, \\ 0 & \text{otherwise} \end{cases}$$

Given $Y = \Pi_E(\mathbf{X})$, estimate \mathbf{X} .

Model #2: Tensor completion

$$\mathbf{X} = \sum_{\ell=1}^r \mathbf{v}_\ell^{\otimes k}$$

Observed entries $E \subseteq [d]^k \equiv \{1, \dots, d\}^k$

$$\Pi_E(\mathbf{X})_{i_1, \dots, i_k} = \begin{cases} X_{i_1, \dots, i_k} & \text{if } (i_1, \dots, i_k) \in E, \\ 0 & \text{otherwise} \end{cases}$$

Given $\mathbf{Y} = \Pi_E(\mathbf{X})$, estimate \mathbf{X} .

Outline

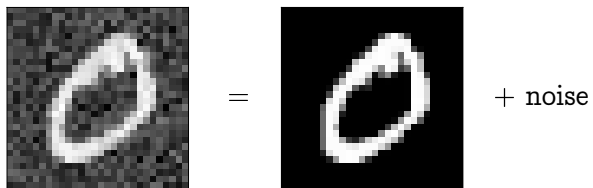
- 1 Why?
- 2 Tensor completion
- 3 Overcomplete tensors
- 4 Spiked tensor model (a.k.a. tensor PCA)
- 5 Conclusion

arXiv:1612.07866, CPAM
arXiv:1411.1076, NIPS

Why?

A cartoon application (not realistic!)

Image denoising

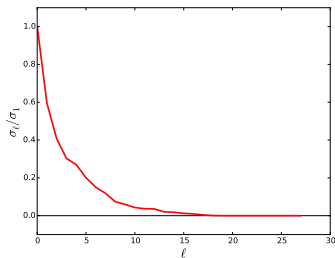


$$Y = X + W$$

Image denoising

Idea (not necessarily a good one): View X as a 28×28 matrix

Singular values:



[Dozens of references]

Singular value thresholding

Noisy image $Y \in \mathbb{R}^{28 \times 28}$

$$Y = \sum_{\ell=1}^{28} \sigma_{\ell} u_{\ell} v_{\ell}^{\top}$$

Denoised image

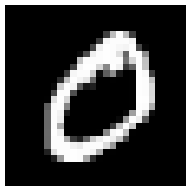
$$Y = \sum_{\ell=1}^{r_*} \sigma_{\ell} u_{\ell} v_{\ell}^{\top}$$

[...; Candés, Sing-Long, Trzasko 2013; Donoho, Gavish 2014; Chatterjee 2015
...]

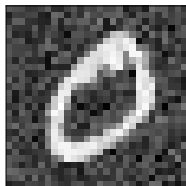
How well does it work?

$$r_0 = 12, \sigma = 30$$

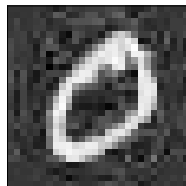
Original



Noisy

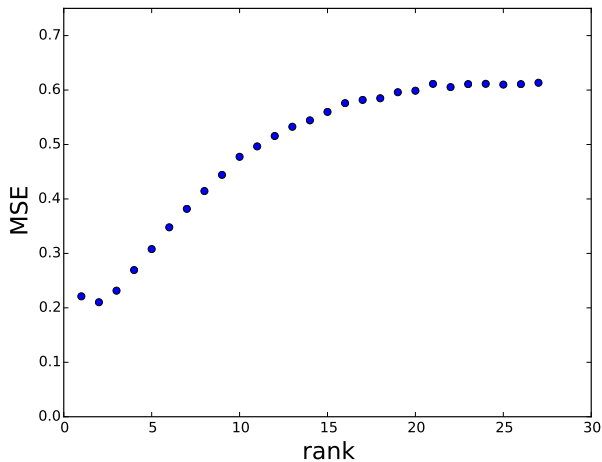


Denoised



How well does it work?

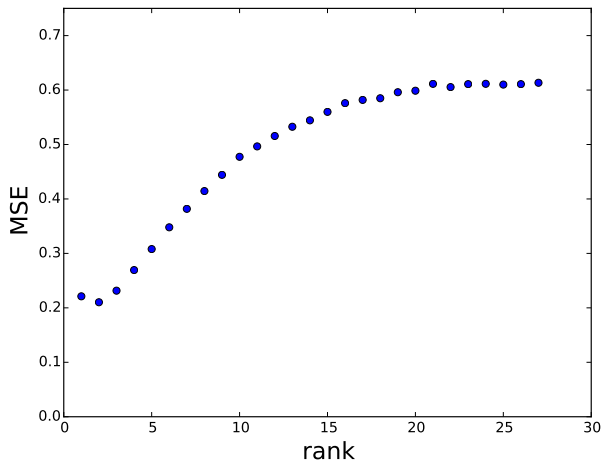
$\sigma = 100$



:(

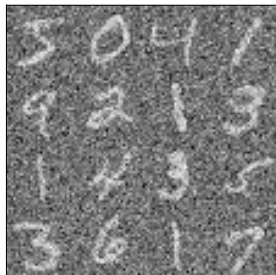
How well does it work?

$\sigma = 100$



:(

In reality we have many similar images



- ▶ $X_1, X_2, \dots, X_n \in \mathbb{R}^{d \times d}$
- ▶ Can we leverage similarities between images?

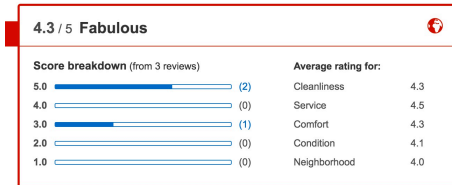
Idea

1. Stack images in a tensor:

$$\mathbf{X} = [\mathbf{X}_1 | \mathbf{X}_2 | \cdots | \mathbf{X}_n] \in \mathbb{R}^d \otimes \mathbb{R}^d \otimes \mathbb{R}^n$$

2. Do tensor denoising

A better application



Collaborative filtering:

Users \times Hotel \times Keyword

Tensor completion

Model

- ▶ Unknown tensor $\mathbf{X} \in (\mathbb{R}^d)^{\otimes k}$:

$$\mathbf{X} = \sum_{\ell=1}^r \mathbf{v}_{\ell}^{\otimes k}$$

- ▶ Entries $E \subseteq [d]^k$, $|E| = n$ uniformly random
- ▶ Observations

$$\mathbf{Y} = \Pi_E(\mathbf{X})$$

Number of parameters $\approx r d$

Model

- ▶ Unknown tensor $\mathbf{X} \in (\mathbb{R}^d)^{\otimes k}$:

$$\mathbf{X} = \sum_{\ell=1}^r \mathbf{v}_\ell^{\otimes k}$$

- ▶ Entries $E \subseteq [d]^k$, $|E| = n$ uniformly random
- ▶ Observations

$$\mathbf{Y} = \Pi_E(\mathbf{X})$$

Number of parameters $\approx r d$

What do we know about matrices? ($k = 2$)

Theorem (Gross 2011)

If $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$, and $\text{rank}(\mathbf{X}) \leq r$ with factors factors are ' μ -incoherent', then we can reconstruct \mathbf{X} exactly from

$$n \geq C(\mu) r (d_1 \vee d_2) (\log(d_1 \vee d_2))^2$$

random entries. This is achieved by nuclear norm minimization.

[Candés, Recht, 2009; Candés, Tao, 2010; Keshavan, Montanari, Oh, 2010; Recht 2011; ...]

What do we know about matrices? ($k = 2$)

Theorem (Gross 2011)

If $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$, and $\text{rank}(\mathbf{X}) \leq r$ with factors factors are ' μ -incoherent', then we can reconstruct \mathbf{X} exactly from

$$n \geq C(\mu) r (d_1 \vee d_2) (\log(d_1 \vee d_2))^2$$

random entries. This is achieved by nuclear norm minimization.

[Candés, Recht, 2009; Candés, Tao, 2010; Keshavan, Montanari, Oh, 2010; Recht 2011; ...]

Let us try to redo the same with tensors!

- ▶ Number of unknown parameters $\approx rd$
- ▶ For $r = O(1)$ and nice factors, maximum likelihood succeeds for

$$n \geq Crd \log d$$

[Do not know a reference]

- ▶ Generalizations of nuclear norm

$$n \geq C(r_{\boxplus}^{k-1} \vee r_{\boxplus}^{(k-1)/2} d^{1/2}) d(\log d)^2$$

[Yuan, Zhang, 2015; 2016]

Evaluating tensor nuclear norm is NP-hard!

Let us try to redo the same with tensors!

- ▶ Number of unknown parameters $\approx rd$
- ▶ For $r = O(1)$ and nice factors, maximum likelihood succeeds for

$$n \geq Crd \log d$$

[Do not know a reference]

- ▶ Generalizations of nuclear norm

$$n \geq C(r_{\boxplus}^{k-1} \vee r_{\boxplus}^{(k-1)/2} d^{1/2}) d(\log d)^2$$

[Yuan, Zhang, 2015; 2016]

Evaluating tensor nuclear norm is NP-hard!

Let us try to redo the same with tensors!

- ▶ Number of unknown parameters $\approx rd$
- ▶ For $r = O(1)$ and nice factors, maximum likelihood succeeds for

$$n \geq Crd \log d$$

[Do not know a reference]

- ▶ Generalizations of nuclear norm

$$n \geq C(r_{\boxplus}^{k-1} \vee r_{\boxplus}^{(k-1)/2} d^{1/2}) d(\log d)^2$$

[Yuan, Zhang, 2015; 2016]

Evaluating tensor nuclear norm is NP-hard!

Let us try to redo the same with tensors!

- ▶ Number of unknown parameters $\approx rd$
- ▶ For $r = O(1)$ and nice factors, maximum likelihood succeeds for

$$n \geq Crd \log d$$

[Do not know a reference]

- ▶ Generalizations of nuclear norm

$$n \geq C(r_{\boxplus}^{k-1} \vee r_{\boxplus}^{(k-1)/2} d^{1/2}) d(\log d)^2$$

[Yuan, Zhang, 2015; 2016]

Evaluating tensor nuclear norm is NP-hard!

Let us try to redo the same with tensors!

- ▶ Number of unknown parameters $\approx rd$
- ▶ For $r = O(1)$ and nice factors, maximum likelihood succeeds for

$$n \geq Crd \log d$$

[Do not know a reference]

- ▶ Generalizations of nuclear norm

$$n \geq C(r_{\boxplus}^{k-1} \vee r_{\boxplus}^{(k-1)/2} d^{1/2}) d(\log d)^2$$

[Yuan, Zhang, 2015; 2016]

Evaluating tensor nuclear norm is NP-hard!

Idea #1: What about reducing to $k = 2$?

$$k = a + b$$

$$\text{unfold}^{a \times b} : (\mathbb{R}^d)^{\otimes k} \rightarrow \mathbb{R}^{d^a} \otimes \mathbb{R}^{d^b}$$

$$Y \rightarrow \text{unfold}^{a \times b}(Y)$$

$$\text{unfold}^{a \times b}(Y)_{i_1 \dots i_a; j_1 \dots j_b} \equiv Y_{i_1 \dots i_a j_1 \dots j_b}$$

- ▶ $M = \text{unfold}^{a \times b}(Y)$
- ▶ Do matrix completion (e.g. nuclear norm minimization) of M
- ▶ Fold back

[Tomioka, Hayashi, Kashima, 2010; Tomioka, Suzuki, Hayashi, Kashima, 2011; Liu, Musialski, Wonka, Ye, 2013; Gandy, Recht, Yamada, 2011]

Idea #1: What about reducing to $k = 2$?

$$k = a + b$$

$$\text{unfold}^{a \times b} : (\mathbb{R}^d)^{\otimes k} \rightarrow \mathbb{R}^{d^a} \otimes \mathbb{R}^{d^b}$$

$$Y \rightarrow \text{unfold}^{a \times b}(Y)$$

$$\text{unfold}^{a \times b}(Y)_{i_1 \dots i_a; j_1 \dots j_b} \equiv Y_{i_1 \dots i_a j_1 \dots j_b}$$

- ▶ $M = \text{unfold}^{a \times b}(Y)$
- ▶ Do matrix completion (e.g. nuclear norm minimization) of M
- ▶ Fold back

[Tomioka, Hayashi, Kashima, 2010; Tomioka, Suzuki, Hayashi, Kashima, 2011; Liu, Musialski, Wonka, Ye, 2013; Gandy, Recht, Yamada, 2011]

Idea #1: What about reducing to $k = 2$?

$$k = a + b$$

$$\text{unfold}^{a \times b} : (\mathbb{R}^d)^{\otimes k} \rightarrow \mathbb{R}^{d^a} \otimes \mathbb{R}^{d^b}$$

$$Y \rightarrow \text{unfold}^{a \times b}(Y)$$

$$\text{unfold}^{a \times b}(Y)_{i_1 \dots i_a; j_1 \dots j_b} \equiv Y_{i_1 \dots i_a j_1 \dots j_b}$$

- ▶ $M = \text{unfold}^{a \times b}(Y)$
- ▶ Do matrix completion (e.g. nuclear norm minimization) of M
- ▶ Fold back

[Tomioka, Hayashi, Kashima, 2010; Tomioka, Suzuki, Hayashi, Kashima, 2011; Liu, Musialski, Wonka, Ye, 2013; Gandy, Recht, Yamada, 2011]

Idea #1: What about reducing to $k = 2$?

Corollary

If $\mathbf{X} \in (\mathbb{R}^d)^{\otimes k}$, $\text{rank}(\mathbf{X}) \leq r$, is such that $\text{unfold}^{a \times b}(\mathbf{X})$ satisfies incoherence, then it can be reconstructed whp from n random entries, provided

$$n \geq C r d^{a \vee b} (\log d)^2$$

Insights (?)

- ▶ Optimal choice $a = \lfloor k/2 \rfloor$, $b = \lceil k/2 \rceil$.
- ▶ Gap: $rd \ll n \ll rd^{\lceil k/2 \rceil}$.
- ▶ Unfolding cannot beat the barrier $n \gtrsim rd^{\lceil k/2 \rceil}$.

Idea #1: What about reducing to $k = 2$?

Corollary

If $\mathbf{X} \in (\mathbb{R}^d)^{\otimes k}$, $\text{rank}(\mathbf{X}) \leq r$, is such that $\text{unfold}^{a \times b}(\mathbf{X})$ satisfies incoherence, then it can be reconstructed whp from n random entries, provided

$$n \geq C r d^{a \vee b} (\log d)^2$$

Insights (?)

- ▶ Optimal choice $a = \lfloor k/2 \rfloor$, $b = \lceil k/2 \rceil$.
- ▶ Gap: $rd \ll n \ll rd^{\lceil k/2 \rceil}$.
- ▶ Unfolding cannot beat the barrier $n \gtrsim rd^{\lceil k/2 \rceil}$.

Idea #2: Non-convex optimization

$$\mathbf{X} = \mathbf{v}_0^{\otimes k}, \|\mathbf{v}_0\|_2 = 1$$

Maximum likelihood

$$\begin{aligned} &\text{maximize } \mathcal{L}_n(\boldsymbol{\theta}) = \frac{1}{2} \|\Pi_E(\mathbf{Y} - \boldsymbol{\theta}^{\otimes k})\|_F^2, \\ &\text{subject to } \|\boldsymbol{\theta}\|_2 = 1. \end{aligned}$$

Heuristic analysis

$$\nabla \mathcal{L}_n(\boldsymbol{\theta}) = -k \Pi_E(\mathbf{Y}) \{\boldsymbol{\theta}^{\otimes(k-1)}\} + \text{smaller terms}$$

Idea #2: Non-convex optimization

$$\mathbf{X} = \mathbf{v}_0^{\otimes k}, \|\mathbf{v}_0\|_2 = 1$$

Maximum likelihood

$$\begin{aligned} &\text{maximize } \mathcal{L}_n(\boldsymbol{\theta}) = \frac{1}{2} \|\Pi_E(\mathbf{Y} - \boldsymbol{\theta}^{\otimes k})\|_F^2, \\ &\text{subject to } \|\boldsymbol{\theta}\|_2 = 1. \end{aligned}$$

Heuristic analysis

$$\nabla \mathcal{L}_n(\boldsymbol{\theta}) = -k \Pi_E(\mathbf{Y}) \{\boldsymbol{\theta}^{\otimes(k-1)}\} + \text{smaller terms}$$

Idea #2: Non-convex optimization

$$\mathbf{X} = \mathbf{v}_0^{\otimes k}, \|\mathbf{v}_0\|_2 = 1$$

Maximum likelihood

$$\begin{aligned} & \text{maximize } \mathcal{L}_n(\boldsymbol{\theta}) = \frac{1}{2} \|\Pi_E(\mathbf{Y} - \boldsymbol{\theta}^{\otimes k})\|_F^2, \\ & \text{subject to } \|\boldsymbol{\theta}\|_2 = 1. \end{aligned}$$

Heuristic analysis

$$\nabla \mathcal{L}_n(\boldsymbol{\theta}) = -k \Pi_E(\mathbf{Y}) \{\boldsymbol{\theta}^{\otimes(k-1)}\} + \text{smaller terms}$$

Idea #2: Non-convex optimization

$$\begin{aligned}\Pi_E(\mathbf{Y}) &= \mathbb{E}\Pi_E(\mathbf{Y}) + \{\Pi_E(\mathbf{Y}) - \mathbb{E}\Pi_E(\mathbf{Y})\} \\ &= \frac{n}{d^k} \mathbf{v}_0^{\otimes k} + \mathbf{W}\end{aligned}$$

Assume random initialization $\langle \boldsymbol{\theta}, \mathbf{v}_0 \rangle = cd^{-1/2}$, $c = O(1)$:

$$\begin{aligned}\langle \mathbf{v}_0, \nabla \mathcal{L}_n(\boldsymbol{\theta}) \rangle &= -k \langle \Pi_E(\mathbf{Y}), \mathbf{v}_0 \otimes \boldsymbol{\theta}^{\otimes(k-1)} \rangle + \text{smaller terms} \\ &= -\frac{kn}{d^k} \langle \mathbf{v}_0, \boldsymbol{\theta} \rangle^{k-1} - k \langle \mathbf{W}, \mathbf{v}_0 \otimes \boldsymbol{\theta}^{\otimes(k-1)} \rangle + \dots \\ &= -ck \frac{n}{d^{(3k-1)/2}} \pm \frac{n^{1/2}}{d^k}\end{aligned}$$

Gradient points in a random direction unless $n \gtrsim d^{k-1}$

Is there a practical estimator for $n \leq rd^{1.1}$?

Barak, Moitra, 2014, $k = 3$:

- ▶ Under 'Feige's hypothesis,' no polynomial algorithm exists for

$$n \ll d^{3/2}$$

- ▶ Degree-6 Sum-Of-Squares works (approximate reconstruction) if

$$n \geq Cr_*^2 d^{3/2} (\log d)^4$$

- ▶ Complexity $O(d^{15})$

Practical algorithms? Better rank dependency?

Is there a practical estimator for $n \leq rd^{1.1}$?

Barak, Moitra, 2014, $k = 3$:

- ▶ Under 'Feige's hypothesis,' no polynomial algorithm exists for

$$n \ll d^{3/2}$$

- ▶ Degree-6 Sum-Of-Squares works (approximate reconstruction) if

$$n \geq Cr_*^2 d^{3/2} (\log d)^4$$

- ▶ Complexity $O(d^{15})$

Practical algorithms? Better rank dependency?

Is there a practical estimator for $n \leq rd^{1.1}$?

Barak, Moitra, 2014, $k = 3$:

- ▶ Under 'Feige's hypothesis,' no polynomial algorithm exists for

$$n \ll d^{3/2}$$

- ▶ Degree-6 Sum-Of-Squares works (approximate reconstruction) if

$$n \geq Cr_*^2 d^{3/2} (\log d)^4$$

- ▶ Complexity $O(d^{15})$

Practical algorithms? Better rank dependency?

Is there a practical estimator for $n \leq rd^{1.1}$?

Barak, Moitra, 2014, $k = 3$:

- ▶ Under 'Feige's hypothesis,' no polynomial algorithm exists for

$$n \ll d^{3/2}$$

- ▶ Degree-6 Sum-Of-Squares works (approximate reconstruction) if

$$n \geq Cr_*^2 d^{3/2} (\log d)^4$$

- ▶ Complexity $O(d^{15})$

Practical algorithms? Better rank dependency?

Is there a practical estimator for $n \leq rd^{1.1}$?

Barak, Moitra, 2014, $k = 3$:

- ▶ Under 'Feige's hypothesis,' no polynomial algorithm exists for

$$n \ll d^{3/2}$$

- ▶ Degree-6 Sum-Of-Squares works (approximate reconstruction) if

$$n \geq Cr_*^2 d^{3/2} (\log d)^4$$

- ▶ Complexity $O(d^{15})$

Practical algorithms? Better rank dependency?

Is there a practical estimator for $n \leq rd^{1.1}$?

Barak, Moitra, 2014, $k = 3$:

- ▶ Under 'Feige's hypothesis,' no polynomial algorithm exists for

$$n \ll d^{3/2}$$

- ▶ Degree-6 Sum-Of-Squares works (approximate reconstruction) if

$$n \geq Cr_*^2 d^{3/2} (\log d)^4$$

- ▶ Complexity $O(d^{15})$

Practical algorithms? Better rank dependency?

A theorem

$$\mathbf{A1.} \quad \max_{i_1 \dots i_k} |X_{i_1 \dots i_k}|^2 \leq \frac{\alpha}{d^k} \|\mathbf{X}\|_F^2$$

$$\mathbf{A2.} \quad \|\text{unfold}^{a \times b}(\mathbf{X})\|_{\text{op}}^2 \leq \frac{\mu}{r} \|\text{unfold}^{a \times b}(\mathbf{X})\|_F^2$$

Theorem (Montanari, Sun, 2017)

Under the above assumptions there exists a spectral algorithm achieving $\|\widehat{\mathbf{X}}(\mathbf{Y}) - \mathbf{X}\|_F \leq \varepsilon \|\mathbf{X}\|_F$ whp, provided^a $r \leq d^{c(k)}$ and

$$n \geq C(\mu, \alpha, \varepsilon) r d^{k/2} (\log d)^9.$$

^a $c(3) = 3/4$, $c(k) = k/2$ (k even), $c(k) = (k/2) - 1$ ($k \geq 5$ odd).

A theorem

$$\mathbf{A1.} \quad \max_{i_1 \dots i_k} |X_{i_1 \dots i_k}|^2 \leq \frac{\alpha}{d^k} \|\mathbf{X}\|_F^2$$

$$\mathbf{A2.} \quad \|\text{unfold}^{a \times b}(\mathbf{X})\|_{\text{op}}^2 \leq \frac{\mu}{r} \|\text{unfold}^{a \times b}(\mathbf{X})\|_F^2$$

Theorem (Montanari, Sun, 2017)

Under the above assumptions there exists a spectral algorithm achieving $\|\widehat{\mathbf{X}}(\mathbf{Y}) - \mathbf{X}\|_F \leq \varepsilon \|\mathbf{X}\|_F$ whp, provided^a $r \leq d^{c(k)}$ and

$$n \geq C(\mu, \alpha, \varepsilon) r d^{k/2} (\log d)^9.$$

^a $c(3) = 3/4$, $c(k) = k/2$ (k even), $c(k) = (k/2) - 1$ ($k \geq 5$ odd).

A theorem

Theorem (Montanari, Sun, 2017)

Under the above assumptions there exists a spectral algorithm achieving $\|\widehat{\mathbf{X}}(\mathbf{Y}) - \mathbf{X}\|_F \leq \varepsilon \|\mathbf{X}\|_F$ whp, provided^a $r \leq d^{c(k)}$ and

$$n \geq C(\mu, \alpha, \varepsilon) r d^{k/2} (\log d)^9.$$

^a $c(3) = 3/4$, $c(k) = k/2$ (k even), $c(k) = (k/2) - 1$ ($k \geq 5$ odd).

- ▶ FALSE: Optimal choice $a = \lfloor k/2 \rfloor$, $b = \lceil k/2 \rceil$.
- ▶ FALSE: Gap: $rd \ll n \ll rd^{\lceil k/2 \rceil}$.
- ▶ FALSE: Unfolding cannot beat the barrier $n \gtrsim rd^{\lceil k/2 \rceil}$.

A theorem

Theorem (Montanari, Sun, 2017)

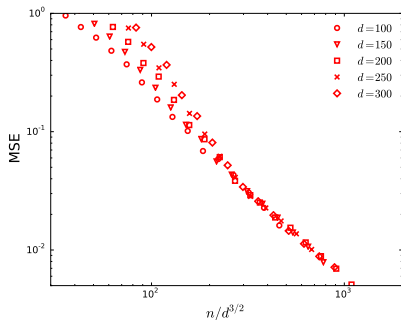
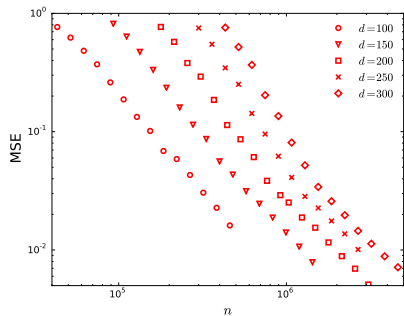
Under the above assumptions there exists a spectral algorithm achieving $\|\widehat{\mathbf{X}}(\mathbf{Y}) - \mathbf{X}\|_F \leq \varepsilon \|\mathbf{X}\|_F$ whp, provided^a $r \leq d^{c(k)}$ and

$$n \geq C(\mu, \alpha, \varepsilon) rd^{k/2}(\log d)^9.$$

^a $c(3) = 3/4$, $c(k) = k/2$ (k even), $c(k) = (k/2) - 1$ ($k \geq 5$ odd).

- ▶ **FALSE:** Optimal choice $a = \lfloor k/2 \rfloor$, $b = \lceil k/2 \rceil$.
- ▶ **FALSE:** Gap: $rd \ll n \ll rd^{\lceil k/2 \rceil}$.
- ▶ **FALSE:** Unfolding cannot beat the barrier $n \gtrsim rd^{\lceil k/2 \rceil}$.

Simulations: $k = 3, r = 4$



Algorithm idea

$$k = a + b, a < b$$

$$M = \text{unfold}^{a \times b}(Y)$$

$$d^a \ll n \ll d^{k/2}.$$

- ▶ Cannot complete M
- ▶ Can estimate the top left singular space!

Algorithm, for $k = 3$ ($\delta = nd^3$)

1. Compute $M = \text{unfold}^{1 \times 2}(X)$ and

$$A = \frac{1}{\delta^2} \Pi_{\text{diag}}^{\perp}(M M^{\top}) + \frac{1}{\delta} \Pi_{\text{diag}}(M M^{\top})$$

2. Eigenvalue decomposition:

$$A = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top}$$

3. Projector

$$Q = \sum_{i=1}^d \mathbf{1}_{\lambda_i \geq \lambda_*} \mathbf{u}_i \mathbf{u}_i^{\top}$$

4. Let $\mathcal{Q} \equiv Q \otimes Q \otimes Q$ and return

$$\widehat{X} = \frac{1}{\delta} \mathcal{Q}(X)$$

Algorithm, for $k = 3$ ($\delta = nd^3$)

1. Compute $M = \text{unfold}^{1 \times 2}(X)$ and

$$A = \frac{1}{\delta^2} \Pi_{\text{diag}}^{\perp}(M M^{\top}) + \frac{1}{\delta} \Pi_{\text{diag}}(M M^{\top})$$

2. Eigenvalue decomposition:

$$A = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top}$$

3. Projector

$$Q = \sum_{i=1}^d \mathbf{1}_{\lambda_i \geq \lambda_*} \mathbf{u}_i \mathbf{u}_i^{\top}$$

4. Let $\mathcal{Q} \equiv Q \otimes Q \otimes Q$ and return

$$\widehat{X} = \frac{1}{\delta} \mathcal{Q}(X)$$

Algorithm, for $k = 3$ ($\delta = nd^3$)

1. Compute $M = \text{unfold}^{1 \times 2}(X)$ and

$$A = \frac{1}{\delta^2} \Pi_{\text{diag}}^{\perp}(M M^{\top}) + \frac{1}{\delta} \Pi_{\text{diag}}(M M^{\top})$$

2. Eigenvalue decomposition:

$$A = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top}$$

3. Projector

$$Q = \sum_{i=1}^d \mathbf{1}_{\lambda_i \geq \lambda_*} \mathbf{u}_i \mathbf{u}_i^{\top}$$

4. Let $\mathcal{Q} \equiv Q \otimes Q \otimes Q$ and return

$$\widehat{X} = \frac{1}{\delta} \mathcal{Q}(X)$$

Algorithm, for $k = 3$ ($\delta = nd^3$)

1. Compute $M = \text{unfold}^{1 \times 2}(X)$ and

$$A = \frac{1}{\delta^2} \Pi_{\text{diag}}^{\perp}(M M^{\top}) + \frac{1}{\delta} \Pi_{\text{diag}}(M M^{\top})$$

2. Eigenvalue decomposition:

$$A = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top}$$

3. Projector

$$Q = \sum_{i=1}^d \mathbf{1}_{\lambda_i \geq \lambda_*} \mathbf{u}_i \mathbf{u}_i^{\top}$$

4. Let $\mathcal{Q} \equiv Q \otimes Q \otimes Q$ and return

$$\widehat{X} = \frac{1}{\delta} \mathcal{Q}(X)$$

Algorithm, for $k = 3$ ($\delta = nd^3$)

1. Compute $M = \text{unfold}^{1 \times 2}(X)$ and

$$A = \frac{1}{\delta^2} \Pi_{\text{diag}}^{\perp}(M M^{\top}) + \frac{1}{\delta} \Pi_{\text{diag}}(M M^{\top})$$

2. Eigenvalue decomposition:

$$A = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top}$$

3. Projector

$$Q = \sum_{i=1}^d \mathbf{1}_{\lambda_i \geq \lambda_*} \mathbf{u}_i \mathbf{u}_i^{\top}$$

4. Let $\mathcal{Q} \equiv Q \otimes Q \otimes Q$ and return

$$\widehat{X} = \frac{1}{\delta} \mathcal{Q}(X)$$

Similar method developed independently by Yuan, Xia, 2017

Back to image denoising

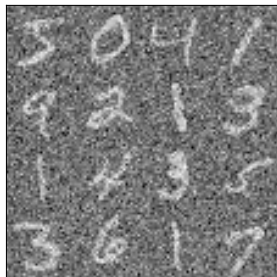
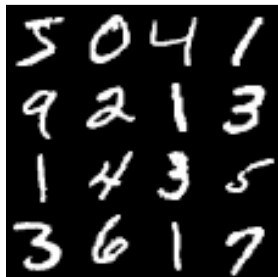
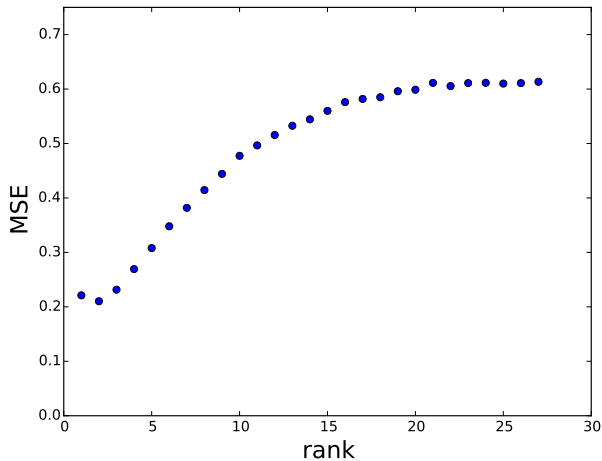
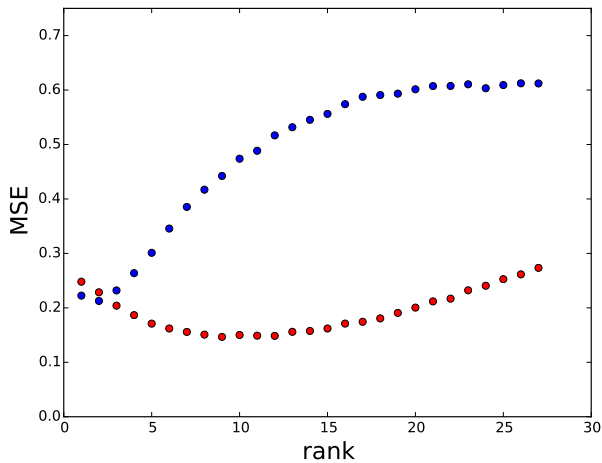


Image-by-image matrix denoising



:(

Tensor denoising



:)

Overcomplete tensors

$$k = 3$$

$$X = \sum_{\ell=1}^r v_{\ell} \otimes v_{\ell} \otimes v_{\ell}$$

X can have rank larger than $d!!!$

Can we use spectral methods?

$$Y = \Pi_E(X)$$

$$M = \text{unfold}^{1 \times 2}(Y)$$

Seems a lost cause:

$$\text{rank}(X) \geq d \quad \Rightarrow \quad \text{rank}(M) = d$$

Structure is lost!

Can we use spectral methods?

$$Y = \Pi_E(\mathbf{X})$$

$$M = \text{unfold}^{1 \times 2}(Y)$$

Seems a lost cause:

$$\text{rank}(\mathbf{X}) \geq d \quad \Rightarrow \quad \text{rank}(M) = d$$

Structure is lost!

Can we use spectral methods?

$$Y = \Pi_E(\mathbf{X})$$

$$M = \text{unfold}^{1 \times 2}(Y)$$

Seems a lost cause:

$$\text{rank}(\mathbf{X}) \geq d \quad \Rightarrow \quad \text{rank}(M) = d$$

Structure is lost!

Idea: Do something before spectral analysis

$$A = (A_{i_1 i_2 ; j_1 j_2})_{i_1, i_2, j_1, j_2 \leq d} \in \mathbb{R}^{d^2 \times d^2},$$
$$A_{i_1 i_2 ; j_1 j_2} \equiv \sum_{l=1}^d Y_{i_1 j_1 l} Y_{i_2 j_2 l}.$$

NOT EQUAL TO:

$$M = \text{unfold}^{1 \times 2}(Y) \in \mathbb{R}^{d \times d^2},$$
$$B = M^T M$$

Claim: $\text{rank}_*(A) \approx r$, $\text{rank}(B) \leq d$

Idea: Do something before spectral analysis

$$\mathbf{A} = (A_{i_1 i_2 ; j_1 j_2})_{i_1, i_2, j_1, j_2 \leq d} \in \mathbb{R}^{d^2 \times d^2},$$
$$A_{i_1 i_2 ; j_1 j_2} \equiv \sum_{\ell=1}^d Y_{i_1 j_1 \ell} Y_{i_2 j_2 \ell}.$$

NOT EQUAL TO:

$$M = \text{unfold}^{1 \times 2}(Y) \in \mathbb{R}^{d \times d^2},$$
$$B = M^T M$$

Claim: $\text{rank}_*(A) \approx r, \text{rank}(B) \leq d$

Idea: Do something before spectral analysis

$$\mathbf{A} = (A_{i_1 i_2 ; j_1 j_2})_{i_1, i_2, j_1, j_2 \leq d} \in \mathbb{R}^{d^2 \times d^2},$$
$$A_{i_1 i_2 ; j_1 j_2} \equiv \sum_{\ell=1}^d Y_{i_1 j_1 \ell} Y_{i_2 j_2 \ell}.$$

NOT EQUAL TO:

$$\mathbf{M} = \text{unfold}^{1 \times 2}(\mathbf{Y}) \in \mathbb{R}^{d \times d^2},$$
$$\mathbf{B} = \mathbf{M}^\top \mathbf{M}$$

Claim: $\text{rank}_*(\mathbf{A}) \approx r$, $\text{rank}(\mathbf{B}) \leq d$

$$\mathbf{X} = \sum_{\ell=1}^r \mathbf{v}_\ell^{\otimes 3}$$

Theorem (Montanari, Sun, 2017)

Assume $(\mathbf{v}_\ell)_{\ell \leq r} \sim_{iid} \mathcal{N}(0, \mathbf{I}_d)$, and $r \leq d^2$. Then there exists a spectral algorithm achieving $\|\widehat{\mathbf{X}} - \mathbf{X}\|_F \leq \varepsilon \|\mathbf{X}\|_F$ whp, for n random entries, provided

$$n \geq C(\varepsilon)(d \vee r)d^{3/2}(\log d)^{c_0}$$

Similar ideas: Hopkins, Schramm, Shi, Steurer 2016; Raghavendra, Schramm 2016

$$\mathbf{X} = \sum_{\ell=1}^r \mathbf{v}_\ell^{\otimes 3}$$

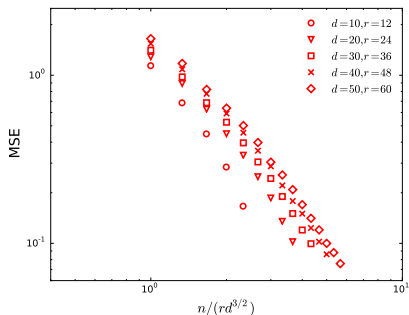
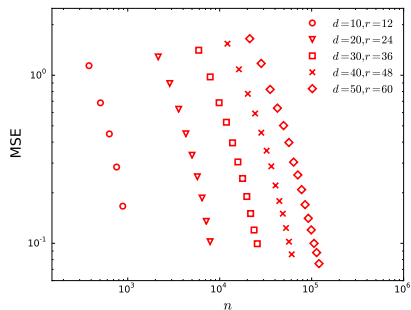
Theorem (Montanari, Sun, 2017)

Assume $(\mathbf{v}_\ell)_{\ell \leq r} \sim_{iid} \mathcal{N}(0, \mathbf{I}_d)$, and $r \leq d^2$. Then there exists a spectral algorithm achieving $\|\widehat{\mathbf{X}} - \mathbf{X}\|_F \leq \varepsilon \|\mathbf{X}\|_F$ whp, for n random entries, provided

$$n \geq C(\varepsilon)(d \vee r)d^{3/2}(\log d)^{c_0}$$

Similar ideas: Hopkins, Schramm, Shi, Steurer 2016; Raghavendra, Schramm 2016

Simulations: $k = 3, r = 4$



Spiked tensor model (a.k.a. tensor PCA)

Reminder: A much simpler model

$$Y = \lambda v_0^{\otimes k} + W$$

Signal: $v_0 \in S^{d-1} \equiv \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_2 = 1\}$.

Noise: $(W_{i_1, i_2, \dots, i_k})_{i_1 < i_2 < \dots < i_k} \sim_{iid} N(0, 1/d)$

SNR: λ

Given Y , estimate v_0

[Montanari, Richard, 2015]

A lot of information from statistical physics

Gibbs measure

$$\mu_{\beta, \lambda}(d\theta) = \frac{1}{Z(\beta, \lambda)} \exp \left\{ \beta \langle Y, \theta^{\otimes k} \rangle \right\} \mu_0(d\theta)$$

- ▶ $\beta = \infty$: Maximum Likelihood
- ▶ $\beta = \lambda/k!$: Bayes posterior
- ▶ $\mu_0(\cdot)$: Uniform measure on S^{d-1}

[Crisanti, Sommers, 1992, 1995; Auffinger, Ben Arous, Cerny, 2013; Chen, 2013; Subag, 2016; Krzakala, Lelarge, Miolane, Zdeborova, 2016; ...]

Theorem

There exists $\lambda_{\text{Bayes}}(k)$ (explicit!) such that:

$$\lim_{n \rightarrow \infty} \mathbb{E} |\langle v_0, v_{\text{Bayes}}(Y) \rangle| = \begin{cases} 0 & \text{if } \lambda < \lambda_{\text{Bayes}}(k), \\ > 0 & \text{if } \lambda > \lambda_{\text{Bayes}}(k). \end{cases}$$

A lot of information from statistical physics

Gibbs measure

$$\mu_{\beta, \lambda}(d\theta) = \frac{1}{Z(\beta, \lambda)} \exp \left\{ \beta \langle Y, \theta^{\otimes k} \rangle \right\} \mu_0(d\theta)$$

- ▶ $\beta = \infty$: Maximum Likelihood
- ▶ $\beta = \lambda/k!$: Bayes posterior
- ▶ $\mu_0(\cdot)$: Uniform measure on S^{d-1}

[Crisanti, Sommers, 1992, 1995; Auffinger, Ben Arous, Cerny, 2013; Chen, 2013; Subag, 2016; Krzakala, Lelarge, Miolane, Zdeborova, 2016; ...]

Theorem

There exists $\lambda_{\text{Bayes}}(k)$ (explicit!) such that:

$$\lim_{n \rightarrow \infty} \mathbb{E} |\langle v_0, v_{\text{Bayes}}(\mathbf{Y}) \rangle| = \begin{cases} 0 & \text{if } \lambda < \lambda_{\text{Bayes}}(k), \\ > 0 & \text{if } \lambda > \lambda_{\text{Bayes}}(k). \end{cases}$$

What about polynomial-time estimators?

Theorem (Montanari, Richard, 2014; Hopkins, Shi, Steurer, 2015)

There exists poly-time estimator achieving $\mathbb{E}\{|\langle \hat{\theta}^{\text{Poly}}, \theta_0 \rangle|\} \geq 1 - \varepsilon$, provided

$$\lambda \geq C(\varepsilon)n^{(k-2)/4}.$$

No poly-time algorithm known for $1 \ll \lambda \ll n^{(k-2)/4}$.

What about polynomial-time estimators?

Theorem (Montanari, Richard, 2014; Hopkins, Shi, Steurer, 2015)

There exists poly-time estimator achieving $\mathbb{E}\{|\langle \hat{\theta}^{\text{Poly}}, \theta_0 \rangle|\} \geq 1 - \varepsilon$, provided

$$\lambda \geq C(\varepsilon)n^{(k-2)/4}.$$

No poly-time algorithm known for $1 \ll \lambda \ll n^{(k-2)/4}$.

What about polynomial-time estimators?

Theorem (Montanari, Richard, 2014; Hopkins, Shi, Steurer, 2015)

There exists poly-time estimator achieving $\mathbb{E}\{|\langle \hat{\theta}^{\text{Poly}}, \theta_0 \rangle|\} \geq 1 - \varepsilon$, provided

$$\lambda \geq C(\varepsilon)n^{(k-2)/4}.$$

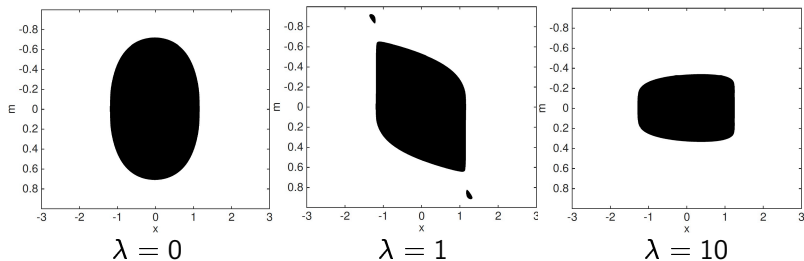
No poly-time algorithm known for $1 \ll \lambda \ll n^{(k-2)/4}$.

A case study in failure: Maximum likelihood

$$\begin{aligned} & \text{maximize} && \langle Y, \boldsymbol{\theta}^{\otimes k} \rangle \\ & \text{subject to} && \boldsymbol{\theta} \in S^{d-1} \end{aligned}$$

$$N(x, m) \equiv \#\left\{ \text{critical points with } \langle Y, \boldsymbol{\theta}^{\otimes k} \rangle \approx x, \langle \mathbf{v}_0, \boldsymbol{\theta} \rangle \approx m \right\}$$

A peek at complexity



$$\mathbb{E}N(x, m) = e^{d\Phi(x, m) + o(d)},$$
$$\Phi(x, m) = \text{explicit expression}$$

[Ben Arous, Mei, Montanari, Nica, in progress]

Conclusion

Conclusion

- ▶ Tensors are useful for modeling multivariate data
- ▶ Estimation requires entirely new ideas
- ▶ Information-computation gap
- ▶ Many open problems

Thanks! Happy birthday Dave!

Conclusion

- ▶ Tensors are useful for modeling multivariate data
- ▶ Estimation requires entirely new ideas
- ▶ Information-computation gap
- ▶ Many open problems

Thanks! Happy birthday Dave!