GRAPH LEARNING FROM MULTI-ATTRIBUTE SMOOTH SIGNALS

Jitendra K. Tugnait

Department of Electrical & Computer Engineering Auburn University, Auburn, AL 36849, USA tugnajk@auburn.edu

ABSTRACT

We consider the problem of estimating the structure of an undirected weighted graph underlying a set of smooth multi-attribute signals. Most existing methods for graph estimation are based on single-attribute models where one associates a scalar data variable with each node of the graph, and the problem is to infer the graph topology that captures the relationships between these variables. An example is image graphs for grayscale texture images for modeling dependence of a pixel on neighboring pixels. In multi-attribute graphical models, each node represents a vector, as for example, in color image graphs where one has three variables (RGB color components) per pixel node. In this paper, we extend the single attribute approach of Kalofolias (2016) to multi-attribute data. An alternating direction method of multipliers (ADMM) algorithm is presented to optimize the objective function to infer the graph topology. Numerical results based on synthetic as well as real data are presented.

Keywords: Sparse graph learning; graph estimation; graph Laplacian; undirected graph; multi-attribute data.

1. INTRODUCTION

An undirected simple weighted graph is denoted $\mathcal{G}=(V,\mathcal{E},\mathbf{W})$ where $V=\{1,2,\cdots,p\}=[p]$ is the set of p nodes, $\mathcal{E}\subseteq[p]\times[p]$ is the set of undirected edges, and $\mathbf{W}\in\mathbb{R}^{p\times p}$ stores the nonnegative 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_1,\cdots,x_p , ($\mathbf{x}=[x_1\ x_2\ \cdots\ x_p]^\top$), 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_1, \cdots, 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_ℓ , $\ell \in [p]$, $\ell \neq i$, $\ell \neq j$. In particular, Gaussian graphical models (GGMs) are CIGs where x is multivariate Gaussian. Suppose x has positive-definite covariance matrix Σ with inverse

covariance matrix $\Omega = \Sigma^{-1}$. Then Ω_{ij} , the (i,j)-th element of Ω , is zero iff x_i and x_j are conditionally independent. Such models for x have been extensively studied. Given n samples of x, in high-dimensional settings where $p \gg 1$ and/or n is of the order of p, one estimates Ω under some sparsity constraints; see [4–8]. In these graphs each node represents a scalar random variable. In many applications, there may be more than one random variable associated with a node. This class of graphical models has been called multi-attribute graphical models in [9–11]. In [12,13] image graphs for grayscale texture images are inferred for modeling dependence of a pixel on neighboring pixels; here one has one variable per pixel node. These approaches do not apply to color images where one has three variables (RGB color components) per pixel node. Image graphs for color images is an example of multi-attribute graphical models.

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 [1,14,15]. Suppose we are given n samples $\{x(t)\}_{t=1}^n$ of the p data variables x_1, x_1, \cdots, x_p , with $x(t) = [x_1(t) \ x_2(t) \ \cdots \ x_p(t)]^{\top}$. Define the $p \times n$ matrix

$$\boldsymbol{X} = \begin{bmatrix} x_1(1) & x_1(2) & \cdots & x_1(n) \\ x_2(1) & x_2(2) & \cdots & x_2(n) \\ \vdots & \vdots & \ddots & \vdots \\ x_p(1) & x_p(2) & \cdots & x_p(n) \end{bmatrix} . \tag{1}$$

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

$$\frac{1}{2} \sum_{i,j=1}^{p} W_{ij} \| \boldsymbol{X}_{i.} - \boldsymbol{X}_{j.} \|_{2}^{2} = \text{tr}(\boldsymbol{X}^{\top} \boldsymbol{L} \boldsymbol{X})$$
 (2)

where X_i denotes the ith row of X, L = D - W is the (combinatorial) graph Laplacian (matrix), and D is the diagonal weighted degree matrix with $D_{ii} = \sum_{j=1}^p W_{ij}$. In the words of [14], if two vectors X_i and X_j from a smooth set of signals are associated with two well-connected nodes i and j (i.e., W_{ij} is large), they are expected to have a small distance $\|X_i - X_j\|_2$ so that $\operatorname{tr}(X^\top LX)$ is small. In particular, if X_i does not vary with i, $\operatorname{tr}(X^\top LX) = 0$. Based on the smoothness constraint, graph learning from data X becomes equivalent to estimation of the graph Laplacian matrix L [1, 14]. The approaches given in [1, 14, 15] apply only to single attribute models. In this paper our objective is learn graphs from multi-attribute data.

Graph Laplacian matrix has been extensively used for embedding, manifold learning, clustering and semi-supervised learn-

This work was supported by NSF Grants CCF-1617610 and ECCS-2040536.

ing [16, 17]; see [1, 14] for further references to applications to web page categorization with graph information, graph regularized sparse coding, and matrix completion. Graph-based transform coding and its potential application to signal/image compression is discussed in [18] where learning of the graph Laplacian matrix plays a key role.

Another set of approaches are based on statistical considerations under the graph Laplacian constraint [1, 12, 13] where Laplacian L, after regularization, plays the role of inverse covariance Ω ; L is a symmetric, nonnegative-definite matrix but with non-positive off-diagonal entries. These approaches too apply only to single attribute models (and when off-diagonal entries of inverse covariance are non-positive).

In this paper, we extend the single attribute approach of Kalofolias (2016) to multi-attribute data. An alternating direction method of multipliers (ADMM) algorithm is presented to optimize the objective function to infer the graph topology. Numerical results based on synthetic as well as real data are presented.

Notation: Given a matrix $A \in \mathbb{R}^{p \times p}$, $\operatorname{tr}(A)$ denotes its trace. For a matrix $B \in \mathbb{R}^{p \times q}$, we define its Frobenius norm as $\|B\|_F = \sqrt{\operatorname{tr}(B^\top B)}$, and denote its (i,j)-th element by B_{ij} . We also denote B_{ij} by $[B]_{ij}$. Also, B_i . and $B_{.j}$ denote column vectors comprising ith row and jth column, respectively, of B. Notation $\mathbf{1}_p$ denotes a column of p ones, $B^+ = \max(B,0)$, $B^- = \max(-B,0)$ and $\operatorname{abs}(B) = B^+ + B^-$, where "max" operation is elementwise.

2. BACKGROUND AND SYSTEM MODELS

2.1. Statistical Models

The conditional dependency structure among p random variables x_1, x_1, \cdots, x_p , components of $\boldsymbol{x} \in \mathbb{R}^p$, is represented using an undirected graph $\mathcal{G} = (V, \mathcal{E})$, where V = [1, p] is the set of p nodes corresponding to the p random variables x_i s, and $\mathcal{E} \subseteq V \times V$ is the set of undirected edges describing conditional dependencies among x_i s. There is no edge between nodes i and j iff x_i and x_j are conditionally independent given the remaining p-2 variables x_ℓ , $\ell \in [1, p]$, $\ell \neq i$, $\ell \neq j$ [3, p. 60]. We will call \mathcal{G} a single-attribute graphical model for \boldsymbol{x} .

Now consider p jointly Gaussian random vectors $\mathbf{z}_i \in \mathbb{R}^m$, $i=1,2,\cdots,p$. We associate \mathbf{z}_i with the ith node of an undirected graph $\mathcal{G}=(V,\mathcal{E})$ where V=[1,p] is the set of p nodes, and $\mathcal{E}\subseteq V\times V$ is the set of edges that describe the conditional dependencies among vectors $\{\mathbf{z}_i,\ i\in V\}$. As in the scalar case (m=1), there is no edge between node i and node j in \mathcal{G} (i.e., $\{i,j\}\not\in\mathcal{E}$) iff random vectors \mathbf{z}_i and \mathbf{z}_j are conditionally independent given all the remaining random vectors \mathbf{z}_ℓ corresponding to the remaining p-2 nodes in V, i.e., for $\ell\in V\setminus\{j,k\}$ [10,11]. This is the **multi-attribute Gaussian graphical model**. The term multi-attribute Gaussian graphical model has been used in [9] for such models. In this paper, we will use a similar set-up under a smoothness constraint, but without statistical considerations.

Define the mp-vector

$$\boldsymbol{x} = [\boldsymbol{z}_1^{\top} \ \boldsymbol{z}_2^{\top} \ \cdots \ \boldsymbol{z}_p^{\top}]^{\top} \in \mathbb{R}^{mp}$$
. (3)

Suppose we have n i.i.d. observations $\boldsymbol{x}(t), t = 1, 2, \cdots, n$, of \boldsymbol{x} . The objective in statistical multi-attribute graph learning is to estimate the inverse covariance matrix $(\mathbb{E}\{\mathbf{x}\mathbf{x}^{\top}\})^{-1}$ and to determine if $\{i,j\} \in \mathcal{E}$, given data $\{\boldsymbol{x}(t)\}_{t=1}^n$.

Let us associate \boldsymbol{x} with an "enlarged" graph $\bar{\mathcal{G}}=(\bar{V},\bar{\mathcal{E}})$, where $\bar{V}=[1,mp]$ and $\bar{\mathcal{E}}\subseteq\bar{V}\times\bar{V}$. Now $[\boldsymbol{z}_j]_\ell$, the ℓ th component of \boldsymbol{z}_j associated with node j of $\mathcal{G}=(V,\mathcal{E})$, is the random variable

 $x_q = [x]_q$, where $q = (j-1)m + \ell$, $j = 1, 2, \cdots, p$ and $\ell = 1, 2, \cdots, m$. The random variable x_q is associated with node q of $\bar{\mathcal{G}} = (\bar{V}, \bar{\mathcal{E}})$. Corresponding to the edge $\{j, k\} \in \mathcal{E}$ in the multi-attribute $\mathcal{G} = (V, \mathcal{E})$, there are m^2 edges $\{q, r\} \in \bar{\mathcal{E}}$ specified by q = (j-1)m+s and r = (k-1)m+t, where $s = 1, 2, \cdots, m$ and $t = 1, 2, \cdots, m$. The graph $\bar{\mathcal{G}} = (\bar{V}, \bar{\mathcal{E}})$ is a single-attribute graph. In order for $\bar{\mathcal{G}}$ to reflect the conditional independencies encoded in \mathcal{G} , we must have the equivalence

$$\begin{aligned} \{j,k\} \not\in \mathcal{E} &\Leftrightarrow \bar{\mathcal{E}}^{(jk)} \cap \bar{\mathcal{E}} = \emptyset \text{, where} \\ \bar{\mathcal{E}}^{(jk)} &= & \{\{q,r\} \mid q = (j-1)m + s, \, r = (k-1)m + t, \\ s,t &= 1,2,\cdots,m \} \text{.} \end{aligned} \tag{5}$$

Let $R_{xx}=\mathbb{E}\{xx^{\top}\}\succ \mathbf{0}$ and $\Omega=R_{xx}^{-1}$. Define the (j,k)th $m\times m$ subblock $\Omega^{(jk)}$ of Ω as

$$[\mathbf{\Omega}^{(jk)}]_{st} = [\mathbf{\Omega}]_{(j-1)m+s,(k-1)m+t}, \ s,t=1,2,\cdots,m.$$
 (6)

It is established in [11, Sec. 2.1] that

$$\mathbf{\Omega}^{(jk)} = \mathbf{0} \Leftrightarrow \mathbf{z}_j \text{ and } \mathbf{z}_k \text{ are conditionally independent}$$

$$\Leftrightarrow \{j, k\} \notin \mathcal{E}. \tag{7}$$

Since $\Omega^{(jk)}=\mathbf{0}$ is equivalent to $[\Omega]_{qr}=0$ for every $\{q,r\}\in \bar{\mathcal{E}}^{(jk)}$, and since, by [2, Proposition 5.2], $[\Omega]_{qr}=0$ iff x_q and x_r are conditionally independent, hence, iff $\{q,r\}\not\in \bar{\mathcal{E}}$, it follows that equivalence (4) holds true.

2.2. Smoothness-Based Graph Learning: Single-Attribute Mod-

Here we review [14]. With reference to (1) and (2), it is established in [14] that

$$\operatorname{tr}(\boldsymbol{X}^{\top} \boldsymbol{L} \boldsymbol{X}) = \frac{1}{2} \operatorname{tr}(\boldsymbol{W} \boldsymbol{Z}) \tag{8}$$

where
$$\boldsymbol{W}, \boldsymbol{Z} \in \mathbb{R}^{p \times p}, \ Z_{ij} = \|\boldsymbol{X}_{i.} - \boldsymbol{X}_{j.}\|_2^2$$
 (9)

and \boldsymbol{W} is the weight matrix (or the weighted adjacency matrix) with $\boldsymbol{L} = \boldsymbol{D} - \boldsymbol{W}, \, \boldsymbol{W} = \boldsymbol{W}^{\top}, \, W_{ij} \geq 0$ and $W_{ii} = 0$ for $1 \leq i, j \leq p$. Instead of performing a penalized minimization of $\operatorname{tr}(\boldsymbol{X}^{\top}\boldsymbol{L}\boldsymbol{X})$ to estimate \boldsymbol{L} , [14] minimizes a penalized $\operatorname{tr}(\boldsymbol{W}\boldsymbol{Z})$ w.r.t. \boldsymbol{W} for graph learning. Given \boldsymbol{W} , one has unique \boldsymbol{L} and the edge-set \mathcal{E} . [14] shows that minimization of $\operatorname{tr}(\boldsymbol{W}\boldsymbol{Z})$ is computationally more efficient than minimizing $\operatorname{tr}(\boldsymbol{X}^{\top}\boldsymbol{L}\boldsymbol{X})$ (as is done in [15]).

Define the space W_p of all valid $p \times p$ weight matrices W

$$\mathcal{W}_p = \left\{ \mathbf{W} \in \mathbb{R}^{p \times p} : \mathbf{W} = \mathbf{W}^\top, W_{ij} \ge 0, \right.$$
$$W_{ii} = 0, \ 1 \le i, j \le p \right\}$$
(10)

The following optimization problem is solved in [14]:

$$\min_{\boldsymbol{W} \in \mathcal{W}_p} \operatorname{tr}(\boldsymbol{W}\boldsymbol{Z}) + \frac{\beta}{2} \|\boldsymbol{W}\|_F^2 - \alpha \sum_{i=1}^p \ln\left(\sum_{j=1}^p W_{ij}\right)$$
(11)

with parameters $\alpha>0$ and $\beta\geq0$ controlling the "shape" (number and weights) of the edges. In (11), $\operatorname{tr}(\boldsymbol{W}\boldsymbol{Z})$ is the main cost but minimizing it alone w.r.t. \boldsymbol{W} is ill-posed ($\boldsymbol{W}=\mathbf{0}$ minimizes it). The sum $\sum_{j=1}^p W_{ij}$ is the node degree at node i and the term $\ln\left(\sum_{j=1}^p W_{ij}\right)$ is the logarithmic barrier on the node degrees. This

implies that the degrees are forced to be positive but do not prevent individual edges to vanish. It is noted in [14] that using only the logarithmic barrier ($\beta=0$) leads to very sparse graphs, and changing α only changes the scale of the solution, not the sparsity pattern. Therefore, the term $\frac{\beta}{2} \| \mathbf{W} \|_F^2$ is added to control the graph sparsity: larger values of β lead to more dense connectivity. In [14], optimization is carried for fixed $\alpha=1$ (search over β to a get a desired edge density) and then one scales \mathbf{W} to obtain a desired $\| \mathbf{W} \|$; we refer the reader to [14] for further details.

In the next section we modify optimization problem (11) to apply to multi-attribute data using the enlarged graph $\bar{\mathcal{G}}=(\bar{V},\bar{\mathcal{E}})$ discussed in Sec. 2.1.

3. GRAPH LEARNING FROM MULTI-ATTRIBUTE DATA

Consider (similar to Sec. 2.1, but without any probabilistic modeling) p data vector sequences $\mathbf{z}_i(t) \in \mathbb{R}^m$, $i=1,2,\cdots,p$, each vector with n samples $t=1,2,\cdots,n$. We associate \mathbf{z}_i with the ith node of a weighted undirected graph $\mathcal{G}=(V,\mathcal{E},\mathbf{W})$ where $V=[1,p], \mathcal{E}\subseteq V\times V$ and $\mathbf{W}\in \mathcal{W}_p$ (defined in (10)). Define $\mathbf{x}(t)\in\mathbb{R}^{mp}$ based on $\mathbf{z}_i(t)$ s, as in (3). Then we have n samples $\{\mathbf{x}(t)\}_{t=1}^n$. Based on $\mathcal{G}=(V,\mathcal{E},\mathbf{W})$, construct an enlarged weighted undirected graph $\bar{\mathcal{G}}=(\bar{V},\bar{\mathcal{E}},\bar{\mathbf{W}})$ with $\bar{\mathbf{W}}\in\mathcal{W}_{mp}$, similar to the enlarged graph $(\bar{V},\bar{\mathcal{E}})$ based on (V,\mathcal{E}) as discussed in Sec. 2.1. Then the ℓ th component $[\mathbf{z}_j(t)]_\ell$ of the data vector $\mathbf{z}_j(t)$ associated with node j of graph $\mathcal{G}=(V,\mathcal{E},\mathbf{W})$ is now the scalar data variable $\mathbf{x}_q(t)=[\mathbf{x}(t)]_q$ associated with node q of the enlarged graph $\bar{\mathcal{G}}=(\bar{V},\bar{\mathcal{E}},\bar{\mathbf{W}})$, where $q=(j-1)m+\ell, j=1,2,\cdots,p$ and $\ell=1,2,\cdots,m$. Corresponding to each edge $\{j,k\}$ in \mathcal{G} , there are m^2 edges in $\bar{\mathcal{G}}$ (see (4) and (5)).

Define the $mp \times n$ matrix

$$\bar{X} = \begin{bmatrix}
x_1(1) & x_1(2) & \cdots & x_1(n) \\
x_2(1) & x_2(2) & \cdots & x_2(n) \\
\vdots & \vdots & \ddots & \vdots \\
x_{mp}(1) & x_{mp}(2) & \cdots & x_{mp}(n)
\end{bmatrix}$$

$$= \begin{bmatrix}
z_1(1) & z_1(2) & \cdots & z_1(n) \\
z_2(1) & z_2(2) & \cdots & z_2(n) \\
\vdots & \vdots & \ddots & \vdots \\
z_p(1) & z_p(2) & \cdots & z_p(n)
\end{bmatrix}$$
(12)

Define $\bar{Z} \in \mathbb{R}^{(mp)\times(mp)}$ with (i,j)th elements $\bar{Z}_{ij} = \|\bar{X}_{i.} - \bar{X}_{j.}\|_2^2$. Following (8) and (9), and using the graph $\bar{\mathcal{G}} = (\bar{V}, \bar{\mathcal{E}}, \bar{W})$, consider

$$\operatorname{tr}(\bar{X}^{\top}\bar{L}\bar{X}) = \frac{1}{2}\operatorname{tr}(\bar{W}\bar{Z}) = \frac{1}{2}\sum_{i=1}^{mp}\sum_{j=1}^{mp}\bar{W}_{ij}\bar{Z}_{ij}$$
(14)

$$= \frac{1}{2} \sum_{i=1}^{p} \sum_{j=1}^{p} \left[\sum_{s=1}^{m} \sum_{k=1}^{m} \bar{W}_{qr} \bar{Z}_{qr} \right], \quad q = (i-1)m + s, \\ r = (j-1)m + k.$$
 (15)

Now observe that the underlying graph for multi-attribute data is a p-node graph $\mathcal G$, not the enlarged mp-node graph $\bar{\mathcal G}$, implying that the strength of edge $\{i,j\}$ between nodes i and j of $\mathcal G$ which is W_{ij} , should be the strength of connection of all m^2 edges between nodes q=(i-1)m+s and r=(j-1)m+k, $(s,k=1,2,\cdots,m)$, of the enlarged graph $\bar{\mathcal G}$, corresponding to nodes i and j of $\mathcal G$. We enforce this by requiring that for any $\{i,j\}$,

$$\bar{W}_{qr} = W_{ij} \text{ for } q = (i-1)m+s, r = (j-1)m+k, \\ s, k = 1, 2, \cdots, m.$$
 (16)

Then we have

$$\operatorname{tr}(\bar{\boldsymbol{W}}\bar{\boldsymbol{Z}}) = \sum_{i=1}^{p} \sum_{j=1}^{p} W_{ij} \left[\sum_{s=1}^{m} \sum_{k=1}^{m} \bar{Z}_{qr} \right] = \operatorname{tr}(\boldsymbol{W}\tilde{\boldsymbol{Z}})$$
(17)

where $\tilde{\boldsymbol{Z}} \in \mathbb{R}^{p \times p}$ with (i, j)th component

$$\tilde{Z}_{ij} = \sum_{s,k=1}^{m} \|\bar{X}_{q.} - \bar{X}_{r.}\|_{2}^{2} = \sum_{s,k=1}^{m} \left(\frac{1}{n} \sum_{t=1}^{n} (x_{q}(t) - x_{r}(t))^{2} \right)$$
(18)

and q, r are as in (16).

Following [14], we have the regularized/penalized optimization problem (as in (11))

$$\min_{\boldsymbol{W} \in \mathcal{W}_p} \operatorname{tr}(\boldsymbol{W}\tilde{\boldsymbol{Z}}) + \frac{\beta}{2} \|\boldsymbol{W}\|_F^2 - \alpha \sum_{i=1}^p \ln\left(\sum_{j=1}^p W_{ij}\right)$$
(19)

We provide a solution in Sec. 3.1.

3.1. Optimization

We will use the alternating direction method of multipliers (ADMM) approach [20] with variable splitting. The solution of [14] is different. Using variable splitting, consider (recall $W_{ii}=0$)

$$\min_{\boldsymbol{W} \in \mathcal{W}_{p}, \boldsymbol{\theta} \in \mathbb{R}^{p}} \left\{ \operatorname{tr}(\boldsymbol{W}\tilde{\boldsymbol{Z}}) - \alpha \sum_{i=1}^{p} \ln(\theta_{i}) + \frac{\beta}{2} \|\boldsymbol{W}\|_{F}^{2} \right\}$$
(20)

subject to
$$\theta_i = \sum_{j=1, j \neq i}^{p} W_{ij}, i = 1, 2, \dots, p.$$
 (21)

The scaled augmented Lagrangian for this problem is [20] ($\mathbf{1}_p$ is a column of all ones)

$$L_{\rho} = \operatorname{tr}(\boldsymbol{W}\tilde{\boldsymbol{Z}}) - \alpha \sum_{i=1}^{p} \ln(\theta_{i}) + \frac{\beta}{2} \|\boldsymbol{W}\|_{F}^{2} + \frac{\rho}{2} \|\boldsymbol{\theta} - \boldsymbol{W}\boldsymbol{1}_{p} + \boldsymbol{u}\|_{2}^{2}$$
(22)

where $\boldsymbol{u} \in \mathbb{R}^p$ is the dual variable, and $\rho > 0$ is the penalty parameter. Given the results $\boldsymbol{\theta}^{(k)}, \boldsymbol{W}^{(k)}, \boldsymbol{u}^{(k)}$ of the kth iteration, in the (k+1)th iteration, an ADMM algorithm executes the following three updates:

(a) With
$$L_a(\boldsymbol{W}) := \operatorname{tr}(\boldsymbol{W}\tilde{\boldsymbol{Z}}) + \frac{\beta}{2} \|\boldsymbol{W}\|_F^2 + \frac{\rho}{2} \|\boldsymbol{\theta}^{(k)} - \boldsymbol{W}\boldsymbol{1}_p + \boldsymbol{u}^{(k)}\|_2^2$$
, let $\boldsymbol{W}^{(k+1)} \leftarrow \operatorname{arg\,min}_{\boldsymbol{W} \in \mathcal{W}_p} L_a(\boldsymbol{W})$.

(b) With
$$L_b(\boldsymbol{\theta}) := -\alpha \sum_{i=1}^p \ln(\theta_i) + \frac{\rho}{2} \|\boldsymbol{\theta} - \boldsymbol{W}^{(k+1)} \mathbf{1}_p + \boldsymbol{u}^{(k)}\|_2^2$$
, let $\boldsymbol{\theta}^{(k+1)} \leftarrow \arg\min_{\boldsymbol{\theta}} L_b(\boldsymbol{\theta})$.

(c)
$$u^{(k+1)} \leftarrow u^{(k)} + \theta^{(k+1)} - W^{(k+1)} \mathbf{1}_n$$
.

In update (a) above, cost $L_a(\boldsymbol{W})$ is separable row-wise in \boldsymbol{W} : $L_a(\boldsymbol{W}) = \sum_{i=1}^p L_{ai}(\check{\boldsymbol{W}}_{i.})$ where $\check{\boldsymbol{W}}_{i.} \in \mathbb{R}^{p-1}$ is obtained from $\boldsymbol{W}_{i.}$ by deleting its ith row (recall that $W_{ii} = 0$), similarly, $\check{\boldsymbol{Z}}_{i.} \in \mathbb{R}^{p-1}$ is obtained from $\check{\boldsymbol{Z}}_{i.}$ by deleting its ith row, and

$$L_{ai}(\check{\mathbf{W}}_{i.}) = \check{\mathbf{W}}_{i.}^{\top} \check{\mathbf{Z}}_{i.} + \frac{\beta}{2} \check{\mathbf{W}}_{i.}^{\top} \check{\mathbf{W}}_{i.} + \frac{\rho}{2} (\theta_{i}^{(k)} - \check{\mathbf{W}}_{i.}^{\top} \mathbf{1}_{p-1} + u_{i}^{(k)})^{2}$$

$$= \check{\mathbf{W}}_{i.}^{\top} \left(\frac{\beta}{2} + \frac{\rho}{2} \mathbf{1}_{p-1} \mathbf{1}_{p-1}^{\top} \right) \check{\mathbf{W}}_{i.} + \frac{\rho}{2} (\theta_{i}^{(k)} - u_{i}^{(k)})^{2}$$

$$+ \left(\check{\mathbf{Z}}_{i.} - \rho(\theta_{i}^{(k)} + u_{i}^{(k)}) \mathbf{1}_{p-1} \right)^{\top} \check{\mathbf{W}}_{i.}$$
(23)

We need to minimize $L_{ai}(\check{\boldsymbol{W}}_{i.})$ w.r.t. $\check{\boldsymbol{W}}_{i.}$ subject to $W_{ij} \geq 0$. An iterative solution is given in [19]. [19] minimizes $(1/2)\boldsymbol{y}^{\top}\boldsymbol{Q}\boldsymbol{y} - \boldsymbol{y}^{\top}\boldsymbol{h}$ w.r.t. $\boldsymbol{y} \in \mathbb{R}^q$ subject to $y_{\ell} \geq 0 \ \forall \ell$, where \boldsymbol{Q} is symmetric, positive-definite and \boldsymbol{h} is arbitrary. The monotonically convergent iterative solution of [19] is

$$y_{\ell} \leftarrow y_{\ell} \left[\frac{2[\boldsymbol{Q}^{-}\boldsymbol{y}]_{\ell} + h_{\ell}^{+} + \delta}{[\text{abs}(\boldsymbol{Q})\boldsymbol{y}]_{\ell} + h_{\ell}^{-} + \delta} \right]$$
(24)

where $0 < \delta \ll 1$. In our problem of minimization of $L_{ai}(\check{\boldsymbol{W}}_{i.})$ w.r.t. $\check{\boldsymbol{W}}_{i.}$ subject to $W_{ij} \geq 0$, we have $\boldsymbol{Q} = \frac{\beta}{2} + \frac{\rho}{2} \mathbf{1}_{p-1} \mathbf{1}_{p-1}^{\top}$ independent of row i of \boldsymbol{W} , and $\boldsymbol{h} = \check{\boldsymbol{Z}}_{i.} - \rho(\theta_i^{(k)} + u_i^{(k)}) \mathbf{1}_{p-1}$ which depends upon row i of \boldsymbol{W} . Since all elements of \boldsymbol{Q} for our problem are positive, $\boldsymbol{Q}^- = \boldsymbol{0}$ and $\boldsymbol{Q}^+ = \boldsymbol{Q}$.

Now we turn to update (b). Notice that $L_b(\theta)$ is completely separable w.r.t. each component of θ : $L_b(\theta) = \sum_{i=1}^p L_{bi}(\theta_i)$, where $L_{bi}(\theta_i) = -\alpha \ln(\theta_i) + (\rho/2)(\theta_i - \boldsymbol{W}_{i.}^{(k+1)} \boldsymbol{1}_p + u_i^{(k)})^2$. Setting

$$0 = \frac{\partial L_{bi}(\theta_i)}{\partial \theta_i} = -\frac{\alpha}{\theta_i} + \rho \left(\theta_i - c_i\right)$$
where $c_i = \mathbf{W}_i^{(k+1)} \mathbf{1}_p - u_i^{(k)}$. (25)

This leads to a quadratic equation in θ_i . Since the *i*th node degree θ_i must be positive, we take

$$\theta_i^{(k+1)} = \frac{1}{2} \left(c_i + \sqrt{c_i + (4\alpha/\rho)} \right)$$
 (26)

where c_i is specified in (25).

Finally, update (c) follows from [20].

We initialize as $u^{(0)} = 0$, $W^{(0)} = 0$ and $\theta = 1_p$. For all numerical results presented later, we used $\rho = 2$ and set $\alpha = 1$.

As noted in [14], the objective function (19) is proper, convex, and lower-semicontinuous. Hence the ADMM algorithm is guaranteed to converge [20, Sec. 3.2], in the sense that we have primal residual convergence to 0, dual residual convergence to 0, and objective function convergence to the optimal value.

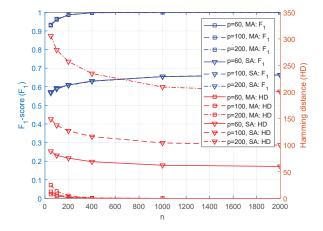


Fig. 1: F_1 -score and Hamming distance for the chain graph with m=3 attributes per node. The label MA refers to multi-attribute data, and SA refers to single attribute data using just the first attribute at each node. $\alpha=1$, and $\beta=10$ for MA graph learning and $\beta=1$ for SA graph learning.

4. NUMERICAL EXAMPLES

4.1. Synthetic Data

We consider a chain graph (an example in [11]) where p nodes are connected in succession. We use a probabilistic/statistical model to generate multi-attribute data, but apply signal smoothness-based graph learning. Using the notation of (6), we set $[\Omega^{(jk)}]_{st}$ $-0.5^{|s-t|}$ for $j=k=1,\cdots,p$ and $s\neq t=1,\cdots,m$ (offdiagonal terms are non-positive). For $j \neq k$, if the two nodes are not connected, we have $\Omega^{(jk)} = 0$, and if nodes j and k are connected in the chain graph, then $[\Omega^{(jk)}]_{st} = -0.2$ if $s \neq t$, and $[\Omega^{(jk)}]_{st} = 0$ if s = t. Then we fill in the diagonal terms of Ω to make it a valid graph Laplacian. This is simulation example 3 in [11, Sec. 5.1] except that off-diagonal terms of Ω are non-positive. Now add γI to Ω with γ picked to make minimum eigenvalue of $\mathbf{\Omega} + \gamma \mathbf{I}$ equal to 0.001. With $\mathbf{\Phi}\mathbf{\Phi}^{\top} = (\mathbf{\Omega} + \gamma \mathbf{I})^{-1}$, we generate $x = \Phi w$ with $w \in \mathbb{R}^{mp}$ as Gaussian $w \sim \mathcal{N}(0, I)$. We generate n i.i.d. observations from x, with m = 3, $p \in \{60, 100, 200\}$, $n \in \{50, 100, 200, 400, 1000, 2000\}$. The true value of $|\mathcal{E}| = p-1$.

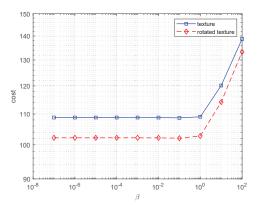


Fig. 2: Cross-validation cost (19) versus β , for β selection for image graph fitting

Simulation results based on 100 runs are shown in Fig. 1. We applied our proposed multi-attribute (labeled MA in figures) graph learning approach to estimate W, hence the edgeset \mathcal{E} , along with the restriction of this approach to a single attribute (labeled SA in figures; we used only the first attribute at each node, results for other attributes were very similar). We fixed $\alpha = 1$ and selected β over a grid of values to maximize the F_1 -score*. The selected values were $\beta = 10$ for multi-attribute and $\beta = 1$ single attribute; the results are not very sensitive to the choice of β . Fig. 1 shows the F_1 -score and the Hamming distance between the true and estimated edge sets \mathcal{E}_0 and $\hat{\mathcal{E}}$, respectively. It is seen that multi-attribute based graph learning is significantly better than single attribute graph learning, since for same number of nodes and sample size, the former yields significantly higher F_1 -score and significantly lower Hamming distance. As sample size n increases, the Hamming distance tends to zero for multi-attribute based graph learning, indicating that the the edgeset \mathcal{E} is learned perfectly.

Since computation the F_1 -score requires knowledge of true \mathcal{E}_0 , in practice, one would select β via cross-validation or some infor-

^{*}It is defined as $F_1=2 \times \operatorname{precision} \times \operatorname{recall}/(\operatorname{precision} + \operatorname{recall})$ where $\operatorname{precision} = |\hat{\mathcal{E}} \cap \mathcal{E}_0|/|\hat{\mathcal{E}}|$, $\operatorname{recall} = |\hat{\mathcal{E}} \cap \mathcal{E}_0|/|\mathcal{E}_0|$, and \mathcal{E}_0 and $\hat{\mathcal{E}}$ denote the true and estimated edge sets, respectively.

mation criterion-based approach. This is illustrated for the real data set later.

Algorithm 1 ADMM Algorithm for Multi-Attribute Graph learning **Input:** Number of samples n, number of nodes p, number of attributes m, data $\{x(t)\}_{t=1}^n$, $x \in \mathbb{R}^{mp}$, regularization parameters α and β , ADMM penalty parameter ρ , and $\delta = 10^{-16}$.

Output: estimated weight matrix W and edge-set \mathcal{E}

- 1: Calculate the (squared) distance matrix $\tilde{Z} \in \mathbb{R}^{p \times p}$ with (i, j)th 1: Calculate the (squared) distance matrix $\mathbf{Z} \in \mathbb{R}^n$ where component $\tilde{Z}_{ij} = \sum_{s,k=1}^m \left(\frac{1}{n}\sum_{t=1}^n (x_q(t) - x_r(t))^2\right)$, where q = (i-1)m + s and r = (j-1)m + k. 2: $\mathbf{Q} = \frac{\beta}{2} + \frac{\rho}{2} \mathbf{1}_{p-1} \mathbf{1}_{p-1}^{\top}$, $\mathbf{P} = \mathbf{Q}^{-1}$. 3: Initialize: $\mathbf{u}^{(0)} = \mathbf{0} \in \mathbb{R}^p$, $\mathbf{W}^{(0)} = \mathbf{0} \in \mathbb{R}^{p \times p}$ and $\mathbf{\theta} = \mathbf{1}_p = \mathbf{0}$
- column of p ones.
- 4: for $k = 0, 1, 2, \ldots$, until convergence, do
- 5: for $i = 1, 2, \dots, p$ do
- Define $\tilde{Z}_{i} \in \mathbb{R}^{p-1}$ from \tilde{Z}_{i} by deleting its *i*th row: 6:

$$\check{\boldsymbol{Z}}_{ij} = \left\{ \begin{array}{ll} \tilde{\boldsymbol{Z}}_{ij}, & 1 \leq j \leq i-1 \\ \tilde{\boldsymbol{Z}}_{i(j+1)}, & i \leq j \leq p-1 \,. \end{array} \right.$$

7:
$$h = \check{Z}_{i.} - \rho(\theta_i^{(k)} + u_i^{(k)}) \mathbf{1}_{p-1}$$
8: $y^{(0)} = \max(Ph, \mathbf{0}) \in \mathbb{R}^{p-1}$
9: $\mathbf{for} \ q = 1, 2, \dots$, until convergence, \mathbf{do}
10: $\mathbf{for} \ \ell = 1, 2, \cdots, p-1 \ \mathbf{do}$
11: $y_{\ell}^{(q)} = y_{\ell}^{(q-1)} \left[\frac{h_{\ell}^{+} + \delta}{[Qy^{(q-1)}]_{\ell} + h_{\ell}^{-} + \delta} \right]$
12: $\mathbf{end} \ \mathbf{for}$
13: $\mathbf{end} \ \mathbf{for}$
14: Update i th row of W as

$$W_{ij}^{(k+1)} = \begin{cases} y_j, & 1 \le j \le i - 1 \\ 0, & j = i \\ y_{j-1}, & i + 1 \le j \le p. \end{cases}$$

$$\begin{array}{lll} \text{15:} & \textbf{end for} \\ \text{16:} & \textbf{for } i=1,2,\cdots,p \ \ \textbf{do} \\ \text{17:} & c_i=\boldsymbol{W}_i^{(k+1)}\mathbf{1}_p-u_i^{(k)} \\ \text{18:} & \theta_i^{(k+1)}=\frac{1}{2}\left(c_i+\sqrt{c_i+(4\alpha/\rho)}\right) \\ \text{19:} & \textbf{end for} \\ \text{20:} & \boldsymbol{u}^{(k+1)}\leftarrow\boldsymbol{u}^{(k)}+\boldsymbol{\theta}^{(k+1)}-\boldsymbol{W}^{(k+1)}\mathbf{1}_p \\ \text{21:} & \textbf{end for} \end{array}$$

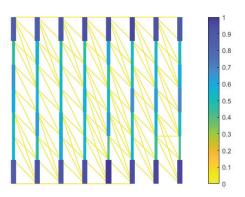
4.2. Real Data Example: Graphs of Color Texture Images

Following [12, 13] where grayscale texture images from a USC database are considered, we now consider color textures from the Amsterdam Library of Textures (ALOT) [21]. We use two versions of the image 108 from [21], c111.png and c111r60.png (400×400) patches shown in Figs. 3 and 4), photographed from different angles. The 400×400 patches were partitioned into non-overlapping 8×8 blocks, vectorized into 64-pixel columns, 3 colors associated with each pixel. Thus, we have m=3, p=64 and n=2500. The data were centered and mean-square value normalized to one before processing. To select β (α =1 throughout), we use cross-validation with criterion C_{cv} = the objective function (19), as follows. We random sample 70% of n (=2500) samples as training data for model fitting and use the remaining 30% samples as test data to compute C_{cv} . We average C_{cv} over 30 random samples for β values over a grid; the value of β is selected to minimize average C_{cv} . With the

selected $\hat{\beta}$, we perform model fitting over the entire dataset (with $\alpha = 1$). The average cost C_{cv} for graph learning for the two color images displayed in Fig. 3(a) and Fig. 4(a), are shown in Fig. 2. For both textures, we get $\hat{\beta} = 0.1$. These values were used for final graph model fitting. We take the final estimated W resulting from (multi-attribute) ADMM algorithm as edge weights, normalize maximum value to one, and show the resulting graphs (arranged as 8 × 8 nodes) with colored edge weights and link thickness also reflecting edge weight. Compare Figs. 3(a) and 3(b), and Figs. 4(a) and 4(b), respectively, to note that the strong link weights follow the texture orientation: vertical in Figs. 3(a) and 3(b), and horizontal in Figs. 4(a) and 4(b). This provides significant "visual" support for the fitted graphs.



(a) a texture (from [21])



(b) image graph of (a): number of edges =458

Fig. 3: Color texture graph example

5. CONCLUSIONS

The problem of how to learn the graph combinatorial Laplacian matrix defining the structure of an undirected weighted graph underlying a set of smooth multi-attribute signals, was addressed. We extended the single attribute approach of Kalofolias (2016) to multiattribute data. In a typical single-attribute graph modeling problem, one associates a scalar data variable with each node of the graph, whereas in multi-attribute graphical models, each node represents a

vector. Image graphs for grayscale images for modeling dependence of a pixel on neighboring pixels is an example of single-attribute graph modeling. These approaches do not apply to color images where one has three variables (RGB color components) per pixel node. Image graphs for color images is an example of multi-attribute graphical models. An ADMM algorithm was presented to optimize the objective function to infer the graph topology. We tested the proposed approach on synthetic as well as real data (color image graphs).

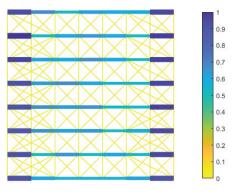
6. REFERENCES

- 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 (Methodological)*, 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 Research*, vol. 9, pp. 485-516, 2008.
- [9] J. Chiquet, G. Rigaill and M. Sundquist, "A multiattribute Gaussian graphical model for inferring multiscale regulatory networks: an application in breast cancer." In: Sanguinetti G., Huynh-Thu V. (eds), Gene Regulatory Networks. Methods in Molecular Biology, vol 1883. Humana Press, New York, NY, 2019.
- [10] M. Kolar, H. Liu and E.P. Xing, "Markov network estimation from multi-attribute data," in *Proc. 30th Intern. Conf. Machine Learning (ICML)*, Atlanta, GA, 2013.
- [11] M. Kolar, H. Liu and E.P. Xing, "Graph estimation from multi-attribute data," *J. Machine Learning Research*, vol. 15, pp. 1713-1750, 2014.
- [12] 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.
- [13] 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.
- [14] V. Kalofolias, "How to learn a graph from smooth signals," in *Proc. 19th Intern. Conf. Artificial Intelligence & Statistics* (AISTATS), Cadiz, Spain, 2016.

- [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] M. Belkin and P. Niyogi, "Laplacian eigenmaps and spectral techniques for embedding and clustering.," in *Proc. NIPS*, vol. 14, pp. 585-591, 2001.
- [17] 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.
- [18] G. Fracastoro, D. Thanou and P. Frossard, "Graph transform optimization with application to image compression," *IEEE Trans. Image Process.*, vol. 29, pp. 419-432, 2020.
- [19] X. Xiao and D. Chen, "Multiplicative iteration for nonnegative quadratic programming," arXiv:1406.1008v1 [math.NA], 4 June 2014.
- [20] 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.
- [21] Amsterdam Library of Textures (ALOT), http://aloi.science.uva.nl/public_alot.



(a) "rotated" texture (from [21])



(b) image graph of (a): number of edges =640

Fig. 4: Another color texture graph example