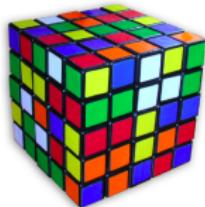


Beyond matrices: nonparametric tensor completion via sign series

Miaoyan Wang

Department of Statistics, UW-Madison

Joint work with Chanwoo Lee (3rd-year PhD student)



Research in my group

Statistical machine learning:

- ▶ Structured tensor decomposition, latent factor models

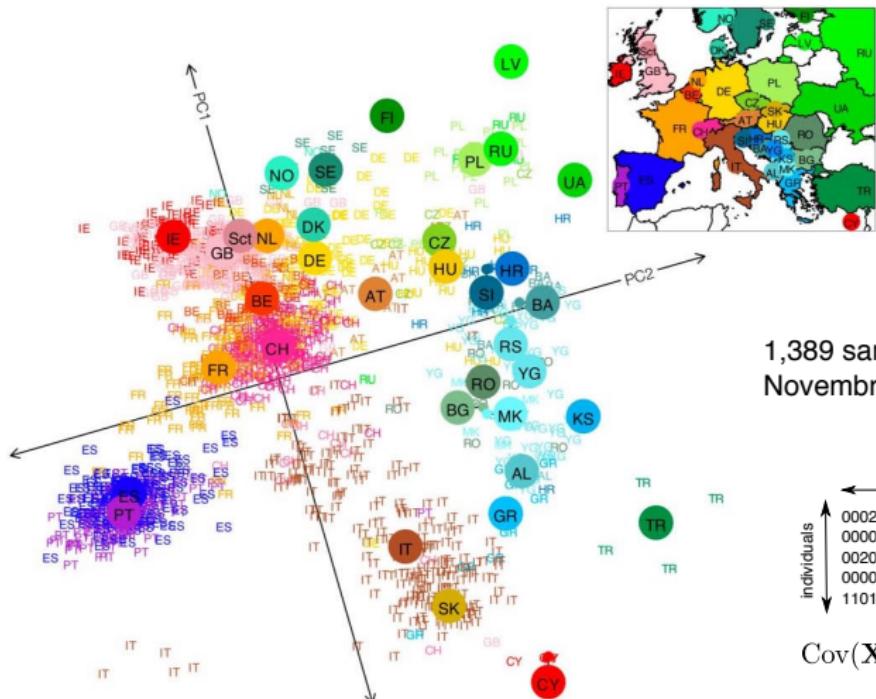
Genetics and genomics:

- ▶ gene expression analyses, genetic association studies

Foundations of data science:

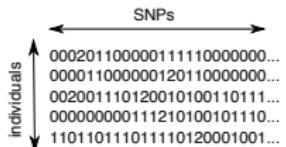
- ▶ Statistical-computational tradeoff for big data analytics

A successful story: PCA of Europeans

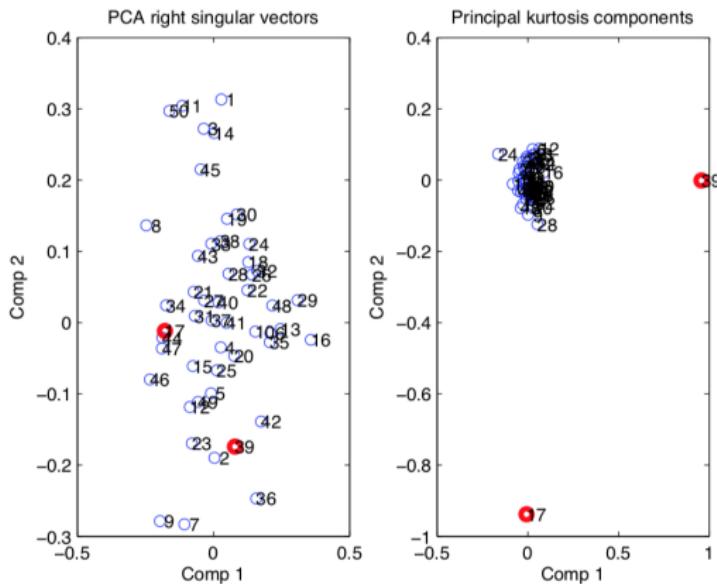


1,389 samples, ~ 200k SNPs
Novembre et al. (2008)

$$\text{Cov}(\mathbf{X}) = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^T$$



Matrix methods are powerful, however...



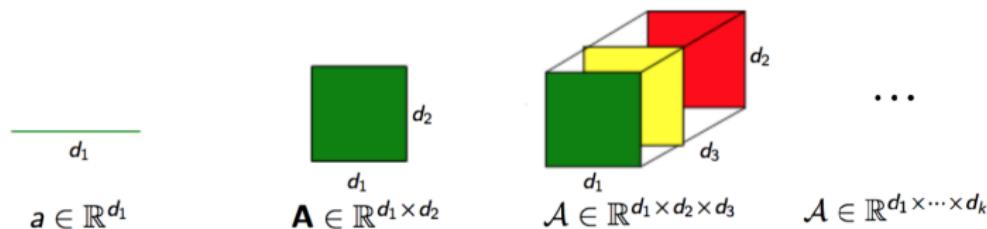
All Gaussian except points 17 and 39.

left: matrix PCA; right: principal components of kurtosis.

Figure credit: Jason Morton and Lek-Heng Lim (2009/2015).

What is a tensor?

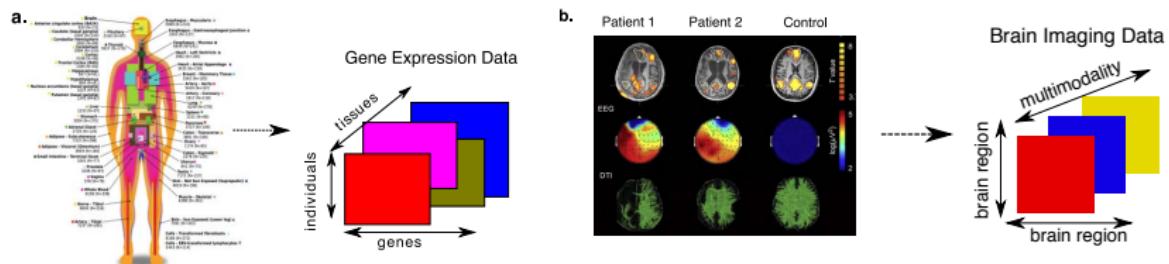
- ▶ Tensors are generalizations of vectors and matrices:



- ▶ An order- k tensor $\mathcal{A} = [[a_{i_1 \dots i_k}]] \in \mathbb{R}^{d_1 \times \dots \times d_k}$ is a hypermatrix with dimensions (d_1, \dots, d_k) and entries $a_{i_1 \dots i_k} \in \mathbb{R}$.
- ▶ This talk will focus on tensor of order 3 or greater, also known as **higher-order tensors**.

Tensors in genomics

- ▶ Many datasets come naturally in a multiway form.
- ▶ Multi-tissue, multi-individual gene expression measures could be organized as an order-3 tensor $\mathcal{A} = [[a_{git}]] \in \mathbb{R}^{n_G \times n_I \times n_T}$.



Tensors in statistical modeling

“Tensors are the new matrices” that tie together a wide range of areas:

- ▶ Longitudinal social network data $\{\mathbf{Y}_t : t = 1, \dots, n\}$
- ▶ Spatio-temporal transcriptome data
- ▶ Joint probability table of a set of variables $\mathbb{P}(X_1, X_2, X_3)$
- ▶ Higher-order moments in topic models
- ▶ Markov models for the phylogenetic tree $K_{1,3}$

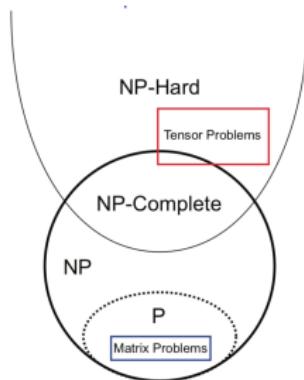
M. Yuan et al 2017, P. Hoff 2015, Montanari-Richard 2014

Anandkumar et al 2014, Mossel et al 2004, P. McCullagh 1987

Talk outline

Prohibitive Computational Complexity

Most higher-order tensor problems are NP-hard [Hillar & Lim, 2013].



Fortunately, the tensors sought in statistical and machine learning applications are often **specially structured**:

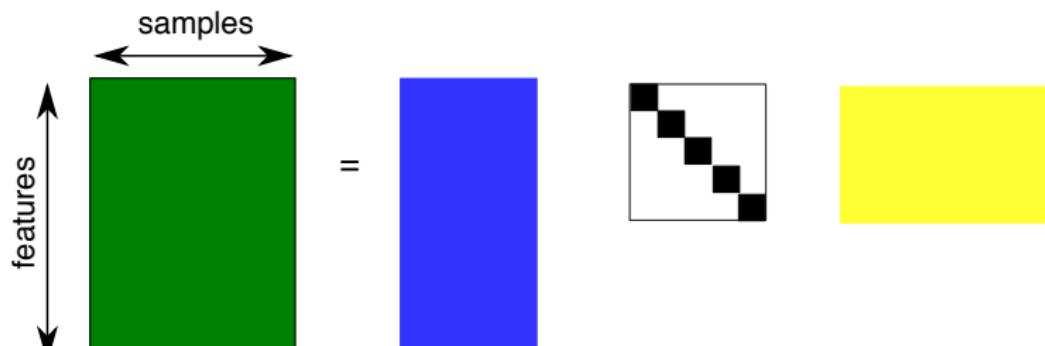
- ▶ Low-rankness
- ▶ Sparsity
- ▶ Non-negativity
- ▶ ...

This talk is based on

Beyond the Signs: Nonparametric Tensor Completion via Sign Series. Lee and W.

- ▶ arXiv:
- ▶ Software: Coming up soon

Review: Matrix SVD for biclustering



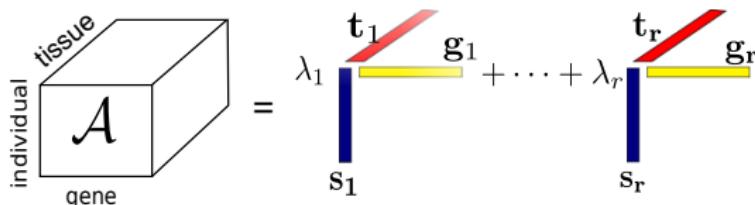
$$\begin{aligned}\mathbf{X} &= \mathbf{U} \quad \Lambda \quad \mathbf{V}^T \\ &= \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{v}_i^T\end{aligned}$$

- ▶ Columns of \mathbf{U} describe patterns across samples
- ▶ Columns of \mathbf{V}^T describe patterns across genes

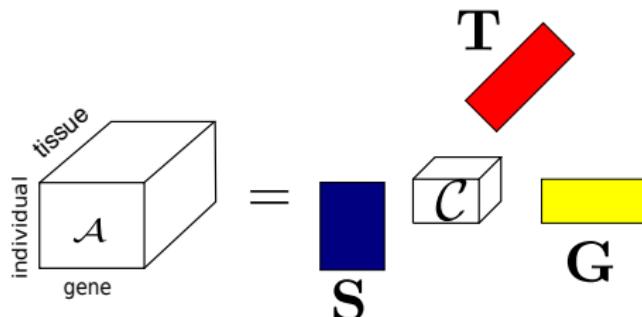
Y Kluger et al, Genome Research (2003). 13(4): 703-71
Data Science Specialization (COURSERA) by Brian Caffo and Jeff Leek

Various notions of low-rankness

- Canonical polyadic (CP) low-rankness: $\mathcal{A} = \sum_{r=1}^R \lambda_r \mathbf{s}_r \otimes \mathbf{g}_r \otimes \mathbf{t}_r.$

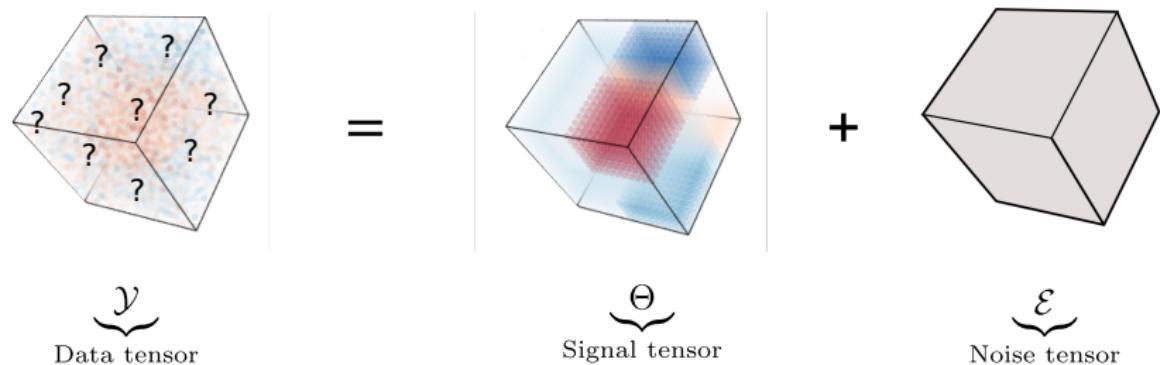


- Tucker low-rankness: $\mathcal{A} = \mathcal{C} \times_1 \mathbf{S} \times_2 \mathbf{G} \times_3 \mathbf{T}.$



- Others: tensor train [Oseledet '11], tensor block model [W. & Zeng '19; Han, Luo, W. et al '20], etc.

Setup: the signal plus noise model



We focus on the two problems:

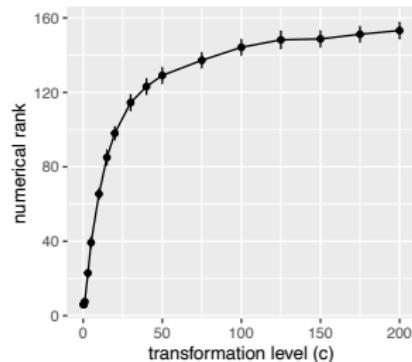
1. Nonparametric tensor estimation: How to estimate the signal tensor Θ under **a wide range of structures?**
2. Tensor completion: How many **observed tensor entries** do we need in order for consistent recovery?

Inadequacies of low-rank models

- ▶ Low rank models (Jain and Oh, 2014; Montanari and Sun, 2018).
- ▶ Two limitations of classical low-rank models:
 1. Sensitivity to order-preserving transformation.
 2. Inadequacy for special structures.

Inadequacies of low-rank models

- ▶ Tensor rank is sensible to order-preserving transformation.

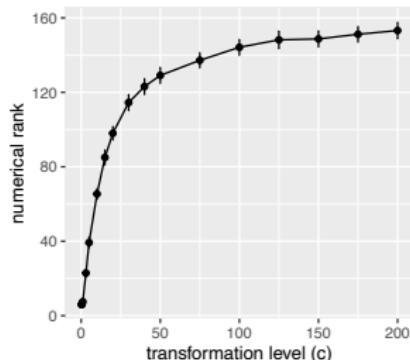


$$\Theta = \frac{1}{1 + \exp(-c(\mathcal{Z}))}, \quad \text{where}$$
$$\mathcal{Z} = \mathbf{a}^{\otimes 3} + \mathbf{b}^{\otimes 3} + \mathbf{c}^{\otimes 3}$$

⇒ Θ is high-rank but \mathcal{Z} is low-rank.

Inadequacies of low-rank models

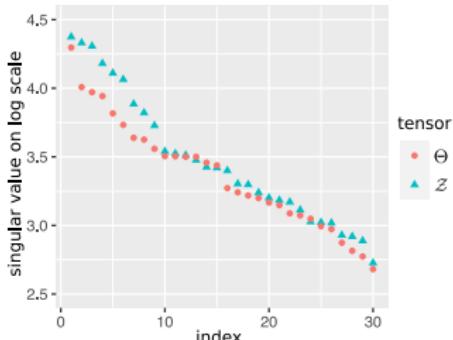
- ▶ Tensor rank is sensible to order-preserving transformation.



$$\Theta = \frac{1}{1 + \exp(-c(\mathcal{Z}))}, \quad \text{where}$$
$$\mathcal{Z} = \mathbf{a}^{\otimes 3} + \mathbf{b}^{\otimes 3} + \mathbf{c}^{\otimes 3}$$

⇒ Θ is high-rank but \mathcal{Z} is low-rank.

- ▶ Low-rank model fails to address several important structures.



$$\Theta = \log(1 + \mathcal{Z}), \quad \text{where}$$
$$\mathcal{Z} = [\mathcal{Z}(i, j, k)] = \frac{1}{d} \max(i, j, k).$$

⇒ Both Θ and \mathcal{Z} are full rank.

The matrix analogy of Θ was studied in the context of graphon analysis by Chan and Airola (2014).

Sign rank

- ▶ Key ideas: we use a local (nonparametric) notion of “low-rankness” that allows a broader family of signal tensors.
- ▶ Two tensors are sign equivalent, denoted $\Theta \simeq \Theta'$, if $\text{sgn}(\Theta) = \text{sgn}(\Theta')$.
- ▶ Define the sign rank by

$$\text{srank}(\Theta) = \min\{\text{rank}(\Theta') : \Theta' \simeq \Theta, \Theta' \in \mathbb{R}^{d_1 \times \dots \times d_K}\}.$$

$$\Theta = \begin{array}{c} \text{A 3D tensor with } 3 \text{ layers, each } 2 \times 2 \\ \text{The first layer has a red square at position } (1,1) \end{array}, \quad \text{sgn}(\Theta) = \begin{array}{c} \text{A 2D matrix with } 2 \text{ columns and } 2 \text{ rows} \\ \text{The top-left entry is dark red, all other entries are blue} \end{array} \implies \begin{array}{l} \text{rank}(\Theta) = d \\ \text{srank}(\Theta) = 2 \end{array}$$

- ▶ For any strictly monotonic function $g: \mathbb{R} \rightarrow \mathbb{R}$ with $g(0) = 0$,

$$\text{srank}(\Theta) \leq \text{rank}(g(\Theta)).$$

Sign representable tensors

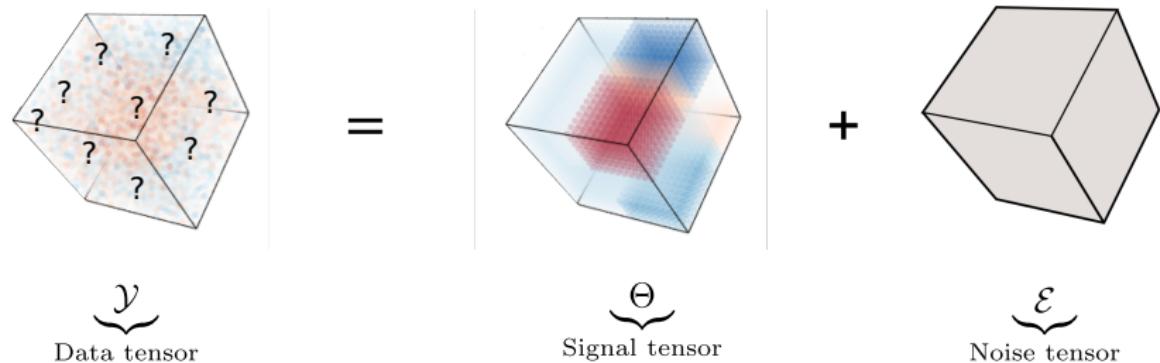
Sign representable tensors

A tensor Θ is called **r-sign representable** if the tensor $(\Theta - \pi)$ has sign rank bounded by r for all $\pi \in [-1, 1]$.

- ▶ Most existing structured tensors belong to sign representable family:
 - ▶ **Low-rank** CP tensors, Tucker tensors, stochastic tensor block models.
 - ▶ **High-rank** tensors from GLM, single index models.
 - ▶ Earlier example $\Theta(i_1, \dots, i_K) = \log(1 + \max(i_1, \dots, i_K))$ is 2-sign representable \Rightarrow conclusion extends to general max/min hypergraphon models.
- ▶ We propose the signal tensor family

$$\Theta \in \mathcal{P}_{\text{sgn}}(r) := \{\Theta : \text{srank}(\Theta - \pi) \leq r \text{ for all } \pi \in [-1, 1]\}.$$

Recall the goal

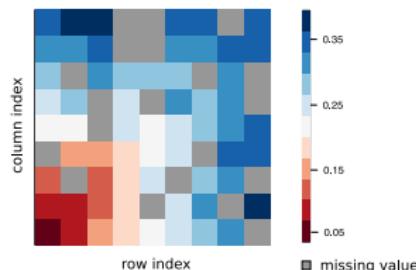


We focus on the two problems:

1. Nonparametric tensor estimation: How to estimate the signal tensor Θ under **both low- and high-rank** models?
2. Tensor completion: How many **observed tensor entries** do we need in order for consistent recovery?

Denoising and completion of high-rank tensors

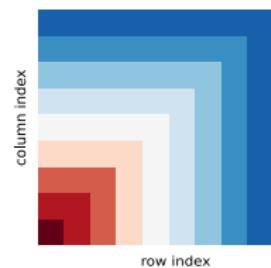
a



noisy and incomplete observation

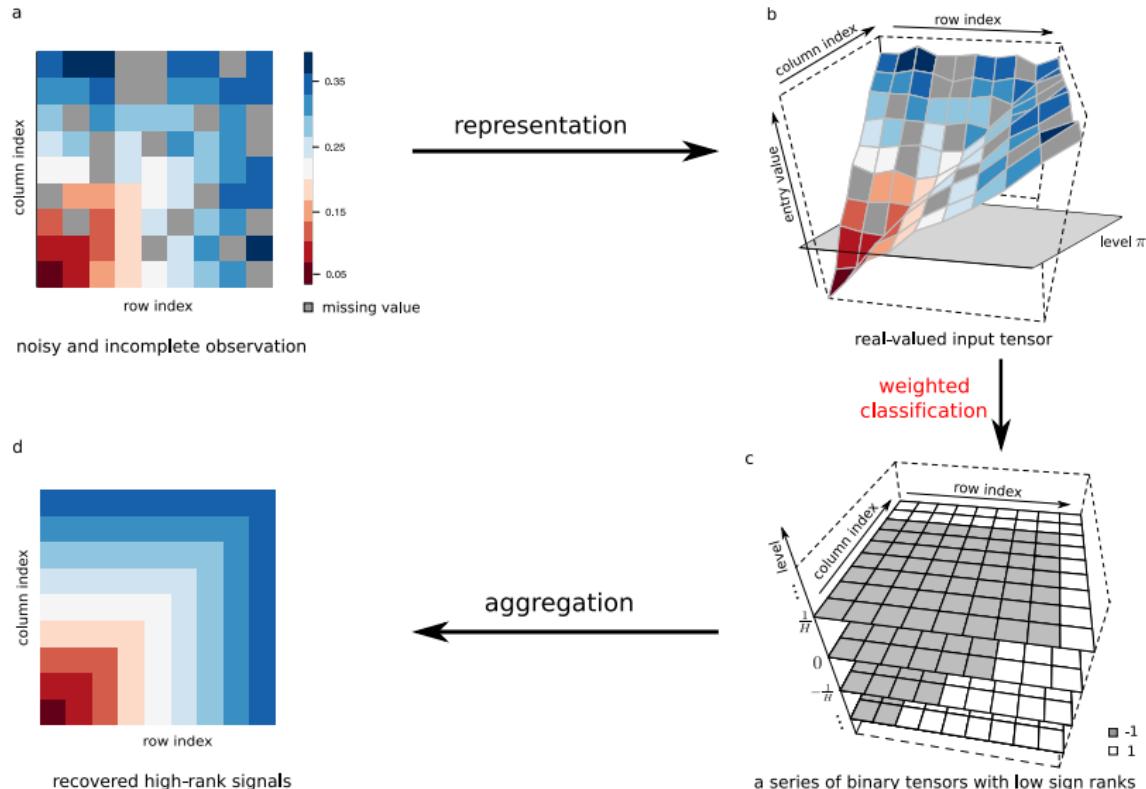


d



recovered high-rank signals

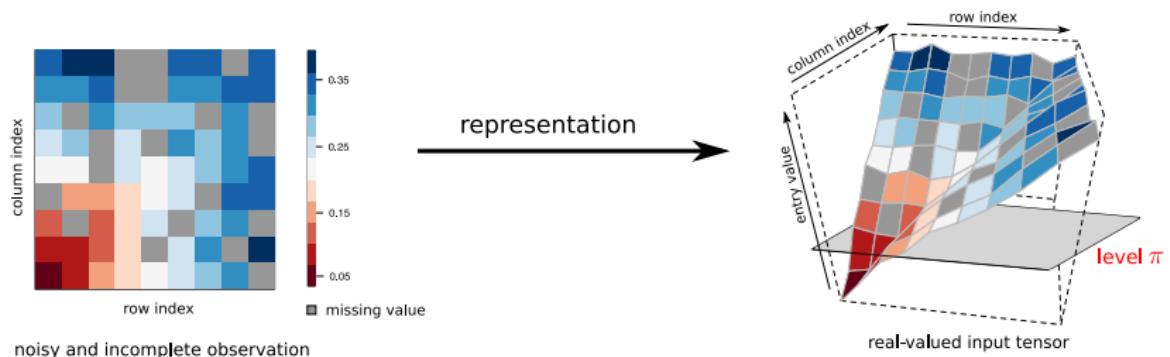
Sign signal helps!



Dichotomized representation

- ▶ We observe **an incomplete noisy tensor** $\mathcal{Y}_\Omega \in [-1, 1]^{d_1 \times \dots \times d_K}$ with observed index set $\Omega \in [d_1] \times \dots \times [d_K]$ under uniform sampling scheme.
- ▶ We dichotomize the data into **a series of sign tensors**:

$$\{\text{sgn}(\mathcal{Y}_\Omega - \pi)\}_{\pi \in \mathcal{H}}, \quad \text{where } \mathcal{H} = \left\{-1, \dots, -\frac{1}{H}, 0, \frac{1}{H}, \dots, 1\right\}.$$



Sign estimation via weighted classification

- ▶ We estimate $\text{sgn}(\Theta - \pi)$ through $\text{sgn}(\mathcal{Y}_\Omega - \pi)$ via weighted classification.
- ▶ Objective function of weighted classification is

$$L(\mathcal{Z}, \mathcal{Y}_\Omega - \pi) = \frac{1}{|\Omega|} \sum_{\omega \in \Omega} \underbrace{|\mathcal{Y}(\omega) - \pi|}_{\text{weight}} \times \underbrace{|\text{sgn}(\mathcal{Z}(\omega)) - \text{sgn}(\mathcal{Y}(\omega) - \pi)|}_{\text{classification loss}}$$



Identification for weighted classification

α -smoothness of signal tensor

For fixed π , Θ is α -smooth if there exist $\alpha = \alpha(\pi) > 0, c = c(\pi) > 0$, s.t.

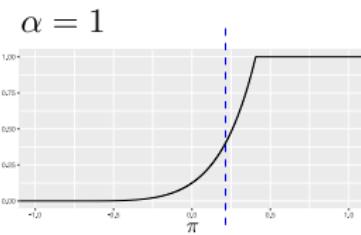
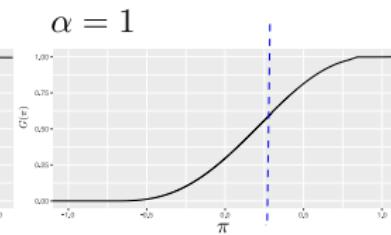
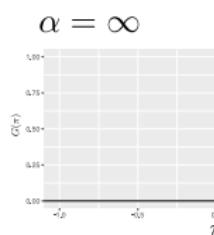
$$\sup_{0 \leq t < \rho(\pi, \mathcal{N})} \frac{\mathbb{P}_{\omega \sim \Pi}[|\Theta(\omega) - \pi| \leq t]}{t^\alpha} \leq c,$$

where $\rho(\pi, \mathcal{N}) = \min_{\pi' \in \mathcal{N}} |\pi - \pi'|$ and $\mathcal{N} = \{\pi : \mathbb{P}(\Theta(\omega) = \pi) \neq 0\}$. If α and c are global constants for all* π 's, we call Θ is α -globally smooth.

* except for a finite number of π 's.

Intuition: sign recovery is harder at levels where point mass concentrates.

Rate depends on the behavior of CDF function $G(\pi) = \mathbb{P}_{\omega \sim \Pi}[\Theta(\omega) \leq \pi]$.



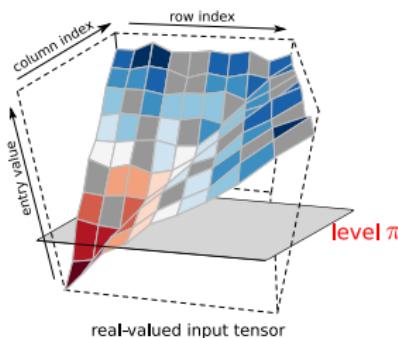
Identification for weighted classification

- If Θ is α -smooth ($\alpha > 0$), we have **an unique optimizer** such that

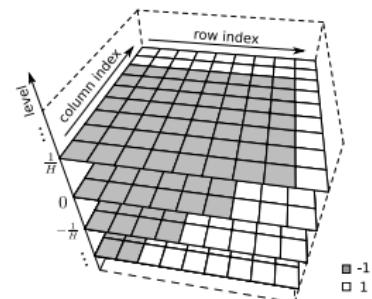
$$\text{sgn}(\Theta - \pi) = \arg \min_{\mathcal{Z}: \text{rank}(\mathcal{Z}) \leq r} \mathbb{E}_{\mathcal{Y}_\Omega} L(\mathcal{Z}, \mathcal{Y}_\Omega - \pi).$$

- We obtain a series of optimizers $\{\hat{\mathcal{Z}}_\pi\}_{\pi \in \mathcal{H}}$ as

$$\hat{\mathcal{Z}}_\pi = \arg \min_{\mathcal{Z}: \text{rank}(\mathcal{Z}) \leq r} L(\mathcal{Z}, \mathcal{Y}_\Omega - \pi).$$



weighted classification



a series of binary tensors with low sign ranks

* Uniqueness up to sign equivalence, meaning the optimizer $\Theta_{\text{opt}} \simeq \text{sgn}(\Theta - \pi)$.

Sign tensor estimation error

- ▶ For two tensors Θ_1, Θ_2 , define $\text{MAE}(\Theta_1, \Theta_2) = \mathbb{E}_{\omega \in \Pi} |\Theta_1(\omega) - \Theta_2(\omega)|$.

Sign tensor estimation for fixed π (L. and Wang, 2021)

Suppose $\Theta \in \mathcal{P}_{\text{sgn}}(r)$ and $\Theta(\omega)$ is α -smooth for fixed π . Let $d_{\max} = \max_{k \in [K]} d_k$. Then, with very high probability over \mathcal{Y}_Ω ,

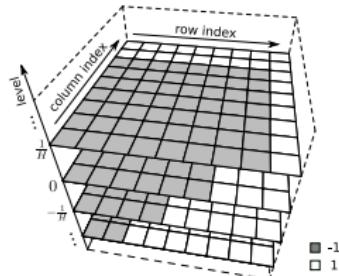
$$\text{MAE}(\text{sgn} \hat{\mathcal{Z}}_\pi, \text{sgn}(\Theta - \pi)) \lesssim \left(\frac{d_{\max} r}{|\Omega|} \right)^{\frac{\alpha}{\alpha+2}}.$$

- ▶ The sign estimation error shows a polynomial decay with $|\Omega|$.
- ▶ Best rate attains for stochastic tensor block models ($\alpha = \infty$).

From sign to signal estimation

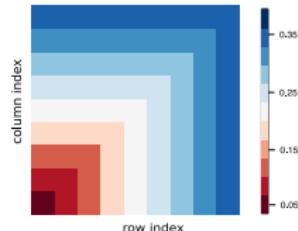
- ▶ Aggregation of sign tensors from weighted classification yields the possibly high-rank signal tensor estimate:

$$\hat{\Theta} = \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{sgn} \hat{\mathcal{Z}}_\pi.$$



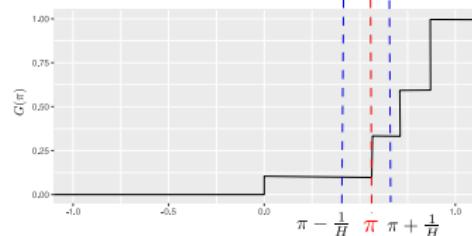
a series of binary tensors with low sign ranks

aggregation →



recovered high-rank signals

- ▶ Signal tensor estimation is robust to a few off-target classifications.



Tensor estimation error

Tensor estimation error (L. and Wang 2021)

Suppose $\Theta \in \mathcal{P}_{\text{sgn}}(r)$ and $\Theta(\omega)$ is α -globally smooth. Then, with very high probability over \mathcal{Y}_Ω ,

$$\text{MAE}(\hat{\Theta}, \Theta) \lesssim \underbrace{\left(\frac{d_{\max} r}{|\Omega|} \right)^{\frac{\alpha}{\alpha+2}}}_{\text{error inherited from sign estimation}} + \underbrace{\frac{1}{H}}_{\text{Bias}} + \underbrace{\frac{H d_{\max} r}{|\Omega|}}_{\text{Variance}}.$$

In particular, setting $H \asymp \left(\frac{|\Omega|}{d_{\max} r} \right)^{1/2}$ yields the error bound

$$\text{MAE}(\hat{\Theta}, \Theta) \lesssim \left(\frac{d_{\max} r}{|\Omega|} \right)^{\frac{\alpha}{\alpha+2} \vee \frac{1}{2}}.$$

- ▶ Sample requirement for tensor completion:

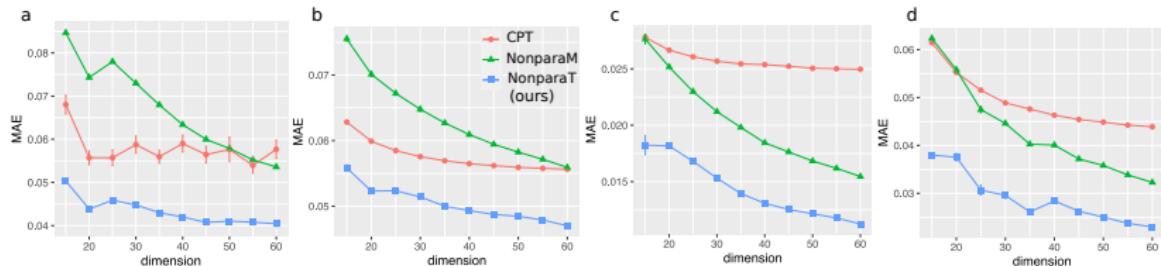
$$\text{MAE}(\hat{\Theta}, \Theta) \rightarrow 0, \text{ as } \frac{|\Omega|}{d_{\max} r} \rightarrow \infty.$$

Comparison of estimation error versus dimension

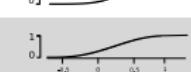
- We simulate signal tensors under a wide range of complexity.

Simulation	Signal Tensor Θ	Rank	Sign Rank	Global α	CDF	Noise
1	$\mathcal{C} \times M_1 \times M_2 \times M_3$	3^3	$\leq 3^3$	∞		Uniform $[-0.3, 0.3]$
2	$ \mathbf{a} \otimes \mathbf{1} \otimes \mathbf{1} - \mathbf{1} \otimes \mathbf{a} \otimes \mathbf{1} $	d	≤ 3	1		Normal $\mathcal{N}(0, 0.15)$
3	$\log(0.5 + Z_{\max})$	$\geq d$	2	1		Uniform $[-0.1, 0.1]$
4	$2.5 - \exp(Z_{\min}^{1/3})$	$\geq d$	2	1		Normal $\mathcal{N}(0, 0.15)$

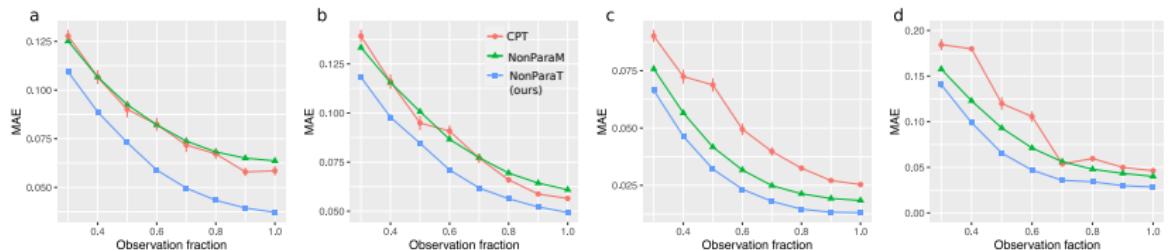
- Our method (**NonparaT**) achieves the best performance, whereas the second best method is low-rank CP tensor (**CPT**) for models 1-2, and matrix version of our method (**NonParaM**) for models 3-4.



Estimation error versus observation fraction

Simulation	Signal Tensor Θ	Rank	Sign Rank	Global α	CDF	Noise
1	$\mathcal{C} \times M_1 \times M_2 \times M_3$	3^3	$\leq 3^3$	∞		Uniform $[-0.3, 0.3]$
2	$ \mathbf{a} \otimes \mathbf{1} \otimes \mathbf{1} - \mathbf{1} \otimes \mathbf{a} \otimes \mathbf{1} $	d	≤ 3	1		Normal $\mathcal{N}(0, 0.15)$
3	$\log(0.5 + Z_{\max})$	$\geq d$	2	1		Uniform $[-0.1, 0.1]$
4	$2.5 - \exp(Z_{\min}^{1/3})$	$\geq d$	2	1		Normal $\mathcal{N}(0, 0.15)$

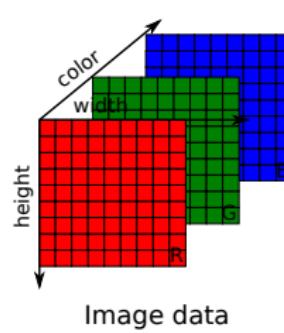
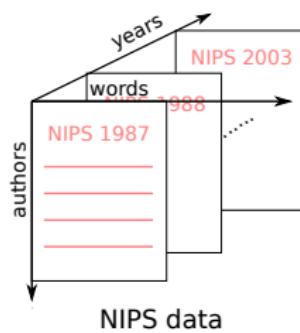
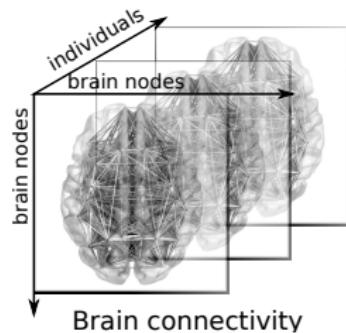
- Our method (NonparaT) achieves the best performance in completion.



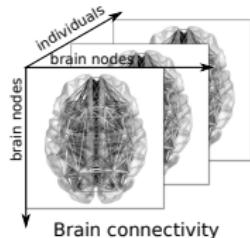
Data application

We apply our method to three datasets:

- ▶ The human brain connectivity data consists of 68 brain regions for 114 individuals along with their IQ scores.
- ▶ The NIPS dataset consists of word occurrence counts in papers published from 1987 to 2003.
- ▶ The 3-channel image data is from licensed google images.

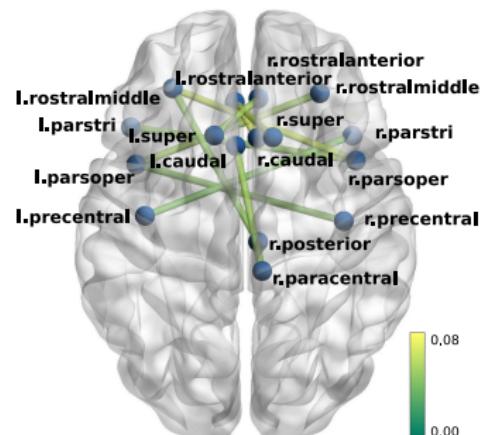


Data application: Brain connectivity

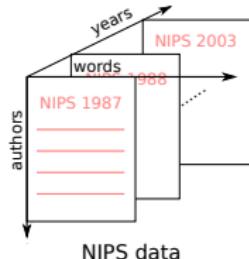


- ▶ The MRN-114 human brain connectivity data consists of 68 brain regions for 114 individuals along with their IQ scores (Wang et al., 2017).
- ▶ Data tensor $\mathcal{Y} \in \{0, 1\}^{68 \times 68 \times 114}$.

- ▶ We examine the estimated signal tensor $\hat{\Theta}$.
- ▶ Top 10 brain edges based on regression analysis show inter-hemisphere connections.

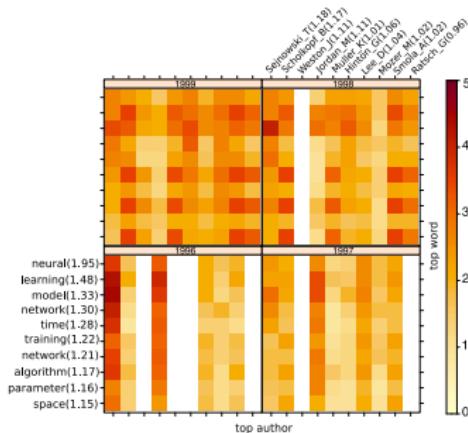


Data application: NIPS



- ▶ The NIPS dataset consists of word occurrence counts in papers published from 1987 to 2003 (Globerson et al., 2007).
- ▶ Data tensor $\mathcal{Y} \in \mathbb{R}^{100 \times 200 \times 17}$.

- ▶ We examine the estimated signal tensor $\hat{\Theta}$.
- ▶ Most frequent words is consistent with the active topics.
- ▶ Strong heterogeneity among word occurrences across authors and years.



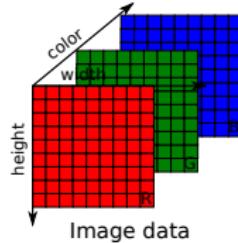
Data application: Brain connectivity + NIPS

- ▶ Our method has lower test error compared to low-rank tensor method.

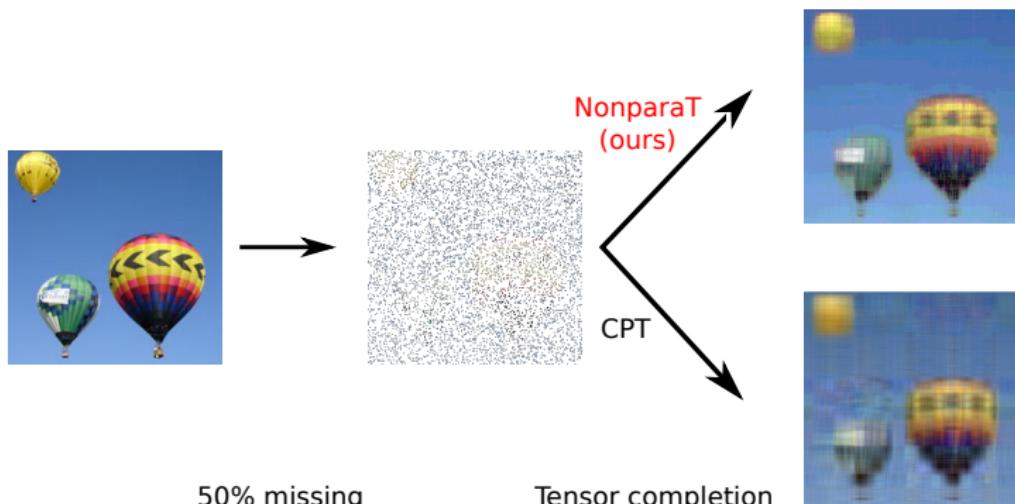
MRN-114 brain connectivity dataset					
Method	$r = 3$	$r = 6$	$r = 9$	$r = 12$	$r = 15$
NonparaT (Ours)	0.18(0.001)	0.14(0.001)	0.12(0.001)	0.12(0.001)	0.11(0.001)
Low-rank CPT	0.26(0.006)	0.23(0.006)	0.22(0.004)	0.21(0.006)	0.20(0.008)
NIPS word occurrence dataset					
Method	$r = 3$	$r = 6$	$r = 9$	$r = 12$	$r = 15$
NonparaT (Ours)	0.18(0.002)	0.16(0.002)	0.15(0.001)	0.14(0.001)	0.13(0.001)
Low-rank CPT	0.22(0.004)	0.20(0.007)	0.19(0.007)	0.17(0.007)	0.17(0.007)
Naive imputation (Baseline)	0.32(.001)				

Table: MAE comparison in the brain data and NIPS data on cross-validation (5 repetitions 5 folds). Standard errors are reported in parenthesis.

Data application: Image



- ▶ The original data is from licensed google image file.
- ▶ $\mathcal{Y} \in [0, 1]^{217 \times 217 \times 3}$.
- ▶ We sample 50% entries in the original image tensor and check completion performance.



Summary

Tensor analysis provides a rich source of

- ▶ fundamental problems in data science.
- ▶ new tools for long-standing questions.
- ▶ potentials for new applications.

Our general strategy is to carve out a broad range of **specially-structured tensors** that are useful in practice, and to develop efficient statistical methods for analyzing these high-dimensional tensor data.

References:

- ▶ Beyond the sign: nonparametric tensor completion from sign series.

Acknowledgment: NSF DMS 1915978 and Grant from Wisconsin Alumni Research Foundation.

Appendix

Algorithm 1 Nonparametric tensor completion

Input: Noisy and incomplete data tensor \mathcal{Y}_Ω , rank r , resolution parameter H .

```
1: for  $\pi \in \mathcal{H} = \{-1, \dots, -\frac{1}{H}, 0, \frac{1}{H}, \dots, 1\}$  do
2:   Random initialization of tensor factors  $\mathbf{A}_k = [\mathbf{a}_1^{(k)}, \dots, \mathbf{a}_r^{(k)}] \in \mathbb{R}^{d_k \times r}$  for all  $k \in [K]$ .
3:   while not convergence do
4:     for  $k = 1, \dots, K$  do
5:       Update  $\mathbf{A}_k$  while holding others fixed:
6:        $\mathbf{A}_k \leftarrow \arg \min_{\mathbf{A}_k \in \mathbb{R}^{d_k \times r}} \sum_{\omega \in \Omega} |\mathcal{Y}(\omega) - \pi| F(\mathcal{Z}(\omega) \text{sgn}(\mathcal{Y}(\omega) - \pi)),$ 
7:       where  $F(\cdot)$  is the large-margin loss, and  $\mathcal{Z} = \sum_{s \in [r]} \mathbf{a}_s^{(1)} \otimes \dots \otimes \mathbf{a}_s^{(K)}$  is a rank- $r$  tensor.
8:     end for
9:   end while
10:  Return  $\mathcal{Z}_\pi \leftarrow \sum_{s \in [r]} \mathbf{a}_s^{(1)} \otimes \dots \otimes \mathbf{a}_s^{(K)}$ .
11: end for
```

Output: Estimated signal tensor $\hat{\Theta} = \frac{1}{2H+1} \sum_{\pi \in \mathcal{H}} \text{sgn}(\mathcal{Z}_\pi)$.

References I

- De Lathauwer, L., De Moor, B., and Vandewalle, J. (2000). A multilinear singular value decomposition. *SIAM Journal on Matrix Analysis and Applications*, 21(4):1253–1278.
- Globerson, A., Chechik, G., Pereira, F., and Tishby, N. (2007). Euclidean embedding of co-occurrence data. *Journal of Machine Learning Research*, 8:2265–2295.
- Hitchcock, F. L. (1927). The expression of a tensor or a polyadic as a sum of products. *Journal of Mathematics and Physics*, 6(1-4):164–189.
- Jain, P. and Oh, S. (2014). Provable tensor factorization with missing data. In *Advances in Neural Information Processing Systems*, volume 27, pages 1431–1439.

References II

- Montanari, A. and Sun, N. (2018). Spectral algorithms for tensor completion. *Communications on Pure and Applied Mathematics*, 71(11):2381–2425.
- Wang, L., Durante, D., Jung, R. E., and Dunson, D. B. (2017). Bayesian network–response regression. *Bioinformatics*, 33(12):1859–1866.