

Sylvester Graphical Models for Complex Spatio-Temporal Processes

Yu (Wayne) Wang

Joint work with Byoungwook Jang and Alfred Hero

University of Michigan

Outline

- 1 Multi-indexed Data and Kronecker Structured Graphical Models
- 2 Sylvester Graphical Model
- 3 Application to Solar Flare Prediction
- 4 Summary

Outline

- 1 Multi-indexed Data and Kronecker Structured Graphical Models
- 2 Sylvester Graphical Model
- 3 Application to Solar Flare Prediction
- 4 Summary

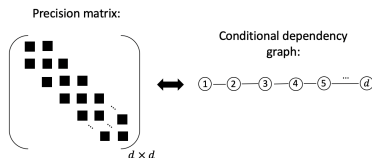
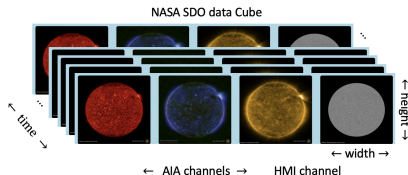
Ultra-high dimensional multi-indexed heterogeneous data

Soln: a structured multiway graphical model

- SyGlasso/SG-PALM sqrt. Kronecker sum decomp.
 - Generative: $AX + XB = Z$
 - Physical interpretability: directly related to PDEs
 - Robustness, low runtime and high accuracy

Applications

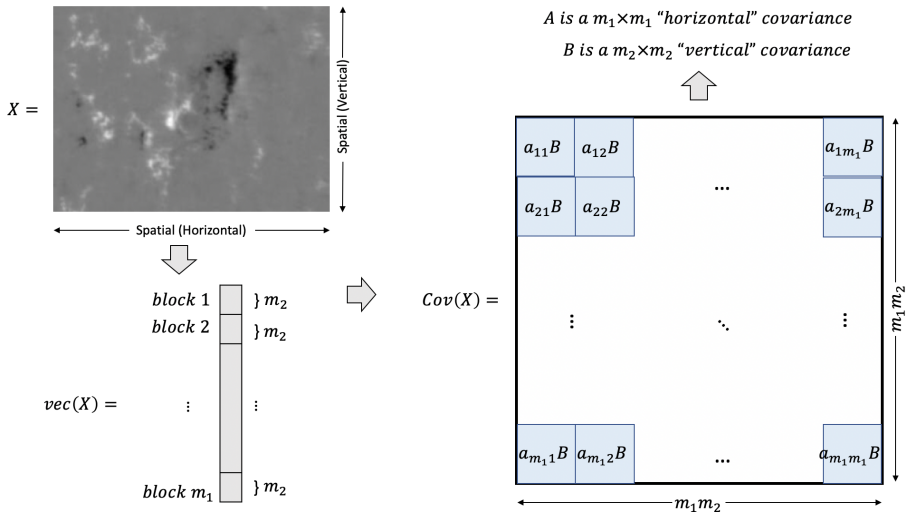
- Astrophysics: solar flare prediction
- Climatology: spatio-temporal weather forecasting
- Neuroscience: EEG analysis
- Radar imaging: STAP, SAR
- Computer vision: video sequence prediction



Challenges

- High dim. $d = \prod_{k=1}^K m_k$ and $\mathcal{X} \in \mathbb{R}^{m_1 \times \dots \times m_K}$
- Non-commutative
- Sample-starved learning and high computational complexity

Kronecker product representation for matrix-variate data



Sparse Kronecker-structured models

Sparsity models:

- Sparse covariance (Σ) models for marginal dependencies:
 - Examples: M-dependent processes, moving average (MA) processes.
- Sparse **precision** ($\Omega = \Sigma^{-1}$) models for conditional dependencies:
 - Examples: Markov random fields, autoregressive (AR) processes.

Kronecker-structured models:

- Kronecker product covariance/inverse covariance
 - Transposable regularized covariance¹
 - KGlasso²
 - GEMINI³
- Kronecker sum covariance/inverse covariance
 - Kronecker sum covariance ($K = 2$) for error-in-variable models⁴
 - Bigraphical Lasso (BiGlasso)⁵
 - Tensor-graphical Lasso (TeraLasso)⁶

¹Genera I Allen and Robert Tibshirani. “Transposable regularized covariance models with an application to missing data imputation”. In: *The Annals of Applied Statistics* 4.2 (2010), p. 764.

²Theodoros Tsiligkaridis, Alfred O Hero III, and Shuheng Zhou. “On convergence of kronecker graphical lasso algorithms”. In: *IEEE transactions on signal processing* 61.7 (2013), pp. 1743–1755.

³Shuheng Zhou. “GEMINI: Graph Estimation with Matrix Variate Normal Instances”. In: *The Annals of Statistics* 42.2 (2014), pp. 532–562.

⁴Mark Rudelson and Shuheng Zhou. “Errors-in-variables models with dependent measurements”. In: *Electronic Journal of Statistics* 11.1 (2017), pp. 1699–1797.

⁵Alfredo Kalaitzis et al. “The bigraphical lasso”. In: *International Conference on Machine Learning*. 2013, pp. 1229–1237.

⁶Kristjan Greenewald, Shuheng Zhou, and Alfred Hero III. “Tensor graphical lasso (TeraLasso)”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 81.5 (2019), pp. 901–931.

Kronecker sum vs Kronecker product

For $\mathbf{A} \in \mathbb{R}^{m_1 \times m_1}$, $\mathbf{B} \in \mathbb{R}^{m_2 \times m_2}$:

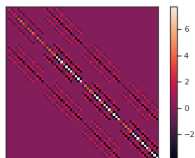
$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & \cdots & a_{1m_1}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m_11}\mathbf{B} & \cdots & a_{m_1m_1}\mathbf{B} \end{bmatrix},$$

and

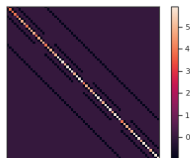
$$\mathbf{A} \oplus \mathbf{B} = \mathbf{I}_{m_2} \otimes \mathbf{A} + \mathbf{B} \otimes \mathbf{I}_{m_1}.$$



Ψ_k



KP $\Omega = \bigotimes_{k=1}^3 \Psi_k$



KS $\Omega = \bigoplus_{k=1}^3 \Psi_k$

KS vs KP

KP models:

- Separable: $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$ and $\det(\mathbf{A} \otimes \mathbf{B}) = (\det \mathbf{A})^{m_2} (\det \mathbf{B})^{m_1}$.
- Biconvex negative log-likelihood: estimation using alternating graphical lasso, e.g., Tlasso⁷.
- Generative representation: $\mathbf{X} = \mathbf{C}^{-1}\mathbf{Z}\mathbf{D}^{-1}$, where $\mathbf{A} = \mathbf{C}\mathbf{C}^T$, $\mathbf{B} = \mathbf{D}\mathbf{D}^T$, and \mathbf{Z} white noise $\Rightarrow \text{Cov}^{-1}(\mathbf{X}) = \mathbf{A} \otimes \mathbf{B}$.

KS models:

- Parsimonious: Cartesian product of graphs avoids explosion of edges⁸.
- Convex negative log-likelihood: estimation using ISTA-type procedure, e.g., TeraLasso⁹.
- No obvious generative representation.

Motivating question: KP + KS?

⁷Xiang Lyu et al. "Tensor Graphical Model: Non-convex Optimization and Statistical Inference". In: *IEEE transactions on pattern analysis and machine intelligence* (2019).

⁸Alfredo Kalaitzis et al. "The bigraphical lasso". In: *International Conference on Machine Learning*. 2013, pp. 1229–1237.

⁹Kristjan Greenewald, Shuheng Zhou, and Alfred Hero III. "Tensor graphical lasso (TeraLasso)". In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 81.5 (2019), pp. 901–931.

Outline

- 1 Multi-indexed Data and Kronecker Structured Graphical Models
- 2 Sylvester Graphical Model**
- 3 Application to Solar Flare Prediction
- 4 Summary

Sylvester Glasso: a generative Kronecker sum model

Let a random tensor $\mathcal{X} \in \mathbb{R}^{m_1 \times \cdots \times m_K}$ be generated by the Sylvester tensor equation^{10,11}:

$$\mathcal{X} \times_1 \Psi_1 + \cdots + \mathcal{X} \times_K \Psi_K = \mathcal{T} \quad \Leftrightarrow \quad (\Psi_1 \oplus \cdots \oplus \Psi_K) \text{vec}(\mathcal{X}) = \text{vec}(\mathcal{T}),$$

where $\Psi_k \in \mathbb{R}^{m_k \times m_k}$, $k = 1, \dots, K$ are sparse matrices, \mathcal{T} is a random tensor of the same dimension as \mathcal{X} , and \times_k a k -mode product.

- If $\text{vec}(\mathcal{T}) \sim \mathcal{N}(0, \mathbf{I}_p)$, then $\text{Cov}^{-1}(\text{vec}(\mathcal{X})) = (\Psi_1 \oplus \cdots \oplus \Psi_K)^2$.
- The Ψ_k 's can be estimated by minimizing the negative log-pseudolikelihood defined as

$$\begin{aligned} \mathcal{L}_\lambda(\{\Psi_k\}_{k=1}^K) = & -\frac{N}{2} \log |(\text{diag}(\Psi_1) \oplus \cdots \oplus \text{diag}(\Psi_K))^2| \\ & + \frac{N}{2} \text{tr}(\mathbf{S} \cdot (\Psi_1 \oplus \cdots \oplus \Psi_K)^2) + \sum_{k=1}^K \lambda_k \|\Psi_k\|_{1,\text{off}}, \end{aligned}$$

using a proximal alternating linearized minimization method, called SG-PALM.

¹⁰Lars Grasedyck. "Existence and computation of low Kronecker-rank approximations for large linear systems of tensor product structure". In: *Computing* 72.3-4 (2004), pp. 247–265.

¹¹Daniel Kressner and Christine Tobler. "Krylov subspace methods for linear systems with tensor product structure". In: *SIAM journal on matrix analysis and applications* 31.4 (2010), pp. 1688–1714.

SG-PALM: optimization convergence¹²

Lemma

The function H (the log det plus tr terms) is convex and continuously differentiable on an open set containing $\text{dom}G$ (G the regularization term) and its gradient is block-wise Lipschitz continuous. Further, the objective function $\mathcal{L}_\lambda(\Psi)$ satisfies the Kurdyka - Łojasiewicz (KL) property with a KL exponent of $\frac{1}{2}$.

Theorem

Let $\{\Psi^{(t)}\}_{t \geq 0}$ be generated by SG-PALM and assume that $\Psi^{(t)} \in \Omega \subset \text{dom} \partial \mathcal{L}_\lambda$. Then, the SG-PALM converges linearly in the sense that

$$\begin{aligned} & \frac{\mathcal{L}_\lambda(\Psi^{(t+1)}) - \min \mathcal{L}_\lambda}{\mathcal{L}_\lambda(\Psi^{(t)}) - \min \mathcal{L}_\lambda} \\ & \leq \left(\frac{\alpha^2 L_{\min}}{4Kc^2(\sum_{j=1}^K L_j)^2 + 4c^2 L_{\max}} + 1 \right)^{-1}, \end{aligned}$$

where $L_{\min} = \min_j L_j$, $L_{\max} = \max_j L_j$, $\alpha > 0$, and $c \in (0, 1)$.

¹²Yu Wang and Alfred Hero. "A Proximal Alternating Linearized Minimization Method for Tensor Graphical Models". In: *Submitted* (2020).

SyGlasso: statistical convergence¹³

Assume: sub-Gaussianity on \mathcal{X} , bounded eigenvalues of Ω , and incoherence condition on the loss function. Further, let $q_k := |\{(i, j) : (\Psi_k)_{i,j} \neq 0, i \neq j\}|$, and

$$\lambda_k = O\left(\sqrt{\frac{m_k \log d}{N}}\right), \quad N > O(\max_k q_k m_k \log d).$$

Theorem (Graph recovery consistency)

There exists a constant $C(\beta) > 0$ ^a such that for any $c_0 > 0$ the following events hold with probability at least $1 - O(\exp(-c_0 \log p))$:

- (Estimation consistency) Any minimizer $\hat{\beta}$ satisfies:

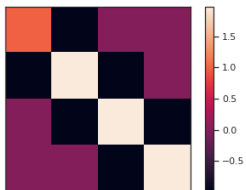
$$\|\hat{\beta} - \beta\|_2 \leq C(\beta) \sqrt{K} \max_k \sqrt{q_k} \lambda_k.$$

- (Sign consistency) If minimal signal strength is satisfied for Ψ_k for each k , then $\text{sign}(\hat{\beta}) = \text{sign}(\beta)$.

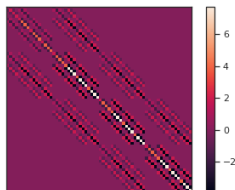
^a β denotes all off-diagonal elements of Ψ_k 's.

¹³Yu Wang, Byoungwook Jang, and Alfred Hero. "Sylvester Graphical Lasso (SyGlasso)". In: *Proceedings of The 23rd International Conference on Artificial Intelligence and Statistics (AISTATS)* (2020).

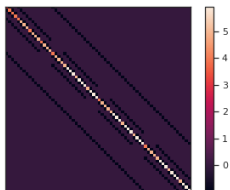
Comparison with KS and KP



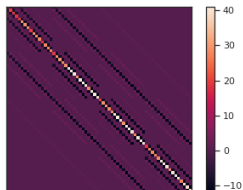
(a) Ψ_k



(b) KP Ω



(c) KS Ω

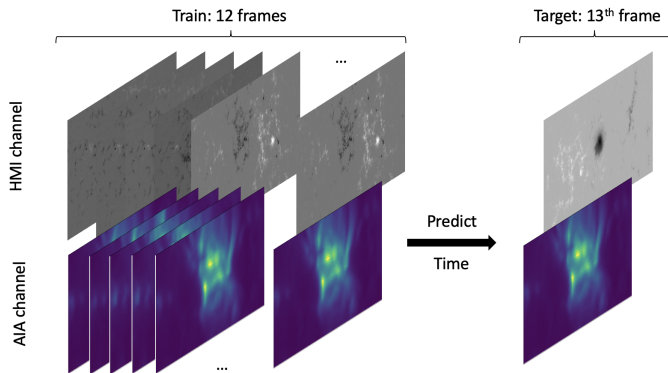


(d) SyGlasso Ω

Outline

- 1 Multi-indexed Data and Kronecker Structured Graphical Models
- 2 Sylvester Graphical Model
- 3 Application to Solar Flare Prediction**
- 4 Summary

Multi-instrument solar imaging data



Raw Data: 319 videos before weak or strong flares. Each sequence is $\mathcal{X} \in \mathbb{R}^{m_{time} \times m_{width} \times m_{height} \times m_{channel}}$, where $m_{time} = 13$ (13-hour window with 1-hour cadence), $m_{width} = 100$, $m_{height} = 50$, and $m_{channel} = 7$ (3 HMI channels and 4 AIA channels).

Goal: Constructing linear forward time series predictors for the last frame (at or right before a flare) by using estimated precision matrix from all previous frames.

Multi-output sparse regression for time series prediction

Consider a spatio-temporal process observed at p time stamps and q locations, i.e., $\mathbf{X} \in \mathbb{R}^{p \times q}$. The estimated precision matrix $\mathbf{\Omega} \in \mathbb{R}^{pq \times pq}$ can be used to construct an optimal linear predictor, for example,

$$\mathbf{y}_t = \mathbf{\Omega}_{2,2}^{-1} \mathbf{\Omega}_{2,1} \mathbf{y}_{t-1:t-(p-1)},$$

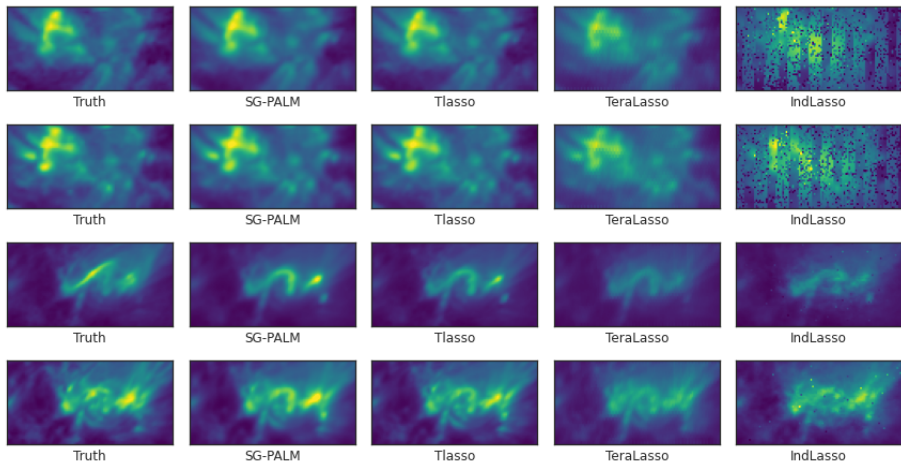
where $\mathbf{y}_t \in \mathbb{R}^q$, $\mathbf{y}_{t-1:t-(p-1)} \in \mathbb{R}^{(p-1)q}$, and $\mathbf{\Omega}_{2,2} \in \mathbb{R}^{q \times q}$, $\mathbf{\Omega}_{2,1} \in \mathbb{R}^{q \times (p-1)q}$ are appropriate submatrices of $\mathbf{\Omega}$.

For the solar flare data, we estimate

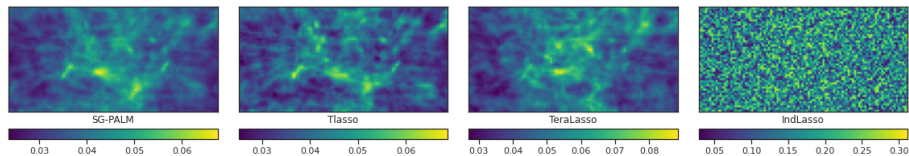
$$\mathbf{\Omega} = (\mathbf{\Psi}_{time} \oplus \mathbf{\Psi}_{height} \oplus \mathbf{\Psi}_{width} \oplus \mathbf{\Psi}_{channel})^2 \in \mathbb{R}^{455000 \times 455000}$$

in training and predict $\mathbf{y}_{13} \in \mathbb{R}^{35000}$ in testing using estimated $\mathbf{\Omega}$.

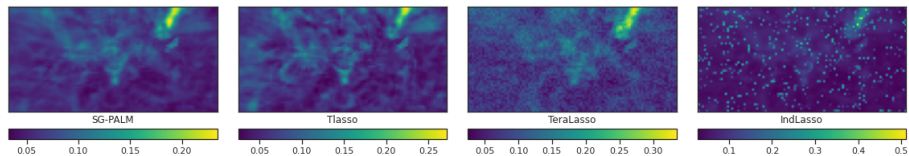
Real vs. predicted images



NRMSE



(a) B-class. Avg. NRMSE (left to right): 0.0379, 0.0386, 0.0579, 0.1628.



(b) MX-class. Avg. NRMSE (left to right): 0.0620, 0.0790, 0.0913, 0.1172.

Outline

- 1 Multi-indexed Data and Kronecker Structured Graphical Models
- 2 Sylvester Graphical Model
- 3 Application to Solar Flare Prediction
- 4 Summary**

Summary

Methodology:

- Kronecker-structured graphical modeling framework inspired by the Sylvester equations in physics.

Theory:

- Convergence (with geometric rate) of the optimization error.
- Convergence of the statistical error.

Applications:

- Multi-modal solar imaging data and flare prediction.

Thank you! Questions?