

# On High-Dimensional Graph Learning Under Total Positivity

Jitendra K. Tugnait

Dept. of Electrical & Computer Eng.  
Auburn University, Auburn, AL 36849, USA

**Abstract**—We consider the problem of estimating the structure of an undirected weighted sparse graphical model of multivariate data under the assumption that the underlying distribution is multivariate totally positive of order 2, or equivalently, all partial correlations are non-negative. Total positivity holds in several applications. The problem of Gaussian graphical model learning has been widely studied without the total positivity assumption where the problem can be formulated as estimation of the sparse precision matrix that encodes conditional dependence between random variables associated with the graph nodes. An approach that imposes total positivity is to assume that the precision matrix obeys the Laplacian constraints which include constraining the off-diagonal elements of the precision matrix to be non-positive. In this paper we investigate modifications to widely used penalized log-likelihood approaches to enforce total positivity but not the Laplacian structure. An alternating direction method of multipliers (ADMM) algorithm is presented for constrained optimization under total positivity and lasso as well as adaptive lasso penalties. Numerical results based on synthetic data show that the proposed constrained adaptive lasso approach significantly outperforms existing Laplacian-based approaches, both statistical and smoothness-based non-statistical.

## I. INTRODUCTION

An undirected simple weighted graph is denoted  $\mathcal{G} = (V, \mathcal{E}, \mathbf{W})$  where  $V = \{1, 2, \dots, p\} = [p]$  is the set of  $p$  nodes,  $\mathcal{E} \subseteq [p] \times [p]$  is the set of undirected edges, and  $\mathbf{W} = \mathbf{W}^\top \in \mathbb{R}^{p \times p}$  stores the non-negative weights  $W_{ij} \geq 0$  associated with the undirected edges. If there is an edge between nodes  $i$  and  $j$ , then edge  $\{i, j\} \in \mathcal{E}$  and  $W_{ij} > 0$ . If there is no edge between nodes  $i$  and  $j$ , then edge  $\{i, j\} \notin \mathcal{E}$  and  $W_{ij} = 0$ . In a simple graph there are no self-loops or multiple edges, so  $\mathcal{E}$  consists of distinct pairs  $\{i, j\}$ ,  $i \neq j$  and  $W_{ii} = 0$ . In an undirected graph, if  $\{i, j\} \in \mathcal{E}$ , then  $\{j, i\} \in \mathcal{E}$ . In graphical models of data variables  $x_1, x_2, \dots, x_p$ , ( $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_p]^\top$ ), a weighted graph  $\mathcal{G} = (V, \mathcal{E}, \mathbf{W})$  (or unweighted  $\mathcal{G} = (V, \mathcal{E})$ ) with  $|V| = p$  is used to capture relationships between the  $p$  variables  $x_i$ s [1]–[3]. If  $\{i, j\} \in \mathcal{E}$ , then  $x_i$  and  $x_j$  are related (similar or dependent) in some sense, with higher  $W_{ij}$  indicating stronger similarity or dependence.

Graphical models provide a powerful tool for analyzing multivariate data [1]–[3]. In a statistical graphical model, the conditional statistical dependency structure among  $p$  random variables  $x_1, x_2, \dots, x_p$ , is represented using an undirected

graph  $\mathcal{G} = (V, \mathcal{E})$ . The graph  $\mathcal{G}$  then is a conditional independence graph (CIG) where there is no edge between nodes  $i$  and  $j$  (i.e.,  $\{i, j\} \notin \mathcal{E}$ ) iff  $x_i$  and  $x_j$  are conditionally independent given the remaining  $p-2$  variables  $x_\ell$ ,  $\ell \in [p]$ ,  $\ell \neq i, \ell \neq j$ . In particular, Gaussian graphical models (GGMs) are CIGs where  $\mathbf{x}$  is multivariate Gaussian. Suppose  $\mathbf{x}$  has positive-definite covariance matrix  $\Sigma$  with inverse covariance matrix (also called precision matrix)  $\Omega = \Sigma^{-1}$ . Then  $\Omega_{ij}$ , the  $(i, j)$ -th element of  $\Omega$ , is zero iff  $x_i$  and  $x_j$  are conditionally independent. Such models for  $\mathbf{x}$  have been extensively studied. Given  $n$  samples of  $\mathbf{x}$ , in *high-dimensional settings* where  $p \gg 1$  and/or  $n$  is of the order of  $p$ , one estimates  $\Omega$  under some sparsity constraints; see [4]–[8].

More recently, several authors have considered Gaussian graphical models under the constraint that the distribution is multivariate totally positive of order 2 (MTP<sub>2</sub>), or equivalently, that all partial correlations are non-negative (see [9], [10] and references therein). Such models are also known as attractive Gaussian random fields [11]. Note that a Gaussian distribution is MTP<sub>2</sub> if and only if its precision matrix  $\Omega$  is an M-matrix, i.e.,  $\Omega_{ij} \leq 0$  for all  $i \neq j$  [12]. As discussed in [9], MTP<sub>2</sub> is a strong form of positive dependence, which is relevant for modeling in various applications including phylogenetics or portfolio selection, where the shared ancestry or latent global market variable often lead to positive dependence among the observed variables.

On the other hand, graphical models for data variables have been inferred from consideration other than statistical, depending upon the intended application, nature of data and available prior information [1]. One class of graphical models are based on signal smoothness [1], [13]–[15]. Suppose we are given  $n$  samples  $\{\mathbf{x}(t)\}_{t=1}^n$  of the  $p$  data variables  $x_1, x_2, \dots, x_p$ , with  $\mathbf{x}(t) = [x_1(t) \ x_2(t) \ \dots \ x_p(t)]^\top$ . Define the  $p \times n$  matrix

$$\mathbf{X} = [\mathbf{x}(1) \ \mathbf{x}(2) \ \dots \ \mathbf{x}(n)]. \quad (1)$$

A measure of smoothness of signal  $\mathbf{x}(t)$  under which the signal takes “similar” values at “neighboring” vertices of a given weighted undirected graph, is the function [1], [13], [14]

$$\frac{1}{2} \sum_{i,j=1}^p W_{ij} \|\mathbf{X}_i - \mathbf{X}_j\|^2 = \text{tr}(\mathbf{X}^\top \mathbf{L} \mathbf{X}) \quad (2)$$

where  $\mathbf{X}_i$  denotes the  $i$ th row of  $\mathbf{X}$ ,  $\mathbf{L} = \mathbf{D} - \mathbf{W}$  is the (combinatorial) graph Laplacian (matrix), and  $\mathbf{D}$  is the diagonal weighted degree matrix with  $D_{ii} = \sum_{j=1}^p W_{ij}$ .

This work was supported by NSF Grants CCF-1617610 and ECCS-2040536. Author’s email: tugnajk@auburn.edu

Graph learning from data  $\mathbf{X}$  then becomes equivalent to estimation of the graph Laplacian matrix  $\mathbf{L}$  [1], [13].

Another set of approaches are based on statistical considerations under the graph Laplacian constraint [1], [16]–[18] where Laplacian  $\mathbf{L}$ , after regularization, plays the role of inverse covariance  $\mathbf{\Omega}$ ;  $\mathbf{L}$  is a symmetric, non-negative-definite matrix but with non-positive off-diagonal entries. Thus, under Gaussian distribution we have an MTP<sub>2</sub> model.

Graph Laplacian matrix has been extensively used for embedding, manifold learning, clustering and semi-supervised learning [19]–[24].

In this paper we investigate modifications to widely used penalized log-likelihood approaches to enforce total positivity but not the Laplacian structure. An alternating direction method of multipliers (ADMM) algorithm is presented for constrained optimization under total positivity and lasso as well as adaptive lasso penalties. Numerical results based on synthetic data show that the proposed constrained adaptive lasso approach significantly outperforms existing Laplacian-based approaches, both statistical [17] and smoothness-based non-statistical [13].

**Notation:** Given  $\mathbf{A} \in \mathbb{R}^{p \times p}$ ,  $\text{tr}(\mathbf{A})$  denotes its trace. For  $\mathbf{B} \in \mathbb{R}^{p \times q}$ , we define its Frobenius norm and the vectorized  $\ell_1$  norm, respectively, as  $\|\mathbf{B}\|_F = \sqrt{\text{tr}(\mathbf{B}^\top \mathbf{B})}$  and  $\|\mathbf{B}\|_1 = \sum_{i,j} |B_{ij}|$  where  $B_{ij}$  is the  $(i, j)$ -th element of  $\mathbf{B}$ . We also denote  $B_{ij}$  by  $[\mathbf{B}]_{ij}$ . Given  $\mathbf{A} \in \mathbb{R}^{p \times p}$ ,  $\mathbf{A}^+ = \text{diag}(\mathbf{A})$  is a diagonal matrix with the same diagonal as  $\mathbf{A}$ ,  $\mathbf{A}^- = \mathbf{A} - \mathbf{A}^+$  is  $\mathbf{A}$  with all its diagonal elements set to zero, and  $\mathbf{A} \succ \mathbf{0}$  denotes that  $\mathbf{A}$  is positive-definite.

## II. SOME EXISTING APPROACHES

### A. Smoothness-Based Graph Learning [13]

With reference to (1) and (2), it follows (see [13]) that  $\text{tr}(\mathbf{X}^\top \mathbf{L} \mathbf{X}) = \frac{1}{2} \text{tr}(\mathbf{W} \hat{\mathbf{Z}})$  where  $\mathbf{W}, \hat{\mathbf{Z}} \in \mathbb{R}^{p \times p}$ ,  $\hat{Z}_{ij} = \|\mathbf{X}_i - \mathbf{X}_j\|_2^2$  and  $\mathbf{W}$  is the weight matrix (or the weighted adjacency matrix) with  $\mathbf{L} = \mathbf{D} - \mathbf{W}$ ,  $\mathbf{W} = \mathbf{W}^\top$ ,  $W_{ij} \geq 0$  and  $W_{ii} = 0$  for  $1 \leq i, j \leq p$ . Instead of performing a penalized minimization of  $\text{tr}(\mathbf{X}^\top \mathbf{L} \mathbf{X})$  to estimate  $\mathbf{L}$ , [13] minimizes a penalized  $\text{tr}(\mathbf{W} \hat{\mathbf{Z}})$  w.r.t.  $\mathbf{W}$  for graph learning. Given  $\mathbf{W}$ , one has unique  $\mathbf{L}$  and the edge-set  $\mathcal{E}$ . In the rest of the paper, we will scale  $\hat{\mathbf{Z}}$  as  $\hat{\mathbf{Z}}/n$  and denote the latter as  $\hat{\mathbf{Z}}$ .

Define the space  $\mathcal{W}_p$  of all valid  $p \times p$  weight matrices  $\mathbf{W}$

$$\mathcal{W}_p = \left\{ \mathbf{W} \in \mathbb{R}^{p \times p} : \mathbf{W} = \mathbf{W}^\top, W_{ij} \geq 0, W_{ii} = 0 \right\} \quad (3)$$

In [13] one looks for  $\min_{\mathbf{W} \in \mathcal{W}_p} f_s(\mathbf{W})$  where

$$f_s(\mathbf{W}) = \text{tr}(\mathbf{W} \hat{\mathbf{Z}}) + \frac{\beta}{2} \|\mathbf{W}\|_F^2 - \alpha \sum_{i=1}^p \ln \left( \sum_{j=1}^p W_{ij} \right) \quad (4)$$

with parameters  $\alpha > 0$  and  $\beta \geq 0$  controlling the “shape”. In (4),  $\text{tr}(\mathbf{W} \hat{\mathbf{Z}})$  is the main cost but minimizing it alone w.r.t.  $\mathbf{W}$  is ill-posed ( $\mathbf{W} = \mathbf{0}$  minimizes it). Using only the logarithmic barrier ( $\beta = 0$ ) leads to very sparse graphs, and changing  $\alpha$  only changes the scale of the solution. The term  $\frac{\beta}{2} \|\mathbf{W}\|_F^2$  controls graph sparsity. Note that the model of [13] unifies all prior smoothness-based models for estimation of Laplacian  $\mathbf{L}$

[15], [25]–[27]. A forward-backward algorithm based on [28] is given in [13] to optimize (4), where optimization is carried for fixed  $\alpha = 1$  and then one scales  $\mathbf{W}$  to obtain a desired  $\|\mathbf{W}\|$ ; a MATLAB implementation is in [29].

### B. Graphical Lasso: Penalized Log-Likelihood [7]

With  $\hat{\mathbf{S}}$  denoting the sample covariance (assume zero-mean:  $\hat{\mathbf{S}} = \frac{1}{n} \sum_{t=1}^n \mathbf{x}(t) \mathbf{x}^\top(t)$ ), seek  $\mathbf{\Omega}$  to yield  $\min_{\mathbf{\Omega} \succ \mathbf{0}} f_L(\mathbf{\Omega})$  where

$$f_L(\mathbf{\Omega}) = \text{tr}(\mathbf{\Omega} \hat{\mathbf{S}}) - \ln(|\mathbf{\Omega}|) + \lambda \|\mathbf{\Omega}^-\|_1, \quad (5)$$

$\lambda \|\mathbf{\Omega}^-\|_1$  is the lasso penalty and  $\lambda > 0$ . Unlike Laplacian  $\mathbf{L}$ , off-diagonal entries of  $\mathbf{\Omega}$  may not be non-positive.

### C. Generalized Graph Laplacian Estimation [17]

In [1], [16]–[18] approaches that make  $\mathbf{\Omega} = \mathbf{L}$  (Laplacian, or some regularized version) in (9) have been considered. In particular, [17] considers

$$\min_{\mathbf{\Theta} \succ \mathbf{0}} \text{tr}(\mathbf{\Theta} \hat{\mathbf{K}}) - \ln(|\mathbf{\Theta}|) \text{ where } \hat{\mathbf{K}} = \hat{\mathbf{S}} + \lambda(\mathbf{I} - \mathbf{1}_p \mathbf{1}_p^\top) \quad (6)$$

with  $\mathbf{\Theta}$  restricted to be a generalized graph Laplacian matrix. Software implementation of this algorithm is available in [30].

### D. Adaptive Lasso [31]

With  $\hat{\mathbf{\Omega}} = \arg \min_{\mathbf{\Omega} \succ \mathbf{0}} f_L(\mathbf{\Omega})$  from Sec. II-B, modify (5) as

$$\min_{\mathbf{\Omega} \succ \mathbf{0}} \text{tr}(\mathbf{\Omega} \hat{\mathbf{S}}) - \ln(|\mathbf{\Omega}|) + \lambda \sum_{i,j=1, i \neq j}^p \frac{\Omega_{ij}}{|\hat{\Omega}_{ij}|}, \quad (7)$$

i.e., use penalty varying with  $(i, j)$  as  $\lambda/|\hat{\Omega}_{ij}|$ .

## III. PROPOSED APPROACH

Define the space  $\mathcal{V}_p$  of all  $p \times p$  matrices  $\mathbf{V}$  that are symmetric with non-positive off-diagonal elements

$$\mathcal{V}_p = \left\{ \mathbf{V} \in \mathbb{R}^{p \times p} : \mathbf{V} = \mathbf{V}^\top, V_{ij} \leq 0, i \neq j \right\} \quad (8)$$

### A. Proposed Constrained Lasso: Penalized Log-Likelihood and Total Positivity

We propose to choose  $\mathbf{\Omega}$  as

$$\min_{\mathbf{\Omega} \succ \mathbf{0}, \mathbf{\Omega} \in \mathcal{V}_p} f_L(\mathbf{\Omega}). \quad (9)$$

We will use ADMM [32] after variable splitting to minimize  $f_L(\mathbf{\Omega})$ . Using variable splitting, consider

$$\min_{\substack{\mathbf{\Omega} \succ \mathbf{0} \\ \mathbf{V} \in \mathcal{V}_p}} \left\{ \text{tr}(\hat{\mathbf{\Sigma}} \mathbf{\Omega}) - \ln(|\mathbf{\Omega}|) + \lambda \|\mathbf{V}^-\|_1 \right\} \text{ subject to } \mathbf{\Omega} = \mathbf{V}. \quad (10)$$

The scaled augmented Lagrangian for this problem is [32]

$$L_\rho = \text{tr}(\hat{\mathbf{\Sigma}} \mathbf{\Omega}) - \ln(|\mathbf{\Omega}|) + \lambda \|\mathbf{V}^-\|_1 + \frac{\rho}{2} \|\mathbf{V} - \mathbf{\Omega} + \mathbf{U}\|_F^2 \quad (11)$$

where  $\mathbf{U}$  is the dual variable, and  $\rho > 0$  is the penalty parameter. Given the results  $\mathbf{\Omega}^{(k)}, \mathbf{V}^{(k)}, \mathbf{U}^{(k)}$  of the  $k$ th iteration, in the  $(k+1)$ st iteration, an ADMM algorithm executes the following three updates:

- (a)  $\mathbf{\Omega}^{(k+1)} \leftarrow \arg \min_{\mathbf{\Omega}} L_a(\mathbf{\Omega}), \quad L_a(\mathbf{\Omega}) := \text{tr}(\hat{\Sigma}\mathbf{\Omega}) - \ln(|\mathbf{\Omega}|) + \frac{\rho}{2} \|\mathbf{V}^{(k)} - \mathbf{\Omega} + \mathbf{U}^{(k)}\|_F^2$   
(b)  $\mathbf{V}^{(k+1)} \leftarrow \arg \min_{\mathbf{V} \in \mathcal{V}_p} L_b(\mathbf{V}), \quad L_b(\mathbf{V}) := \lambda \|\mathbf{V}^{-}\|_1 + \frac{\rho}{2} \|\mathbf{V} - \mathbf{\Omega}^{(k+1)} + \mathbf{U}^{(k)}\|_F^2$   
(c)  $\mathbf{U}^{(k+1)} \leftarrow \mathbf{U}^{(k)} + (\mathbf{V}^{(k+1)} - \mathbf{\Omega}^{(k+1)})$

Solution to update (a) follows from [32, Sec. 6.5] and is given in Algorithm 1.

In update (b) notice that  $L_b(\mathbf{V})$  is completely separable w.r.t. each element  $V_{ij}$ . Therefore, we solve  $V_{ij}^{(k+1)} \leftarrow \arg \min_{V_{ij} \leq 0, i \neq j} J_{ij}(V_{ij})$ , where

$$J_{ij}(V_{ij}) := \lambda |V_{ij}| \mathbf{1}_{i \neq j} + \frac{\rho}{2} (V_{ij} - [\mathbf{\Omega}^{(k+1)} - \mathbf{U}^{(k)}]_{ij})^2$$

We claim that the solution is given by

$$V_{ij}^{(k+1)} = \begin{cases} [\mathbf{\Omega}^{(k+1)} - \mathbf{U}^{(k)}]_{ii} & \text{if } i = j \\ S_{neg}([\mathbf{\Omega}^{(k+1)} - \mathbf{U}^{(k)}]_{ij}, \frac{\lambda}{\rho}) & \text{if } i \neq j \end{cases} \quad (12)$$

where, with  $(a)_+ := \max(0, a)$  and  $(a)_- := \min(0, a)$ ,

$$S_{neg}(a, \beta) := (1 - \beta/|a|)_+ a_-$$

denotes scalar soft thresholding for negative values of  $a$  and hard thresholding for  $a > 0$ . When  $i = j$ , we need to minimize only  $(V_{ij} - [\mathbf{\Omega}^{(k+1)} - \mathbf{U}^{(k)}]_{ij})^2$  w.r.t.  $V_{ii} = V_{ij}$ , thus the given solution follows. For constrained optimization under  $V_{ij} \leq 0$ , after setting  $A_{ij} = [\mathbf{\Omega}^{(k+1)} - \mathbf{U}^{(k)}]_{ij}$ , consider the Lagrangian  $L_v$

$$L_v = \lambda |V_{ij}| + \frac{\rho}{2} (V_{ij} - A_{ij})^2 + \nu V_{ij} \quad (13)$$

where  $\nu \geq 0$  is the Lagrange multiplier for the inequality constraint  $V_{ij} \leq 0$ . With  $v^*$  denoting an optimal solution, the KKT conditions for minimization are

$$0 \in \partial L_v = \lambda t + \rho(v^* - A_{ij}) + \nu \quad (14)$$

$$\nu v^* = 0 \quad (15)$$

$$\nu \geq 0 \quad (16)$$

$$v^* \leq 0 \quad (17)$$

where  $\partial L_v$  denotes the subdifferential of  $L_v$  at  $v^*$  and

$$t = \begin{cases} v^*/|v^*| & \text{if } v^* \neq 0 \\ \in \{u : |u| \leq 1, u \in \mathbb{R}\} & \text{if } v^* = 0. \end{cases} \quad (18)$$

When  $A_{ij} > 0$ , our claimed solution is  $v^* = 0$ . We need to check if  $\nu \geq 0$  and  $0 \in \partial L_v$  for some  $|t| \leq 1$ . The choice  $t = 0$  and  $\nu = \rho A_{ij} > 0$  satisfies the KKT conditions. When  $A_{ij} \leq 0$ , our claimed solution is the well-known soft-thresholding solution which satisfies the KKT conditions with  $\nu = 0$ . If  $|A_{ij}| \leq \rho/\lambda$ , then  $v^* = 0$  and  $t = \rho A_{ij}/\lambda$  satisfies the KKT conditions since  $|t| \leq 1$ . If  $|A_{ij}| > \rho/\lambda$ , then the given solution with  $t = A_{ij}/|A_{ij}|$  satisfies the KKT conditions. This proves that the solution (12) minimizes  $J_{ij}(V_{ij})$ .

A pseudocode for the ADMM algorithm used in this paper is given in Algorithm 1 where we use the stopping (convergence) criterion following [32, Sec. 3.3.1] and varying penalty parameter  $\rho$  following [32, Sec. 3.4.1]. For constrained lasso we take  $\lambda_{ij} = \lambda$  for all  $(i, j)$  in Algorithm 1.

## B. Proposed Constrained Adaptive Lasso

With  $\min_{\mathbf{\Omega} > \mathbf{0}, \mathbf{\Omega} \in \mathcal{V}_p} f_L(\mathbf{\Omega})$  from constrained lasso optimization, modify (9) as

$$\min_{\mathbf{\Omega} > \mathbf{0}, \mathbf{\Omega} \in \mathcal{V}_p} \text{tr}(\mathbf{\Omega}\hat{\mathbf{S}}) - \ln(|\mathbf{\Omega}|) + \lambda \sum_{i,j=1, i \neq j}^p \frac{\Omega_{ij}}{|\hat{\Omega}_{ij}|}, \quad (19)$$

i.e., use penalty varying with  $(i, j)$  as  $\lambda_{ij} = \lambda/|\hat{\Omega}_{ij}|$  where  $\hat{\Omega}_{ij}$  is obtained from proposed constrained lasso. Here we follow [31]. The solution given in Algorithm 1 applies.

---

### Algorithm 1 ADMM Algorithm for Constrained Lasso and Constrained Adaptive Lasso

---

**Input:** Number of samples  $n$ , number of nodes  $p$ , data  $\{\mathbf{x}(t)\}_{t=1}^n$ ,  $\mathbf{x} \in \mathbb{R}^p$ , regularization and penalty parameters  $\lambda_{ij}$  and  $\rho_0$ , tolerances  $\tau_{abs}$  and  $\tau_{rel}$ , variable penalty factor  $\mu$ , maximum number of iterations  $k_{max}$ .  $\lambda_{ij} = \lambda$  for lasso and  $\lambda_{ij} = \lambda/\hat{\Omega}_{ij}$  for adaptive lasso where  $\hat{\Omega}_{ij}$  is the result of lasso.

**Output:** estimated inverse covariance  $\hat{\mathbf{\Omega}}$  and edge-set  $\hat{\mathcal{E}}$

- 1: Calculate sample covariance  $\hat{\mathbf{S}} = \frac{1}{n} \sum_{t=1}^n \mathbf{x}(t)\mathbf{x}^\top(t)$  (after centering  $\mathbf{x}(t)$ ).
- 2: Initialize:  $\mathbf{U}^{(0)} = \mathbf{V}^{(0)} = \mathbf{0}$ ,  $\mathbf{\Omega}^{(0)} = (\text{diag}(\hat{\mathbf{S}}))^{-1}$ , where  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{(p) \times (p)}$ ,  $\rho^{(0)} = \rho_0$
- 3: converged = FALSE,  $k = 0$
- 4: **while** converged = FALSE AND  $k \leq k_{max}$ , **do**
- 5: Eigen-decompose  $\hat{\mathbf{S}} - \rho^{(k)}(\mathbf{V}^{(k)} + \mathbf{U}^{(k)})$  as  $\hat{\mathbf{S}} - \rho^{(k)}(\mathbf{V}^{(k)} + \mathbf{U}^{(k)}) = \mathbf{Q}\mathbf{D}\mathbf{Q}^\top$  with diagonal matrix  $\mathbf{D}$  consisting of eigenvalues. Define diagonal matrix  $\hat{\mathbf{D}}$  with  $\ell$ th diagonal element  $\hat{D}_{\ell\ell} = (-D_{\ell\ell} + \sqrt{D_{\ell\ell}^2 + 4\rho^{(k)}})/(2\rho^{(k)})$ . Set  $\mathbf{\Omega}^{(k+1)} = \mathbf{Q}\hat{\mathbf{D}}\mathbf{Q}^\top$ .
- 6: Define thresholding operator  $S_{neg}(a, \beta) := (1 - \beta/|a|)_+ a_-$  where  $(a)_+ := \max(0, a)$  and  $(a)_- := \min(0, a)$ . The  $(i, j)$ th element of  $\mathbf{V}$  is updated as
$$V_{ij}^{(k+1)} = \begin{cases} [\mathbf{\Omega}^{(k+1)} - \mathbf{U}^{(k)}]_{ii} & \text{if } i = j \\ S_{neg}([\mathbf{\Omega}^{(k+1)} - \mathbf{U}^{(k)}]_{ij}, \frac{\lambda_{ij}}{\rho}) & \text{if } i \neq j \end{cases}$$
- 7: Dual update  $\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} + (\mathbf{V}^{(k+1)} - \mathbf{\Omega}^{(k+1)})$ .
- 8: Check convergence. Set tolerances

$$\tau_{pri} = p\tau_{abs} + \tau_{rel} \max(\|\mathbf{\Omega}^{(k+1)}\|_F, \|\mathbf{V}^{(k+1)}\|_F)$$

$$\tau_{dual} = p\tau_{abs} + \tau_{rel} \|\mathbf{U}^{(k+1)}\|_F / \rho^{(k)}$$

Define  $d_p = \|\mathbf{\Omega}^{(k+1)} - \mathbf{V}^{(k+1)}\|_F$  and  $d_d = \rho^{(k)} \|\mathbf{V}^{(k+1)} - \mathbf{V}^{(k)}\|_F$ . If  $(d_p \leq \tau_{pri})$  AND  $(d_d \leq \tau_{dual})$ , set converged = TRUE.

- 9: Update penalty parameter  $\rho$  :

$$\rho^{(k+1)} = \begin{cases} 2\rho^{(k)} & \text{if } d_p > \mu d_d \\ \rho^{(k)}/2 & \text{if } d_d > \mu d_p \\ \rho^{(k)} & \text{otherwise.} \end{cases}$$

We also need to set  $\mathbf{U}^{(k+1)} = \mathbf{U}^{(k+1)}/2$  for  $d_p > \mu d_d$  and  $\mathbf{U}^{(k+1)} = 2\mathbf{U}^{(k+1)}$  for  $d_d > \mu d_p$ .

---

Model:	Chain Graph: number of nodes $p=100$				
sample size $n$	50	100	200	400	2000
Approach	$F_1$ score ( $\pm\sigma$ )				
Const. Lasso	0.4058 $\pm$ 0.0115	0.4353 $\pm$ 0.0113	0.4630 $\pm$ 0.0107	0.4921 $\pm$ 0.0095	0.5681 $\pm$ 0.0102
Kalofolias [13]	0.2631 $\pm$ 0.0919	0.2656 $\pm$ 0.0997	0.2753 $\pm$ 0.1049	0.2770 $\pm$ 0.1068	0.2941 $\pm$ 0.1139
GGL [17]	0.6727 $\pm$ 0.0218	0.6753 $\pm$ 0.0204	0.6715 $\pm$ 0.0192	0.6760 $\pm$ 0.0217	0.6758 $\pm$ 0.0215
Const. Adap. Lasso	<b>0.9851</b> $\pm$ 0.0089	<b>0.9977</b> $\pm$ 0.0035	<b>0.9999</b> $\pm$ 0.0007	<b>1.000</b> $\pm$ 0.0000	<b>1.0000</b> $\pm$ 0.0000
	Frobenius Error Norm ( $\pm\sigma$ )				
Const. Lasso	0.4752 $\pm$ 0.0116	0.4041 $\pm$ 0.0103	0.3405 $\pm$ 0.0084	0.2785 $\pm$ 0.0077	0.1620 $\pm$ 0.0053
Kalofolias [13]	0.5814 $\pm$ 0.0211	0.5564 $\pm$ 0.0206	0.5445 $\pm$ 0.0192	0.5357 $\pm$ 0.0196	0.5262 $\pm$ 0.0180
GGL [17]	0.2599 $\pm$ 0.0216	0.1799 $\pm$ 0.0135	0.1252 $\pm$ 0.0099	0.0881 $\pm$ 0.0066	0.0387 $\pm$ 0.0031
Const. Adap Lasso	0.2390 $\pm$ 0.0210	0.1623 $\pm$ 0.0126	0.1113 $\pm$ 0.0089	0.0785 $\pm$ 0.0066	0.0345 $\pm$ 0.0030
	Time(s) ( $\pm\sigma$ )				
Const. Lasso	2.4289 $\pm$ 0.3891	1.9707 $\pm$ 0.0868	1.5710 $\pm$ 0.0601	1.2978 $\pm$ 0.0441	0.8413 $\pm$ 0.0286
Kalofolias [13]	0.2706 $\pm$ 0.0070	0.2706 $\pm$ 0.0172	0.2789 $\pm$ 0.0763	0.2652 $\pm$ 0.0030	0.2662 $\pm$ 0.0019
GGL [17]	0.0668 $\pm$ 0.0055	0.0688 $\pm$ 0.0086	0.0682 $\pm$ 0.0051	0.0687 $\pm$ 0.0064	0.0679 $\pm$ 0.0022
Const. Adap Lasso	8.1710 $\pm$ 0.8587	5.7741 $\pm$ 0.2202	4.8850 $\pm$ 1.2391	3.4469 $\pm$ 0.3290	1.5416 $\pm$ 0.1274
Model:	Erdős-Rényi Graph: number of nodes $p=100$				
	$F_1$ score ( $\pm\sigma$ )				
Const. Lasso	0.2363 $\pm$ 0.0135	0.2462 $\pm$ 0.0141	0.2585 $\pm$ 0.0174	0.2745 $\pm$ 0.0152	0.3119 $\pm$ 0.0205
Kalofolias [13]	0.1363 $\pm$ 0.0398	0.1420 $\pm$ 0.0397	0.1491 $\pm$ 0.0399	0.1666 $\pm$ 0.0495	0.1586 $\pm$ 0.0592
GGL [17]	0.4609 $\pm$ 0.0232	0.4901 $\pm$ 0.0232	0.4984 $\pm$ 0.0197	0.5032 $\pm$ 0.0175	0.4927 $\pm$ 0.0218
Const. Adap Lasso	<b>0.7165</b> $\pm$ 0.0575	<b>0.8756</b> $\pm$ 0.0391	<b>0.9584</b> $\pm$ 0.0183	<b>0.9859</b> $\pm$ 0.0087	<b>0.9991</b> $\pm$ 0.0022
	Frobenius Error Norm ( $\pm\sigma$ )				
Const. Lasso	0.9964 $\pm$ 0.0043	0.9898 $\pm$ 0.0098	0.6367 $\pm$ 0.0914	0.4743 $\pm$ 0.0378	0.3720 $\pm$ 0.0212
Kalofolias [13]	0.7621 $\pm$ 0.0316	0.7028 $\pm$ 0.0236	0.6756 $\pm$ 0.0224	0.6562 $\pm$ 0.0200	0.6475 $\pm$ 0.0166
GGL [17]	0.6062 $\pm$ 0.0615	0.4204 $\pm$ 0.0492	0.2866 $\pm$ 0.0317	0.1963 $\pm$ 0.0234	0.0848 $\pm$ 0.0117
Const. Adap Lasso	0.6309 $\pm$ 0.0609	0.4454 $\pm$ 0.0501	0.3119 $\pm$ 0.0379	0.2265 $\pm$ 0.0301	0.1039 $\pm$ 0.0186
	Time(s) ( $\pm\sigma$ )				
Const. Lasso	3.4578 $\pm$ 0.6695	2.6559 $\pm$ 0.2382	2.1640 $\pm$ 0.1588	1.8420 $\pm$ 0.0908	1.2341 $\pm$ 0.0584
Kalofolias [13]	0.2587 $\pm$ 0.0068	0.2570 $\pm$ 0.0031	0.2568 $\pm$ 0.0023	0.2557 $\pm$ 0.0024	0.2582 $\pm$ 0.0021
GGL [17]	0.0652 $\pm$ 0.0047	0.0639 $\pm$ 0.0013	0.0637 $\pm$ 0.0012	0.0636 $\pm$ 0.0015	0.0633 $\pm$ 0.0014
Const. Adap Lasso	9.3087 $\pm$ 0.3813	8.5631 $\pm$ 0.4493	7.0382 $\pm$ 0.4599	5.3631 $\pm$ 0.3213	2.8511 $\pm$ 0.2075

TABLE I: Results for Chain and Erdos-Renyi graphs. Frobenius error norm is  $\|c\hat{\Omega}^- - \Omega_0^-\|_F / \|\Omega_0^-\|_F$  where scalar  $c$  is picked to minimize  $\|c\hat{\Omega}^- - \Omega_0^-\|_F$ . “Const. Lasso” stands for the proposed constrained lasso approach that enforces total positivity, and “Const. Adap lasso” denotes its adaptive lasso version.

---

```

10:    $k \leftarrow k + 1$ 
11: end while
12: For  $i \neq j$ , if  $|V_{ij}| > 0$ , assign edge  $\{i, j\} \in \hat{\mathcal{E}}$ , else
     $\{i, j\} \notin \hat{\mathcal{E}}$ . Inverse covariance estimate is  $\hat{\Omega}$ .

```

---

#### IV. SIMULATION EXAMPLES

We consider Gaussian graphical models based on two graphs: a **chain graph** where  $p$  nodes are connected in succession, and an **Erdős-Rényi** graph where nodes are connected with probability  $p_{er} = 0.03$ . In each model, in the upper triangular  $\Omega$  (inverse covariance),  $\Omega_{ij} = 0$  if  $\{i, j\} \notin \mathcal{E}$ , and  $\Omega_{ij}$  is uniformly distributed over  $[-0.3, -0.1]$  if  $\{i, j\} \in \mathcal{E}$ . With  $\Omega = \Omega^\top$ , we take  $\Omega_{ii} = -\sum_{j=1}^p \Omega_{ij}$  for every  $i$ , yielding the combinatorial Laplacian matrix  $\mathbf{L} = \Omega$ . Now add  $\kappa \mathbf{I}$  to  $\Omega$  with  $\kappa$  picked to make minimum eigenvalue of  $\Omega + \kappa \mathbf{I}$  equal to 0.001, and with  $\Phi \Phi^\top = (\Omega + \kappa \mathbf{I})^{-1}$ , we generate  $\mathbf{x} = \Phi \mathbf{w}$  with  $\mathbf{w} \in \mathbb{R}^p$  as Gaussian  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . We generate  $n$  i.i.d. observations from  $\mathbf{x}$  using  $p = 100$ . Addition of  $\kappa \mathbf{I}$  yields a generalized Laplacian matrix  $\mathbf{L} = \Omega + \kappa \mathbf{I}$  [17]. Given choice of  $\kappa$  leads to a precision matrix (and corresponding covariance matrix) with (matrix inversion) condition number  $> 1000$  for both models, making it a challenging problem.

We apply four methods for estimating the true edgeset  $\mathcal{E}_0$  and true (off-diagonal) inverse covariance  $\Omega_0^-$ . One of the methods is the signal smoothness-based method of [13] which yields the weighted adjacency matrix  $\mathbf{W}$ , equaling  $-\Omega_0^-$  (off-diagonal  $\Omega_0^-$ ) under the Laplacian assumption, hence one of our performance criterion is error in estimating  $\Omega_0^-$ . The chosen methods are

- (1) Proposed constrained lasso approach, labeled “Const. Lasso” in Table I. We use Algorithm 1 with constant  $\lambda$  for all  $(i, j)$ .
- (2) Smoothness-based graph learning [13] solved via the forward-backward algorithm available in [29] (MATLAB function `gsp_learn_graph_log_degrees.m`), labeled “smooth [13]” in Table I. It requires one to set small values in estimated  $\mathbf{W}$  to be set to zero; following [13], [29], all  $\hat{W}_{ij} \leq 10^{-4}$  are set to zero.
- (3) Generalized graph Laplacian (GGL) method of [17], using MATLAB function `estimate_ggl.m` from [30], labeled “GGL [17]” in Table I.
- (4) Proposed constrained adaptive lasso approach, labeled “Const. Adap Lasso” in Table I. We use Algorithm 1 with edge-dependent  $\lambda$ ,  $\lambda_{ij} = \lambda / |\hat{\Omega}_{ij}|$ , where  $\hat{\Omega}_{ij}$  is obtained from proposed constrained lasso.

A performance measure is  $F_1$ -score for efficacy in edge detection. The  $F_1$ -score is defined as  $F_1 = 2 \times \text{precision} \times \text{recall} / (\text{precision} + \text{recall})$  where  $\text{precision} = |\hat{\mathcal{E}} \cap \mathcal{E}_0| / |\hat{\mathcal{E}}|$ ,  $\text{recall} = |\hat{\mathcal{E}} \cap \mathcal{E}_0| / |\mathcal{E}_0|$ , and  $\mathcal{E}_0$  and  $\hat{\mathcal{E}}$  denote the true and estimated edge sets, respectively. Table I shows the simulation results where the run time in seconds was calculated via MATLAB tic-toc functions on a Window 10 operating system with processor Intel(R) Core(TM) i5-6400T CPU @2.20 GHz with 12 GB RAM. For each of the four schemes, tuning parameter  $\lambda$  ( $\beta$  for [13]) was picked for  $n = 200$  via simulations to maximize  $F_1$ -score, then  $\lambda$  was scaled as  $\propto \sqrt{\ln(p)/n}$  [8], [17] while  $\beta$  was kept fixed. The performance measures are  $F_1$ -score for efficacy in edge detection, and normalized Frobenius error norm in estimating  $\Omega_0^-$  (off-diagonal true  $\Omega_0$ ), defined as  $\|c\hat{\Omega}_0^- - \Omega_0^-\|_F / \|\Omega_0^-\|_F$  where  $c$  is selected as follows. We scale estimated  $\hat{\Omega}_0^-$  (when signal-smoothing is used, set  $\hat{\Omega}_0^- = -\hat{W}_0^-$ ), by a scalar  $c$  chosen to minimize mean-square error  $\|\Omega_0^- - c\hat{\Omega}_0^-\|_F^2$ , resulting in  $c = \text{tr}(\Omega_0^- \hat{\Omega}_0^-) / \text{tr}(\hat{\Omega}_0^- \hat{\Omega}_0^-)$ . In practice,  $\Omega_0^-$  is unknown. The above scaling preserves relative weighting among  $\Omega_{ij}$ 's which is what is relevant in applications [13] and is available without knowing  $\Omega_0^-$ .

The condition number  $> 1000$  of covariance and precision matrices for both models makes it a challenging problem. We see that only constrained adaptive lasso performs well, with high  $F_1$  scores.

## V. CONCLUSIONS

We considered the problem of estimating the structure of an undirected weighted sparse graphical model of multivariate data under the assumption of total positivity where all partial correlations are non-negative. We investigated modifications to widely used penalized log-likelihood approaches to enforce total positivity but not the Laplacian structure. An ADMM algorithm was presented for constrained optimization under total positivity and lasso as well as adaptive lasso penalties. Numerical results show that the proposed constrained adaptive lasso approach significantly outperforms existing Laplacian-based approaches.

## REFERENCES

- [1] X. Dong, D. Thanou, M. Rabbat and P. Frossard, "Learning graphs from data," *IEEE Signal Process. Mag.*, pp. 44-63, May 2019.
- [2] S.L. Lauritzen, *Graphical models*. Oxford, UK: Oxford Univ. Press, 1996.
- [3] J. Whittaker, *Graphical Models in Applied Multivariate Statistics*. New York: Wiley, 1990.
- [4] P. Danaher, P. Wang and D.M. Witten, "The joint graphical lasso for inverse covariance estimation across multiple classes," *J. Royal Statistical Society, Series B*, vol. 76, pp. 373-397, 2014.
- [5] N. Meinshausen and P. Bühlmann, "High-dimensional graphs and variable selection with the Lasso," *Ann. Statist.*, vol. 34, no. 3, pp. 1436-1462, 2006.
- [6] K. Mohan, P. London, M. Fazel, D. Witten and S.I. Lee, "Node-based learning of multiple Gaussian graphical models," *J. Machine Learning Research*, vol. 15, 2014.
- [7] J. Friedman, T. Hastie and R. Tibshirani, "Sparse inverse covariance estimation with the graphical lasso," *Biostatistics*, vol. 9, no. 3, pp. 432-441, July 2008.
- [8] O. Banerjee, L.E. Ghaoui and A. d'Aspremont, "Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data," *J. Machine Learning Res.*, vol. 9, pp. 485-516, 2008.

- [9] Y. Wang, U. Roy and C. Uhler, "Learning high-dimensional Gaussian graphical models under total positivity without adjustment of tuning parameters," in *Proc. 23rd Intern. Conf. Artificial Intelligence & Statistics (AISTATS)*, Palermo, Italy, 2020.
- [10] S. Lauritzen, C. Uhler and P. Zwiernik, "Maximum likelihood estimation in Gaussian models under total positivity," *Annals of Statistics*, vol. 47, pp. 1835-1863, 2019.
- [11] M. Slawski and M. Hein, "Estimation of positive definite M-matrices and structure learning for attractive Gaussian Markov random fields," *Linear Algebra and its Applications*, vol. 473, pp. 145-179, 2015.
- [12] S. Karlin and Y. Rinott, "M-matrices as covariance matrices of multinormal distributions," *Linear Algebra and its Applications*, vol. 52, pp. 419-438, 1983.
- [13] V. Kalofolias, "How to learn a graph from smooth signals," in *Proc. 19th Intern. Conf. Artificial Intelligence & Statistics (AISTATS)*, Cadiz, Spain, 2016.
- [14] V. Kalofolias and N. Perraudin, "Large scale graph learning from smooth signals," in *7th Intern. Conf. Learning Representations (ICLR 2019)*, New Orleans, LA, USA, May 6-9, 2019.
- [15] X. Dong, D. Thanou, P. Frossard and P. Vandergheynst "Learning Laplacian matrix in smooth graph signal representations," *IEEE Trans. Signal Process.*, vol. 64, no. 23, pp. 6160-6173, Dec. 1, 2016.
- [16] E. Pavez and A. Ortega, "Generalized Laplacian precision matrix estimation for graph signal processing," in *Proc. IEEE ICASSP 2016*, Shanghai, China, March 2016, pp. 6350-6354.
- [17] H.E. Egilmez, E. Pavez and A. Ortega, "Graph learning from data under Laplacian and structural constraints," *IEEE J. Sel. Topics Signal Process.*, vol. 11, no. 6, pp. 825-841, Sept. 2017.
- [18] E. Pavez, H.E. Egilmez and A. Ortega, "Learning graphs with monotone topology properties and multiple connected components," *IEEE Trans. Signal Process.*, vol. 66, no. 9, pp. 2399-2413, May 1, 2018.
- [19] M. Belkin and P. Niyogi, "Laplacian eigenmaps and spectral techniques for embedding and clustering," in *Proc. NIPS*, vol. 14, pp. 585-591, 2001.
- [20] M. Belkin, P. Niyogi and V. Sindhwani "Manifold regularization: A geometric framework for learning from labeled and unlabeled examples," *J. Machine Learning Research*, vol. 7, pp. 2399-2434, 2006.
- [21] Y. Bengio, O. Delalleau and N. Le Roux, "Label propagation and quadratic criterion," in *Semi-Supervised Learning*, O. Chapelle, B. Schölkopf and A. Zien (Eds), pp. 193-216, MIT Press, 2006.
- [22] U.V. Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, vol. 17, no. 4, pp. 395-416, 2007
- [23] X. Zhu, Z. Ghahramani and J. Lafferty, "Semi-supervised learning using Gaussian fields and harmonic functions," in *Proc. 20th Intern. Conf. Machine Learning (ICML)*, vol. 3, pp. 912-919, 2003.
- [24] X. Zhou and M. Belkin, "Semi-supervised learning by higher order regularization," in *Proc. 14th Intern. Conf. Artificial Intelligence & Statistics (AISTATS)*, Fort Lauderdale, FL, 2011.
- [25] F. Wang and C. Zhang, "Label propagation through linear neighborhoods," *IEEE Trans. Knowledge and Data Engineering*, vol. 20, no. 1, pp. 55-67, 2008.
- [26] S.T. Daitch, J.A. Kelner and D.A. Spielman, "Fitting a graph to vector data," in *Proc. 26th Intern. Conf. Machine Learning (ICML)*, pp. 201-208, Montreal, Canada, 2009.
- [27] B. Lake and J. Tenenbaum, "Discovering structure by learning sparse graphs," *Proc. 32nd Annual Meeting of the Cognitive Science Society (CogSci 2010)*, Portland, Oregon, Aug. 2010, pp. 778-784.
- [28] N. Komodakis and J.C. Pesquet, "Playing with duality: An overview of recent primal-dual approaches for solving large-scale optimization problems," *IEEE Signal Processing Mag.*, vol. 32, pp. 31-54, 2015.
- [29] N. Perraudin, J. Paratte, D. Shuman, V. Kalofolias, P. Vandergheynst and D.K. Hammond, "GSPBOX: A toolbox for signal processing on graphs," *arXiv:1408.5781v2 [cs.IT]*, 15 March 2016.
- [30] H.E. Egilmez, E. Pavez and A. Ortega, "GLL: Graph Laplacian learning package, version 1.0," [Online]. Available: [https://github.com/STACUSC/Graph\\_Learning](https://github.com/STACUSC/Graph_Learning), 2017.
- [31] H. Zou, "The adaptive lasso and its oracle properties," *J. American Statistical Assoc.*, vol. 101, pp. 1418-1429, 2006.
- [32] S. Boyd, N. Parikh, E. Chu, B. Peleato and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1-122, 2010.