

**Title:** The Good, the Bad, the Pragmatic: Tensor Methods for Network Learning

**Author:** Dr. Miaoyan Wang, Assistant Professor in Department of Statistics, Affiliated with Institute for the Foundations of Data Science (IFDS)

**WISPER Record Number:** MSN251318

**Abstract:** The prevailing theme in the proposal is to develop powerful tensor methods for high-dimensional multi-layer network analysis. Rapid developments in modern technologies have made large-scale network datasets readily available. Modern networks are not only large in size, but they also have intricate structure. It is therefore of great importance to find a low-dimensional representation to better understand the key structure buried in noisy observations.

Higher-order tensors provide effective representation of multi-layer networks using multi-way structure. The PI will develop a framework — of tensor models, efficient algorithms, and softwares — to analyze multi-layer networks. Previous literature has advocated unfolding the tensor into a matrix and applying classical methods developed for matrices. Despite the popularity of such techniques, tensor method provides more powerful tools to capture complex structures in data that lower-order methods fail to exploit. The research goal goes beyond the traditional multivariate analysis; we aim to characterize probabilistic distributions over multi-layer edge connections, while taking into accounting the higher-order structures such as transitivity, balance, and community. This will allow researchers to examine complex interactions among entities in a context-specific manner, thereby providing solutions to questions that were previously impossible. The software packages resulting from this proposal, will be released freely, as well as related visualization tools for network analyses.

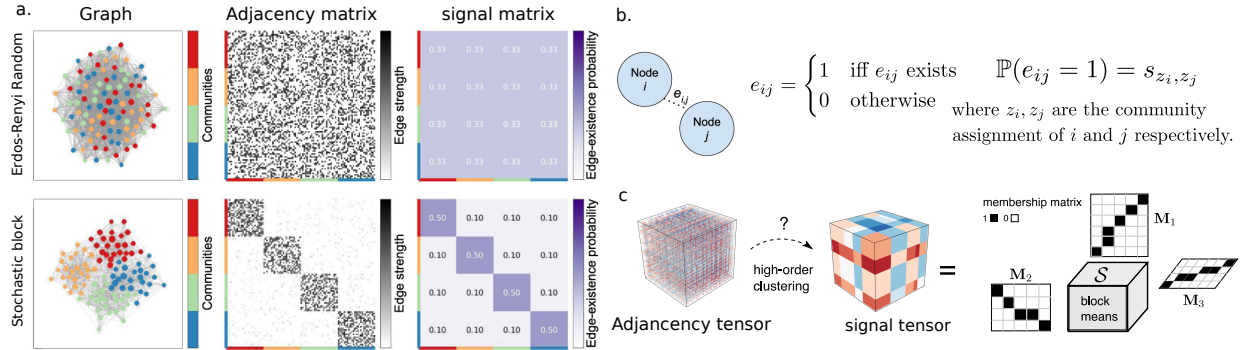
# The Good, the Bad, the Pragmatic: Tensor Methods for Network Learning

My research project is broadly driven by questions related to understanding hidden salient structures of complex network datasets. The questions and techniques that we develop span across the fields of information theory, machine learning, and quantitative social sciences.

**Motivations.** A central theme in modern data analysis is to find a low-dimensional representation to better understand, compress, and convey the key phenomena buried in noisy observations. Modern datasets are not only large in size, but they also have intricate structure. A typical example is in the form of **network** [2, 10, 19], which is quantitative representation of interactions between entities in complex systems. As real-world networks are huge in nature, dimension deduction is crucial for pattern detection and subsequent specialized tasks.

In network studies, researcher are interested in interpretable low-dimensional structure within the high-dimensional relational data. The question goes beyond the traditional multivariate analysis; we aim to characterize probabilistic distributions over pairwise edge connections, while taking into accounting the higher-order structures such as transitivity (“a friend of a friend is a friend”), balance (“an enemy of a friend and an enemy”), and community (“cohesive subgroups of nodes”). **My project targets at developing higher-order tensor methods for analyzing high-dimensional network data.** The resulting prototypes will facilitate automatic detection of hidden salient structure in network data, thereby providing solutions to questions that cannot be addressed by existing methods. I will elaborate two specific directions for social network learning using tensor methods.

**Community detection in multi-layer networks.** A multi-layer network consists of multiple undirected graphs (or adjacency matrices), where each graph represents the connection among the same set of vertices (Fig 1a-b). The dataset is naturally organized as an order-3 tensor with the first two modes being vertices and the third mode being the contexts under which the graph is observed. Multilayer networks arise commonly in longitudinal study and multi-relational analysis. While the community structure in each single-layer network has been widely analyzed in the literature, little work has studied the heterogeneous pattern across multiple layers.

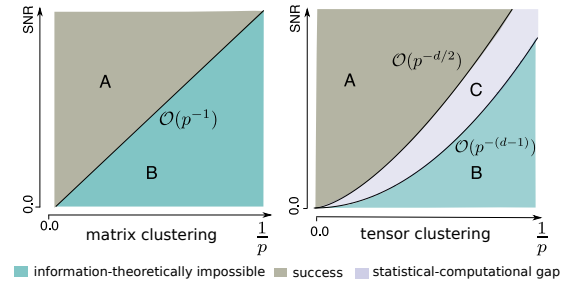


**Figure 1:** (a) Representation of network data using adjacency matrices [4]. (b) Stochastic block model for single-layer network (i.e., matrix). (c) We propose tensor extension of block model for high-order clustering in the context of multi-layer networks [5, 6].

Methods built on tensors provide generalized tools to capture complex data structure that the off-the-shelf methods may fail to exploit. We develop a tensor stochastic block model [6, 17] for simultaneous clustering of entities along each mode. Specifically, let  $\mathcal{Z} = \llbracket z_{i_1, \dots, i_d} \rrbracket \in \{0, 1\}^{p \times \dots \times p}$  denote the order- $d$  adjacency data tensor, where the entries  $z_{i_1, \dots, i_d}$  represent the presence or absence of edge  $(i_1, i_2)$  at context  $(i_3, \dots, i_d)$ . We model the signal tensor  $\mathbb{E}(\mathcal{Z})$  using block structure

$$\mathbb{E}(\mathcal{Z}) = \mathcal{S} \times_1 \mathbf{M}_1 \times_2 \dots \times_d \mathbf{M}_d,$$

for some low-dimensional core tensor  $\mathcal{S}$  and community membership matrices  $\mathbf{M}_1, \dots, \mathbf{M}_d$  (see Figure 1c). The learning goal is to estimate  $d$ -way connection strength tensor  $\mathcal{S}$  and community assignment  $(\mathbf{M}_1, \dots, \mathbf{M}_d)$  from a noisy observation  $\mathcal{Z}$ . We propose a higher-order Lloyd algorithm, which uses alternating optimization for parameter estimations. Our preliminary analysis shows that the tensor algorithm achieves *exact recovery of communities* under less stringent assumptions than existing algorithms. Surprisingly, we find that the learning performance is fully characterized by signal-to-noise ration (SNR) (see Figure 2). In the strong SNR region A, we prove that the our algorithm achieves exact clustering *in polynomial time*. We also show that the estimation error bound of the target tensor is *free of tensor dimension*. This feature is especially appealing in modern large-scale network analysis. In the weak SNR region B, we provide evidence to the



**Figure 2:** Comparison between matrix and tensor methods.



## References

- [1] Quentin Berthet and Philippe Rigollet, Complexity theoretic lower bounds for sparse principal component detection, *Conference on learning theory (COLT)*, 2013, pp. 1046–1066.
- [2] Peter J Bickel and Aiyu Chen, A nonparametric view of network models and Newman–Girvan and other modularities, *Proceedings of the National Academy of Sciences (PNAS)* 106 (2009), no. 50, 21068–21073.
- [3] Matthew Brennan, Guy Bresler, and Wasim Huleihel, Reducibility and computational lower bounds for problems with planted sparse structure, *Conference on learning theory (COLT)*, 2018, pp. 48–166.
- [4] Joshua Faskowitz, Xiaoran Yan, Xi-Nian Zuo, and Olaf Sporns, Weighted stochastic block models of the human connectome across the life span, *Scientific reports* 8 (2018), no. 1, 1–16.
- [5] Jiaxin Hu, Chanwoo Lee, and **Wang, Miaoyan**, Supervised tensor decomposition with interactive side information, *Advances in Neural Information Processing Systems (NeurIPS)* 33 Workshop on Machine Learning and the Physical Sciences, 2020.  
This work wins **Best Student Paper Award** from the Statistical Computing and Graphics Section of American Statistical Association (ASA), 2021.
- [6] Rungang Han, Yuetian Luo, **Miaoyan Wang**, and Anru R Zhang, Exact clustering in tensor block model: Statistical optimality and computational limit, arXiv preprint arXiv:2012.09996 (2020).  
This work wins **Best Student Paper Award** from the Statistical Learning and Data Science Section of the American Statistical Association (ASA), 2021.
- [7] Chanwoo Lee and **Miaoyan Wang**, Tensor denoising and completion based on ordinal observations, *International conference on machine learning (ICML)*, 2020, pp. 5778–5788.
- [8] Chanwoo Lee and **Miaoyan Wang**, Beyond the signs: Nonparametric tensor completion via sign series, arXiv preprint arXiv:2102.00384 (2021).
- [9] Hanbaek Lyu, Deanna Needell, and Laura Balzano, Online matrix factorization for markovian data and applications to network dictionary learning, *Journal of Machine Learning Research* 21 (2020), no. 251, 1–49.
- [10] Karl Rohe, Sourav Chatterjee, and Bin Yu. Spectral clustering and the high-dimensional stochastic blockmodel. *Annals of Statistics*, Vol 39, 4, 1878–1915, 2011.
- [11] **Miaoyan Wang**, Khanh Dao Duc, Jonathan Fischer, and Yun S Song, Operator norm inequalities between tensor unfoldings on the partition lattice, *Linear Algebra and Its Applications* 520 (2017), 44–66.
- [12] **Miaoyan Wang**, Jonathan Fischer, and Yun S Song, Three-way clustering of multi-tissue multi-individual gene expression data using semi-nonnegative tensor decomposition, *The Annals of Applied Statistics* 13 (2019), no. 2, 1103–1127.
- [13] **Miaoyan Wang** and Lexin Li, Learning from binary multiway data: Probabilistic tensor decomposition and its statistical optimality, *Journal of Machine Learning Research* 21 (2020), no. 154, 1–38.
- [14] **Miaoyan Wang**, Fabrice Roux, Claudia Bartoli, Carine Huard-Chauveau, Christopher Meyer, Hana Lee, Dominique Roby, Mary Sara McPeck, and Joy Bergelson, Two-way Mixed-effects methods for joint association analysis using both host and pathogen genomes, *Proceedings of the National Academy of Sciences (PNAS)* 115 (2018), no. 24, E5440–E5449.
- [15] **Miaoyan Wang**, Johanna Jakobsdottir, Albert V. Smith, and Mary Sara McPeck. G-STRATEGY: Optimal selection of individuals for sequencing in genetic association studies. *Genetic Epidemiology*, Vol. 40, No. 6, 446–460, 2016.  
Highlighted as **Editor’s Pick Paper** of this issue. This work wins **ASHG Charles J. Epstein Trainee Award** and **IGES Williams Award**.
- [16] **Miaoyan Wang** and Yun Song, Tensor decompositions via two-mode higher-order SVD (HOSVD), *Artificial intelligence and statistics*, 2017, pp. 614–622.
- [17] **Miaoyan Wang** and Yuchen Zeng, Multiway clustering via tensor block models, *Advances in Neural Information Processing Systems (NeurIPS)* 32 . (2019).
- [18] Yihong Wu and Jiaming Xu, Statistical problems with planted structures: Information-theoretical and computational limits, *Information-Theoretic Methods in Data Science* (2021), 383.
- [19] L Zheng and G Raskutti. Testing for high-dimensional network parameters in auto-regressive models. *Electronic Journal of Statistics*, 13(2): 4977–5043 (2019).