Nonconvex Sparse Graph Learning under Laplacian Constrained Graphical Model Jiaxi Ying, José Vinícius de M. Cardoso, and Daniel P. Palomar

Graphical Model



- A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{W})$ is an intuitive way to represent relationships between entities.
- Nodes: $\mathcal{V} = \{1, 2, \dots, p\}$ correspond to the entities.
- Edges: $\mathcal{E} = \{(1, 2), \dots, (i, j), \dots\}$ encodes conditional dependence between entities.
- Weights: W is weight matrix with W_{ij} the graph weight between node i and node j.

1.1 Laplacian Constrained Gaussian Graphical Model

• Graph Laplacian: $m{L} = m{D} - m{W}$, where $m{D}$: the degree matrix with $D_{ii} = \sum_{i=1}^p W_{ij}$.

$$S_L$$
: the set of all graph Laplacians for connected graphs,

 $\mathcal{S}_L := \{ \boldsymbol{\Theta} \in \mathcal{S}^p_+ | \, \Theta_{ij} = \Theta_{ji} \le 0, \, \forall \, i \ne j, \, \boldsymbol{\Theta} \cdot \mathbf{1} = \mathbf{0}, \, \operatorname{rank}(\boldsymbol{\Theta}) = p - 1 \},$ where 0 and 1 are the constant zero and one vectors.

Definition 1 (L-GMRF). A zero-mean random vector $\boldsymbol{x} = [x_1, \dots, x_p]^\top \in V^{p-1}$ is called a Laplacian constrained Gaussian Markov Random Fields (L-GMRF) with parameters $(\mathbf{0}, \mathbf{\Theta})$ with $\Theta \in \mathcal{S}_L$, if and only if its density function $q_L : V^{p-1} \to \mathbb{R}$ follows

$$q_L(\boldsymbol{x}) = (2\pi)^{-\frac{p-1}{2}} \det^{\star}(\boldsymbol{\Theta})^{\frac{1}{2}} \exp\Big(-\frac{1}{2}\boldsymbol{x}^{\top}\boldsymbol{\Theta}\boldsymbol{x}\Big),$$

where det^{*}: the product of nonzero eigenvalues; $V^{p-1} := \{ \boldsymbol{x} \in \mathbb{R}^p | \boldsymbol{1}^\top \boldsymbol{x} = 0 \}.$

ℓ_1 -norm Analysis and Proposed Method

ℓ_1 -norm Does Not Work 2.1

The ℓ_1 -norm regularized maximum likelihood estimation of Laplacian constrained precision matrices [1, 2] can be formulated as

$$\min_{\boldsymbol{\Theta} \in \mathcal{S}_L} -\log \det(\boldsymbol{\Theta} + \boldsymbol{J}) + \operatorname{tr}(\boldsymbol{\Theta} \boldsymbol{S}) + \lambda \sum_{i>i} |\Theta_{ij}|,$$

where S: sample covariance matrix; λ : regularization parameter; J: constant matrix with each element equal to $\frac{1}{n}$.

Theorem 2. Let $\hat{\Theta} \in \mathbb{R}^{p \times p}$ be the global minimum of (3) with p > 3. Define $s_1 = \max_k S_{kk}$ and $s_2 = \min_{ij} S_{ij}$. If the regularization parameter λ in (3) satisfies $\lambda \in [(2 + 2\sqrt{2})(p + 1)]$ $1)(s_1 - s_2), +\infty)$, then the estimated graph weight $\hat{W}_{ij} = -\hat{\Theta}_{ij}$ obeys

$$\hat{W}_{ij} \ge \frac{1}{(s_1 - (p+1)s_2 + \lambda)p} > 0, \quad \forall i \neq j.$$

Theorem 2 shows that a large regularization parameter of the ℓ_1 -norm will make every graph weight strictly positive and the estimated graph is fully connected.



Figure 1: Graph learning using ℓ_1 -norm with different regularization parameters. The number of nonzero edges in (a), (b), (c) and (d) are 49, 135, 286 and 1225.

For more information visit: https://www.danielppalomar.com; R Package: https://github.com/mirca/sparseGraph

Proposed Method 2.2

The penalized maximum likelihood of the precision matrix with Laplacian structural constraints can be formulated as

 $\min_{\boldsymbol{\Theta} \in \mathcal{S}_L} -\log \det(\boldsymbol{\Theta} + \boldsymbol{J}) + \mathsf{tr}\left(\boldsymbol{\Theta}\boldsymbol{S}\right) + \mathbf{I}$

where h_{λ} : nonconvex sparsity-promoting function such as **SCAD** and **MCP**. To handle the constraint $\Theta \in S_L$, we introduce a linear operator \mathcal{L} [3] that maps a vector \boldsymbol{w} to a Laplacian matrix $\mathcal{L}\boldsymbol{w}$. For example $\boldsymbol{w} = [w_1, w_2, w_3, w_4, w_5, w_6]^{\top}$,

$$\mathcal{L}\boldsymbol{w} = \begin{bmatrix} \sum_{i=1,2,3} w_i & -w_1 & -w_2 & -w_3 \\ -w_1 & \sum_{i=1,4,5} w_i & -w_4 & -w_5 \\ -w_2 & -w_4 & \sum_{i=2,4,6} w_i & -w_6 \\ -w_3 & -w_5 & -w_6 & \sum_{i=3,5,6} w_i \end{bmatrix}$$

With the usage of the linear operator \mathcal{L} , the optimization (4) can be reformulated as

 $\min_{\boldsymbol{w} \ge \boldsymbol{0}} -\log \det(\mathcal{L}\boldsymbol{w} + \boldsymbol{J}) + \mathsf{tr}\left(\boldsymbol{S}\mathcal{L}\boldsymbol{w}\right) -$

We establish a sequence $\{\hat{w}^{(k)}\}_{k>1}$ by solving a sequence of sub-problems

$$\hat{\boldsymbol{w}}^{(k)} = \arg\min_{\boldsymbol{w} \ge \boldsymbol{0}} -\log\det(\mathcal{L}\boldsymbol{w} + \boldsymbol{J}) + \operatorname{tr}(\boldsymbol{S}\mathcal{L}\boldsymbol{w}) + \sum_{i} h'_{\lambda}\left(\hat{w}_{i}^{(k-1)}\right)w_{i}.$$
 (6)

The optimization (6) can be solved by a projected gradient descent algorithm with backtracking line search.

Algorithm 1 Nonconvex Graph Learning (NGL)

Input: Sample covariance $S, \lambda, \hat{w}^{(0)}$;

- $k \leftarrow 1$:
- 1: while Stopping criteria not met do
- 2: Update $z_i^{(k-1)} = h'_\lambda(\hat{w}_i^{(k-1)})$, for $i=1,\ldots,p(p-1)$
- Update $\hat{w}^{(k)} = \arg\min_{w>0} \log\det(\mathcal{L}w + J) + tr$ $k \leftarrow k+1;$
- 5: end while
- Output: $\hat{m{w}}^{(k)}$.

Theoretical Results 2.3

Assumption 3. The function $h_{\lambda} : \mathbb{R} \to \mathbb{R}$ satisfies the following conditions: 1. $h_{\lambda}(0) = 0$, and $h'_{\lambda}(x)$ is monotone and Lipschitz continuous for $x \in [0, +\infty)$; **2.** There exists a $\gamma > 0$ such that $h'_{\lambda}(x) = 0$ for $x \ge \gamma \lambda$; **3.** $h'_{\lambda}(x) = \lambda$ for $x \leq 0$ and $h'_{\lambda}(c\lambda) \geq \lambda/2$, where $c = (2 + \sqrt{2})\lambda_{\max}^2(\mathcal{L}w^*)$ is a constant. **Assumption 4.** The minimal nonzero graph weight satisfies $\min_{i \in S^*} w_i^* \ge (c + \gamma) \lambda \gtrsim \lambda$, where c and γ are defined in Assumption 3.

Theorem 5. Under Assumptions 3 and 4, take the regularization parameter $\lambda =$ $\sqrt{4\alpha c_0^{-1} \log p/n}$ for some $\alpha > 2$. If the sample size n is lower bounded by

 $n \ge \max(94\alpha c_0^{-1}\lambda_{\max}^2(\mathcal{L}\boldsymbol{w}^{\star})s\log p, 8\alpha\log p),$

then with probability at least $1-1/p^{\alpha-2}$, the sequence $\hat{w}^{(k)}$ returned by Algorithm 1. satisfies

$$\|\hat{\boldsymbol{w}}^{(k)} - \boldsymbol{w}^{\star}\| \leq \underbrace{2(3\sqrt{2}+4)\lambda_{\max}^{2}(\mathcal{L}\boldsymbol{w}^{\star})\sqrt{\alpha c_{0}^{-1}s\log p/n}}_{Statistical\ error} + \underbrace{\left(\frac{3}{2+\sqrt{2}}\right)^{k}\|\hat{\boldsymbol{w}}^{(0)} - \boldsymbol{w}^{\star}\|}_{Optimization\ error},$$

where $c_0 = 1/(8 \| \mathcal{L}^* (\mathcal{L} w^* + J)^{-1} \|_{\max}^2)$ is a constant.

Theorem 5 shows that the estimation error is bounded by the optimization error and statistical error. The optimization error decays to zero at a linear rate with respect to the iteration number k. The statistical error is with the order of $\sqrt{s \log p/n}$, and a large sample size n will lead to a small statistical error.

(1)

(2)

$$\sum_{i>j} h_{\lambda}(\Theta_{ij}),$$

(4)

$$+\sum_{i}h_{\lambda}(w_{i}).$$
 (5)

$$)/2;$$

 $(oldsymbol{S}\mathcal{L}oldsymbol{w}) + \sum_i z_i^{(k-1)} w_i;$

Experimental Results Synthetic Data 3.1

The data matrix $X \in \mathbb{R}^{p imes n}$ with each column independently sampled from L-GMRF. The ground-truth graph is a random Barabasi-Albert graph with 50 nodes, and the weights are randomly sampled from U(2,5). The compared methods include the state-of-the-art GLE-ADMM algorithm [2] and the baseline projected gradient descent with ℓ_1 -norm.



Figure 2: Performance measures (a) Number of positive edges, (b) Relative error and (c) F-score as a function of regularization parameter λ .



(c) F-score as a function of the sample size ratio n/p.

3.2 Real-world Data

The data set is 2019-nCoV [4] from 98 Chinese patients affected by the outbreak of 2019-nCoV on early February, 2020. The features include age, gender, and location. The labels represent the life status of patients, alive (green) or no longer alive (red).



Figure 4: The learned graphs using the 2019-nCoV data set by (a) GLE-ADMM, (b) NGL-SCAD (proposed method), and (c) NGL-MCP (proposed method).

References

- of Selected Topics in Signal Processing, vol. 11, no. 6, pp. 825–841, 2017.
- *IEEE Transactions on Signal Processing*, vol. 67, no. 16, pp. 4231–4244, 2019. [3] S. Kumar, J. Ying, J. V. d. M. Cardoso, and D. P. Palomar, "A unified framework for structured graph learning via spectral
- constraints," *Journal of Machine Learning Research*, vol. 21, no. 22, pp. 1–60, 2020.
- 2020. models?" arXiv preprint arXiv:2006.14925, 2020.



THE HONG KONG **UNIVERSITY OF SCIENCE** AND TECHNOLOGY

Figure 3: Performance measures (a) Number of positive edges, (b) Relative error and

[1] H. E. Egilmez, E. Pavez, and A. Ortega, "Graph learning from data under laplacian and structural constraints," *IEEE Journal*

[2] L. Zhao, Y. Wang, S. Kumar, and D. P. Palomar, "Optimization algorithms for graph Laplacian estimation via ADMM and MM,"

[4] T. Wu, X. Ge, G. Yu, and E. Hu, "Open-source analytics tools for studying the covid-19 coronavirus outbreak," medRxiv,

[5] J. Ying, J. V. d. M. Cardoso, and D. P. Palomar, "Does the ℓ_1 -norm learn a sparse graph under laplacian constrained graphical