

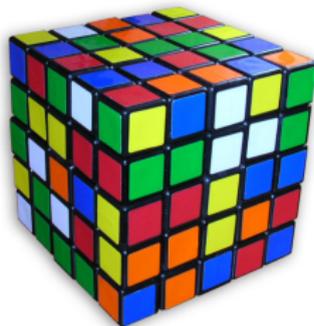
Low Rank Tensor Completion

Ming Yuan

Department of Statistics
Columbia University

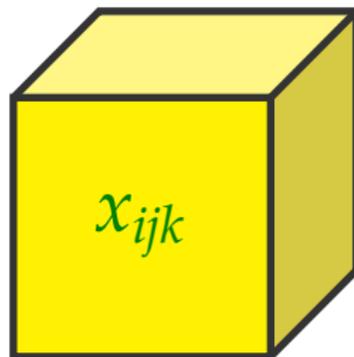
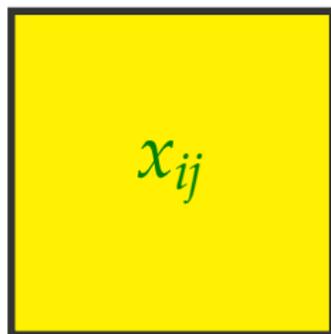
ming.yuan@columbia.edu

<http://www.columbia.edu/~my2550>



(Joint work with Dong Xia and Cun-hui Zhang)

DATA IN THE FORM OF MULTILINEAR ARRAY



- ▶ Spatio-temporal expression data
[e.g. Kang et al. (2011); Parikshak et al. (2013); Hawrylycz et al. (2015)]
- ▶ Imaging (video) data – 3D images, hyper-spectral, and etc.
[e.g. Liu et al. (2009); Li and Li (2010); Gandy et al. (2011); Semerci et al. (2014)]
- ▶ Relational data, recommender system, text mining and etc.
[e.g. Cohen and Collins (2012); Dunlavy et al. (2013); Barak and Moitra (2016)]
- ▶ Latent variable models – topic models, phylogenetic tree, and etc.
[e.g. Anandkumar et al. (2014)]

⋮

- ▶ Algebraic: best low rank approximation may not exist!
[e.g. de Silva and Lim (2008)]
- ▶ Computational: most computations are NP hard!
[e.g. Hillar and Lim (2013)]
- ▶ Probabilistic: different concentration behavior.
[e.g. Y. and Zhang (2016)]

This talk: implications in *tensor completion*

1. Problem

2. Convex Methods

3. Non-convex Methods

Summary

1. Problem

2. Convex Methods

3. Non-convex Methods

Summary

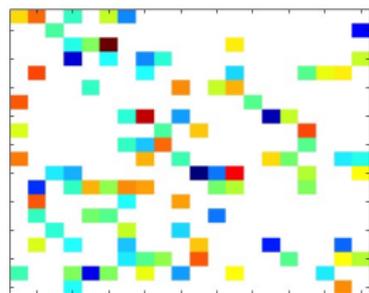
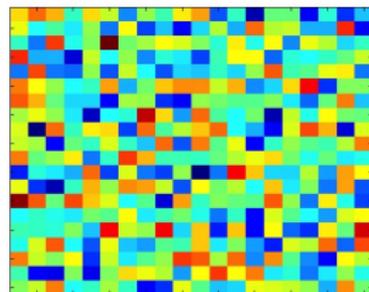
TENSOR COMPLETION

- ▶ Interested in $T \in \mathbb{R}^{d_1 \times \dots \times d_k}$
 - T is of high dimension (d_j s large)
 - T is (approximately) low rank
- ▶ Partial observations:

$$Y_i = T(\omega_i) + \varepsilon_i, \quad i = 1, \dots, n$$

To fix ideas:

- Cubic tensors – $d_1 = \dots = d_k =: d$
 - ω_i s are independently and uniformly sampled
- ▶ Our goal:
- without noise – *exact recovery*
 - with noise – *rates of convergence*



MATRIX COMPLETION ($k = 2$)

- ▶ Also known as *Netflix problem*
- ▶ Incoherence: every observation carries similar amount of information

$$\mu(U) = \frac{d}{r} \max_{1 \leq i \leq d} \|P_U e_i\|^2$$

- ▶ Without measurement error – nuclear norm minimization:

$$\min_M \|M\|_* \quad \text{subject to } M(\omega_i) = T(\omega_i) \quad \forall i$$

Exact recovery if:

$$n \gg rd \cdot \log(d)$$

- ▶ With measurement error – nuclear norm regularization:

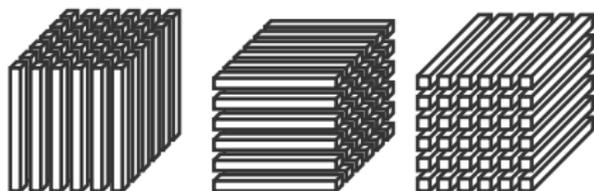
$$\hat{T} = \arg \min_M \left\{ \frac{1}{n} \sum_{i=1}^n [M(\omega_i) - Y_i]^2 + \lambda \|M\|_* \right\}$$

Estimation error:

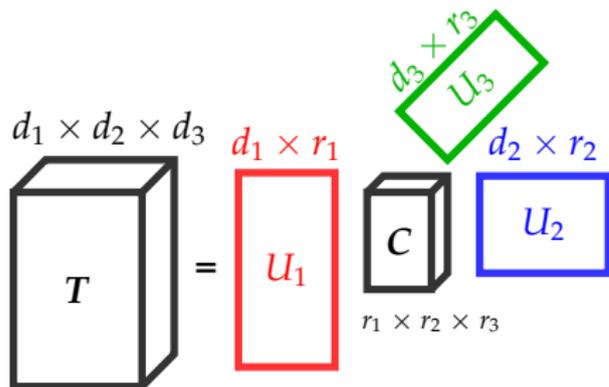
$$\text{MSE}(\hat{T}) := \frac{1}{d^2} \|\hat{T} - T\|_F^2 \lesssim rd \cdot \log(d)/n$$

[e.g., Candes and Recht (2008); Keshavan et al. (2009); Candes and Tao (2010); Gross (2011); Negahban and Wainwright (2011); Recht (2011); Rohde and Tsybakov (2011); Koltchinskii et al. (2012)]

MULTILINEAR RANKS



- ▶ Fibers – vectors obtained by fixing two indices
i.e. mode-1 fibers: $T(:, i_2, \dots, i_k)$
- ▶ Linear space spanned by fibers,
i.e. $\mathcal{L}_1(T) = \text{l.s.}\{T(:, i_2, \dots, i_k) : i_2, \dots, i_k\}$
- ▶ Multilinear ranks



$$r_j = \dim(\mathcal{L}_j(T))$$

- ▶ Tucker decomposition

$$T = (U_1, \dots, U_k) \cdot C$$

$$\mathcal{A}(r) = \{T \in \mathbb{R}^{d \times \dots \times d} : r_j(T) \leq r\} \cong \underbrace{\mathcal{G}(d, r) \times \dots \times \mathcal{G}(d, r)}_{k \text{ times}} \times \mathbb{R}^{r \times \dots \times r}$$

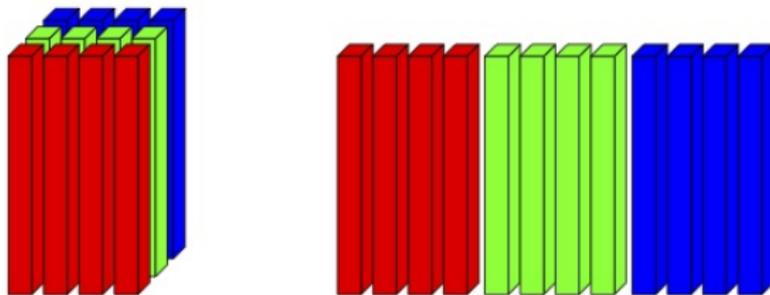
- ▶ Assume same multilinear ranks ($r_1 = \dots = r_3 =: r$) for brevity
- ▶ Any $T \in \mathcal{A}(r)$ – $\text{rank}(T) \in [r, r^{k-1}]$
- ▶ Dimension of $\mathcal{A}(r)$ – $O(r^k + rd)$
- ▶ Gold standard:
 - Exact recovery with $\tilde{O}(r^k + rd)$ noiseless entries
 - Estimation error of the order $\tilde{O}_p((r^k + rd)/n)$
- ▶ *For matrices*: similar bounds are attainable with nuclear norm minimization/regularization

1. Problem

2. Convex Methods

3. Non-convex Methods

Summary



- ▶ Tensor unfolding (matricization)

$$T \in \mathbb{R}^{d \times d \times d} \mapsto \mathcal{M}_1(T) \in \mathbb{R}^{d \times d^2}$$

- ▶ Nuclear norm minimization

$$\min_{A \in \mathbb{R}^{d \times d \times d}} \sum_{j=1}^3 \|\mathcal{M}_j(A)\|_* \quad \text{subject to } A(\omega_i) = T(\omega_i), \forall i$$

- ▶ Sample size requirement

$$|\Omega| \gg rd^2 \text{polylog}(d)$$

[e.g., Signoretto et al. (2010), Tomioka et al. (2010, 2011), Gandy et al. (2011)]

$$\min_{A \in \mathbb{R}^{d \times d \times d}} \|A\|_* \quad \text{subject to } A(\omega_i) = T(\omega_i), \forall i$$

- ▶ Tensor nuclear norm
 - Spectral norm

$$\|A\| = \max_{\|u\|=\|v\|=\|w\|=1} \langle A, u \otimes v \otimes w \rangle$$

- Nuclear norm

$$\|A\|_* = \max_{Y \in \mathbb{R}^{d \times d \times d}: \|Y\| \leq 1} \langle Y, A \rangle.$$

- ▶ Exact recovery with high probability if

$$n \gg (r^{1/2}d^{3/2} + r^2d)\text{polylog}(d)$$

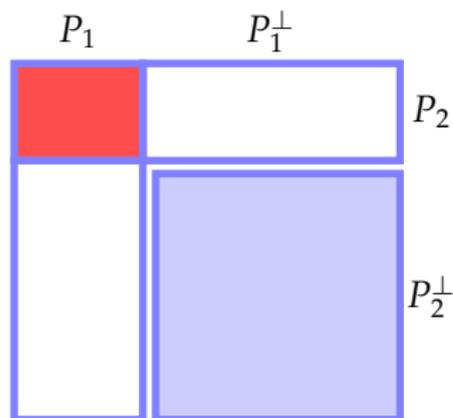
Find a dual certificate:

- ▶ in the subdifferential of $\partial \|\cdot\|_*(T)$
 - characterizing sub-differential of tensor nuclear norm
- ▶ supported on $\{X_1, \dots, X_n\}$
 - concentration inequalities for tensor martingales

MATRIX NUCLEAR NORM

If $M = UDV^\top$, then

$$\left. \begin{array}{l} W = P_U^\perp W P_V^\perp \\ \|W\| \leq 1 \end{array} \right\} \Rightarrow \|X\|_* \geq \|M\|_* + \langle UV^\top + W, X - M \rangle$$



$$\mathcal{P}_M^0 = P_U \otimes P_V$$

$$\mathcal{P}_M^1 = P_U \otimes P_V^\perp$$

$$\mathcal{P}_M^2 = P_U^\perp \otimes P_V$$

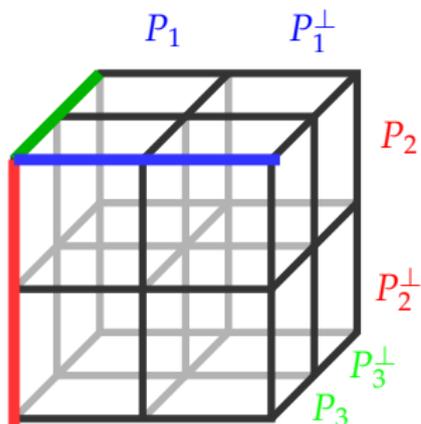
$$\mathcal{P}_M^\perp = P_U^\perp \otimes P_V^\perp$$

[Watson, 1992]

► Projection

$$P_1 \otimes P_2 \otimes P_3[A, B, C] = [P_1A, P_2B, P_3C]$$

► Decomposition of space



$$Q_T^0 = P_T^1 \otimes P_T^2 \otimes P_T^3$$

$$Q_T^1 = P_{T^\perp}^1 \otimes P_T^2 \otimes P_T^3$$

$$Q_T^2 = P_T^1 \otimes P_{T^\perp}^2 \otimes P_T^3$$

$$Q_T^3 = P_T^1 \otimes P_T^2 \otimes P_{T^\perp}^3$$

$$Q_{T^\perp}^0 = P_{T^\perp}^1 \otimes P_{T^\perp}^2 \otimes P_{T^\perp}^3$$

$$Q_{T^\perp}^1 = P_T^1 \otimes P_{T^\perp}^2 \otimes P_{T^\perp}^3$$

$$Q_{T^\perp}^2 = P_{T^\perp}^1 \otimes P_T^2 \otimes P_{T^\perp}^3$$

$$Q_{T^\perp}^3 = P_{T^\perp}^1 \otimes P_{T^\perp}^2 \otimes P_T^3$$

$$\left. \begin{array}{l} W = Q_T^0 W \text{ and } \langle T, W \rangle = \|T\|_* \\ W^\perp = Q_T^\perp W^\perp \text{ and } \|W^\perp\| \leq \frac{1}{2} \end{array} \right\} \implies \|X\|_* \geq \|T\|_* + \langle W + W^\perp, X - T \rangle$$

- ▶ More complex geometry than matrix case
- ▶ For matrices $\|W^\perp\| \leq 1$ is sufficient and necessary
- ▶ For tensor $\|W^\perp\| \leq 1/2$ is only sufficient, *not* necessary:

$$\partial \|\cdot\|_*(T) \supseteq \{W + W^\perp : \|W^\perp\| \leq 1/2\}.$$

* For general k th order tensors, upper bound 1, lower bound $2/k(k-1)$.

CONCENTRATION OF TENSOR MARTINGALE

With probability at least $1 - e^{-t}$,

$$\left\| \frac{d^3}{n} \sum_{i=1}^n A(\omega_i) \mathbf{e}_i - \mathbf{A} \right\| \lesssim \|\mathbf{A}\|_{\max} \cdot \left(\sqrt{\frac{d^4 t}{n}} + \frac{d^3 t}{n} \right) \cdot \text{polylog}(d).$$

- ▶ Contributions from variance and maximum
 - in matrix or vector case – variance dominates
 - for higher order tensors – maximum dominates
- ▶ If \mathbf{A} is *incoherent*, then $\|\mathbf{A}\|_{\max} = O(d^{-3/2})$. Thus

$$n \gg d^{3/2} \text{polylog}(d) \text{ implies } \left\| \frac{d^3}{n} \sum_{i=1}^n A(\omega_i) \mathbf{e}_i - \mathbf{A} \right\| \leq \frac{1}{2}$$

* For general k th order tensors, second term $d^k t/n$.

WHAT HAPPENS WHEN $k > 3$?

- ▶ Following a similar argument leads to sample size requirement

$$n \gtrsim d^{k/2} \text{polylog}(d).$$

- ▶ It can be improved if we incorporate incoherence more explicitly:

$$n \gtrsim (r^{(k-1)/2} d^{3/2} + r^{k-1} d) (\log(d))^2$$

- Depends on the order k only through the rank r
- If $r = O(1)$, then the sample size requirement becomes

$$n \gtrsim d^{3/2} (\log(d))^2$$

INCOHERENT NUCLEAR NORM MINIMIZATION

$$\min_{A \in \mathbb{R}^{d_1 \times \dots \times d_k}} \|X\|_{\star, \delta} \quad \text{subject to } A(\omega_i) = T(\omega_i) \quad i = 1, \dots, n$$

- Incoherent rank-one tensors: $\mathcal{U}(\delta) = \cup_{1 \leq j_1 < j_2 \leq k} \mathcal{U}_{j_1 j_2}(\delta)$ where

$$\mathcal{U}_{j_1 j_2}(\delta) = \{u_1 \otimes \dots \otimes u_k : \|u_j\|_{\ell_2} \leq 1, \forall j; \|u_j\|_{\ell_\infty} \leq \delta_j, \forall j \neq j_1, j_2\}$$

- Incoherent tensor norms:

$$\|X\|_{\circ, \delta} = \sup_{Y \in \mathcal{U}(\delta)} \langle Y, X \rangle, \quad \|X\|_{\star, \delta} = \sup_{\|Y\|_{\circ, \delta} \leq 1} \langle Y, X \rangle$$

- Encourages solution to be incoherent:

- In general, $\|X\|_* \leq \|X\|_{\star, \delta}$
- If $\delta_j \geq \mu_j(X)$, then $\|X\|_* = \|X\|_{\star, \delta}$

1. Problem

2. Convex Methods

3. Non-convex Methods

Summary

POLYNOMIAL-TIME METHODS?

- ▶ Gold standard is $\tilde{O}(r^3 + rd)$
- ▶ Matricization requires $\tilde{O}(rd^2)$
- ▶ Nuclear norm minimization needs $\tilde{O}(r^{1/2}d^{3/2} + r^2d)$
- ▶ *But* tensor nuclear norm is NP hard to compute in the worst case
 - Relaxation – theta norm, sum of squares relaxation
 - Feasible in principle but do not scale well
 - General performance guarantee unclear
- ▶ What about nonconvex methods?
 - Success in some practical examples
 - General performance guarantee unclear

DO THEY WORK?

$$\min_{A \in \mathcal{A}(r)} \frac{1}{n} \sum_{i=1}^n [T(\omega_i) - A(\omega_i)]^2$$

- ▶ Recall that we can write

$$A = (W_1, W_2, W_3) \cdot \mathbf{G}$$

If $n \gg r^3 \text{polylog}(d)$, the above minimization can be equivalently expressed as

$$\min_{W_1, W_2, W_3 \in \mathcal{G}(d, r)} f(W_1, W_2, W_3)$$

- ▶ Smooth optimization techniques to minimize f – practical successes [see, e.g., Vervliet et al., 2014; Kressner et al., 2014]
- ▶ But *why*?

If

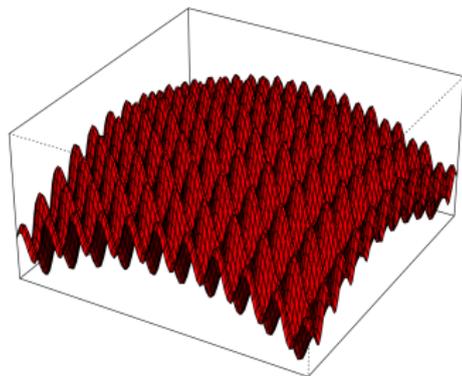
$$n \gg \left(r^{5/2} d^{3/2} + r^4 d \right) \cdot \text{polylog}(d)$$

then f is well-behaved in an *incoherent* neighborhood around truth

WHAT EXACTLY IS THE PROBLEM?

Suppose we want to compute the spectral norm

- ▶ of a random matrix – at most d local optima
- ▶ of a random tensor – $\exp(\Omega(d))$ local optima [Auffinger and Ben Arous (2013)]



- ▶ Polynomial optimization
- ▶ Very smooth but highly nonconvex
- ▶ *If* we can get close to global optimum, ...

Pay a hefty price to get close, pay a little more to get exact!

- ▶ f is minimized at the linear subspace of $\mathcal{L}_1(\mathbf{T})$, $\mathcal{L}_2(\mathbf{T})$ and $\mathcal{L}_3(\mathbf{T})$
- ▶ A first attempt: random initialization – exponentially many tries
- ▶ A second attempt: spectral method
 - $\mathcal{L}_1(\mathbf{T})$ is the column space of $\mathcal{M}_1(\mathbf{T}) \in \mathbb{R}^{d \times d^2}$
 - An unbiased estimate of $\mathcal{M}_1(\mathbf{T})$ is

$$\mathcal{M}_1(\hat{\mathbf{T}}) := \mathcal{M}_1 \left(\frac{d^3}{n} \sum_{i=1}^n T(\omega_i) \mathbf{e}_{\omega_i} \right)$$

- Estimate $\mathcal{L}_1(\mathbf{T})$ by applying SVD to the above estimate
- $n \gg d^2$ to ensure closeness
- ▶ A third attempt: “second order” spectral method

SECOND ORDER SPECTRAL METHOD

- ▶ $\mathcal{L}_1(\mathbf{T})$ is also the eigen-space of $\mathbf{S} := \mathcal{M}_1(\mathbf{T})\mathcal{M}_1(\mathbf{T})^\top \in \mathbb{R}^{d \times d}$
- ▶ $\mathcal{M}_1(\widehat{\mathbf{T}})\mathcal{M}_1(\widehat{\mathbf{T}})^\top$ is a *biased* estimate
- ▶ Unbiased estimate – U-statistic

$$\widehat{\mathbf{S}} := \frac{d^6}{n(n-1)} \sum_{i \neq j} T(\omega_i)T(\omega_j)\mathcal{M}_1(\mathbf{e}_{\omega_i})\mathcal{M}_1(\mathbf{e}_{\omega_j})^\top$$

- ▶ with probability at least $1 - e^{-t}$,

$$\|\widehat{\mathbf{S}} - \mathbf{S}\| \lesssim \|\mathbf{T}\|_{\max}^2 \cdot \left(\frac{d^6 t^2}{n^2} + \frac{d^{9/2} t}{n} \right) \cdot \text{polylog}(d)$$

- For incoherent \mathbf{T} , $\|\widehat{\mathbf{S}} - \mathbf{S}\| = o_p(1)$ if $n \gg d^{3/2} \text{polylog}(d)$
- Sharper concentration around \mathbf{S} than $\mathcal{M}_1(\widehat{\mathbf{T}})$ around $\mathcal{M}_1(\mathbf{T})$
- ▶ Consistent estimates iff

$$n \gg \left(rd^{3/2} + r^2 d \right) \log d.$$

$$Y_i = T(\omega_i) + \varepsilon_i, \quad i = 1, \dots, n$$

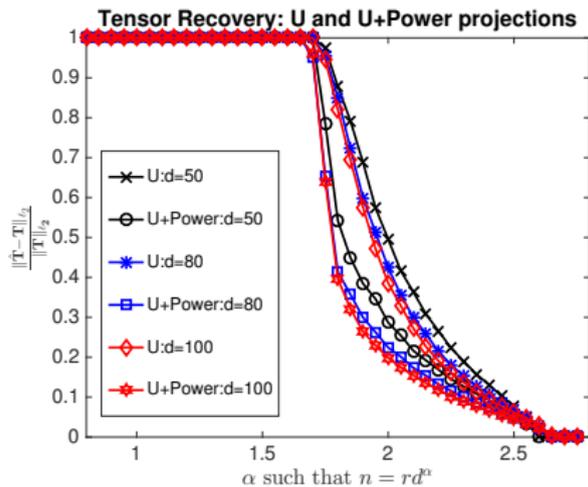
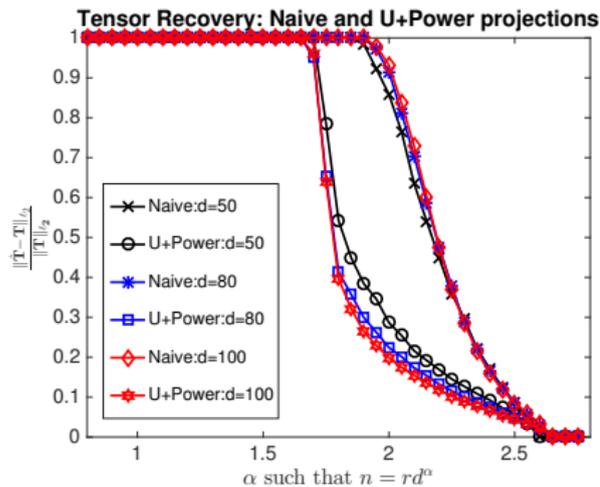
- ▶ More difficult – control the effect of noise
- ▶ But easier – suffices to get close to the target
- ▶ And subtlety – scaling

$$\inf_{\hat{T}} \sup_{T \in \Theta(r_0, \beta_0)} \left(\frac{1}{d^3} \|\hat{T} - T\|_{\ell_p}^p \right)^{1/p} \asymp (\|T\|_{\ell_\infty} + \sigma_\varepsilon) \sqrt{\frac{rd \log(d)}{n}},$$

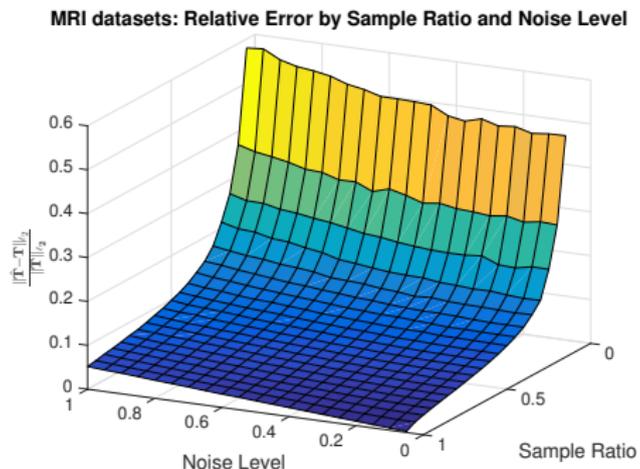
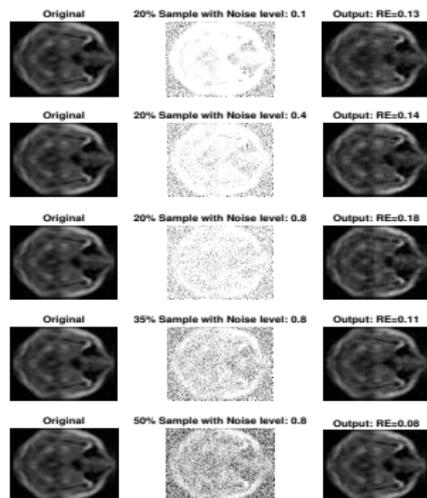
provided that

$$n \gg (rd^{3/2} + r^2d) \log d.$$

EFFECT OF INITIALIZATION AND POWER ITERATION



A DATA* EXAMPLE



* Taken from BrainWeb [Cocosco et al., 1997; $217 \times 181 \times 181$].

1. Problem
2. Convex Methods
3. Non-convex Methods

Summary

CONCLUDING REMARKS

Methods	Tractable?	Sample size requirement
Matricization	Yes	$\tilde{O}(rd^2)$
Nuclear Norm minimization	No	$\tilde{O}(r^{1/2}d^{3/2} + r^2d)$
Nonconvex	Yes	$\tilde{O}(rd^{3/2} + r^2d)$

- ▶ Polynomial-time methods with better dependence on d ? *Possibly no.*
 - “Equivalence” to random k -SAT problem [Barak and Moitra (2016)]
- ▶ Polynomial-time methods with better dependence on r ? *I don't know.*
- ▶ Is it worth the while? *Definitely yes!*
 - Regression [Chen, Raskutti and Y. (2015, 2016)]
 - PCA [Liu, Y. and Zhao (2017)]

⋮