# An Algorithm for Graph-Fused Lasso Based on Graph Decomposition

Feng Yu
Department of Mathematics, University of Central Florida
and
Yi Yang

Department of Mathematics and Statistics, McGill University and

Teng Zhang\*

Department of Mathematics, University of Central Florida

August 8, 2019

#### Abstract

This work proposes a new algorithm for solving the graph-fused lasso (GFL), a method for parameter estimation that operates under the assumption that the signal tends to be locally constant over a predefined graph structure. The proposed method applies the alternating direction method of multipliers (ADMM) algorithm and is based on the decomposition of the objective function into two components. While ADMM has been widely used in this problem, existing works such as network lasso decompose the objective function into the loss function component and the total variation penalty component. In comparison, this work proposes to decompose the objective function into two components, where one component is the loss function plus part of the total variation penalty, and the other component is the remaining total variation penalty. Compared with the network lasso algorithm, this method has a smaller computational cost per iteration and converges faster in most simulations numerically.

Keywords: alternating direction methods of multipliers, graph-fused lasso, nonsmooth convex optimization

<sup>\*</sup>The authors are gratefully supported by National Science Foundation (NSF), grant CNS-1739736.

## 1 Introduction

In this article, we consider graph-fused lasso, an estimation method based on noisy observations and the assumption that the signal tends to be locally constant over a predefined graph structure. Given a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is the set of vertices and  $\mathcal{E}$  is the set of edges, we let  $\mathbf{x}_i \in \mathbb{R}^p$  be the signal that is associated with the *i*-the vertex of the graph, then GFL is defined as the solution to the following optimization problem:

$$\{\hat{\mathbf{x}}_i\}_{i\in\mathcal{V}} = \underset{\{\mathbf{x}_i\}_{i\in\mathcal{G}}\subset\mathbb{R}^p}{\operatorname{arg min}} \sum_{i\in\mathcal{V}} f_i(\mathbf{x}_i) + \lambda \sum_{(s,t)\in\mathcal{E}} \|\mathbf{x}_r - \mathbf{x}_s\|,$$
(1)

where the first component is a loss function for the observation  $\mathbf{x}_i$ , and the second component uses the total variation norm to penalize the difference between the two signals on the edges in the graph.

There have been extensive studies on (1) with p = 1 (i.e.,  $\mathbf{x}_i$  are scalars) and many algorithms have been developed. When graph  $\mathcal{G}$  is a one-dimensional chain graph, then it is the standard fused lasso problem (Tibshirani et al., 2005). For this problem, there exist finite-step algorithms with computational costs of O(n): a taut-string method proposed by Davies and Kovac (2001), a method based on analyzing its dual problem by Condat (2013), a dynamic programming-based approach by Johnson (2013), and a modular proximal optimization approach by Barbero and Sra (2014) all solve the problem with O(n)complexity. When the  $\mathcal{G}$  is a two-dimensional grid graph, it has important applications in image denoising and it often referred to as total-variation denoising (Rudin et al., 1992), and parametric max-flow algorithm (Chambolle and Darbon, 2009) can be used to solve (1) in finite steps. When the graph is a tree, Kolmogorov et al. (Kolmogorov et al., 2016) extended the dynamic programming approach of Johnson to solve the fused lasso problem. While these algorithms can find the exact solution in finite steps, they only apply to some specific graph structure and do not work for general graphs. In addition, they cannot be naturally generalized to the setting of p > 1, which is sometimes called group fused lasso (Bleakley and Vert, 2011).

In addition, many iterative algorithms based on convex optimization algorithms have been proposed to solve (1). For example, Liu et al. (2010) use a projected gradient descent method to solve the dual of (1), and reformulate it as the problem of finding an "appropriate" subgradient of the fused penalty at the minimizer. Chen et al. (2012) propose the smoothing proximal gradient (SPG) method. Lin et al. (2014) proposed an alternating linearization method. Yu et al. (2015) proposed a majorization-minimization method. One of the more popular methods is the alternating direction method of multipliers (ADMM), due to its simplicity and competitive empirical performance. Ye and Xie (2011) and Wahlberg et al. (2012) proposed algorithms based on the ADMM method. However, there is a step of solving a linear system for the  $n \times n$  matrix  $\mathbf{I} + \rho \mathbf{D}^T \mathbf{D}$ , which is usually in the order of  $O(n^2)$ . Zhu (2017) proposed a modified ADMM algorithm that has a smaller computational cost of O(n) in an update step, but this modification generally converges slower. Ramdas and Tibshirani (2015) proposed a special ADMM algorithm that used dynamic programming in one of the update step, which can be used in the trend filtering problem, or when **D** has a diagonal structure. Barbero and Sra (2014) propose a method based on the Douglas-Rachford decomposition for the two-dimensional grid graph, which can be considered as the dual algorithm of ADMM (Eckstein and Bertsekas, 1992). Tansey and Scott (2015) leveraged fast 1D fused lasso solvers in an ADMM method based on graph decomposition, but it can only be applied to the case when p=1. Hallac et al. (2015) proposed the network lasso algorithm based on ADMM that can be applied to any generic graph and any  $p \ge 1$ .

There exist other types of algorithms as well. Friedman et al. (2007) and Arnold and Tibshirani (2016) gave solution path algorithms (tracing out the solution over all  $\lambda \geq 0$ ). Some other algorithms include a working-set/greedy algorithm (Landrieu and Obozinski, 2017) and an algorithm based on an active set search (Kovac and Smith, 2011).

Among all algorithms, the network lasso (Hallac et al., 2015) is particularly interesting since it is scalable to any large graphs and can be applied to the case  $p \ge 1$ . The algorithm proposed in this work follows this direction and can be considered as an improvement of the network lasso algorithm. The main contribution of this work is a novel ADMM method by dividing the objective function into two parts based on graph decomposition so that one of the subgraphs does not contain any two adjacent edges. This method can be applied to any

graph and can be generalized to some other problems such as trend filtering. Compared with the network lasso algorithm in (Hallac et al., 2015), it reduces the computational complexity per iteration and achieves a faster convergence rate.

The rest of this paper is organized as follows. In Section 2 we introduce our proposed method and analyze its computational complexity per iteration as well as convergence rates and establish the advantage of the proposed algorithm theoretically. Then we compare our algorithm with the network lasso algorithm in Section 3, both in simulated data sets and a real-life data set, which verifies the advantage of the proposed algorithm numerically.

## 2 Proposed Method

In this section, we will review the network lasso algorithm (Hallac et al., 2015) for solving (1) in Section 2.1, present our algorithm in Sections 2.2 and 2.3, and analyze its performance in terms of computational cost per iteration and convergence rate in Sections 2.4 and 2.5.

#### 2.1 Review: Network lasso

Hallac et al. (2015) introduce the following "Network Lasso" algorithm: for any  $(s,t) \in \mathcal{E}$ , introduce a pair of variables  $\mathbf{z}_{st}, \mathbf{z}_{ts} \in \mathbb{R}^p$ , which are the copies of  $\mathbf{x}_r$  and  $\mathbf{x}_s$  respectively, and rewrite the problem (1) as follows:

$$\underset{\{\mathbf{x}_i\}_{i \in \mathcal{G}}, \{\mathbf{z}_{st}, \mathbf{z}_{ts}\}_{(s,t) \in \mathcal{E}}}{\arg \min} \sum_{i \in \mathcal{V}} f_i(\mathbf{x}_i) + \lambda \sum_{(s,t) \in \mathcal{E}} \|\mathbf{z}_{st} - \mathbf{z}_{ts}\|, \text{ s.t. } \mathbf{x}_s = \mathbf{z}_{st} \text{ and } \mathbf{x}_t = \mathbf{z}_{ts} \text{ for all } (s,t) \in \mathcal{E}.$$
(2)

Then the standard ADMM routine would apply: let  $\mathbf{u}_{st}$ ,  $\mathbf{u}_{ts}$  be the dual variables for  $\mathbf{x}_s - \mathbf{z}_{st}$  and  $\mathbf{x}_t - \mathbf{z}_{ts}$  respectively, then the augmented Lagrangian is (here x and z represents  $\{\mathbf{x}_i\}_{i\in\mathcal{G}}$ ):

$$L_{\rho}(x, z, u) = \sum_{i \in \mathcal{V}} f_i(\mathbf{x}_i) + \sum_{(s,t) \in \mathcal{E}} \left( \lambda \|\mathbf{z}_{st} - \mathbf{z}_{ts}\| + \mathbf{u}_{st}^T (\mathbf{x}_s - \mathbf{z}_{st}) + \mathbf{u}_{ts}^T (\mathbf{x}_t - \mathbf{z}_{ts}) + \frac{\rho}{2} \|\mathbf{x}_s - \mathbf{z}_{st}\|^2 + \frac{\rho}{2} \|\mathbf{x}_t - \mathbf{z}_{ts}\|^2 \right)$$
(3)

and the algorithm can be written as

$$x^{(k+1)} = \arg\min_{x} L_{\rho}(x, z^{(k)}, u^{(k)})$$
(4)

$$z^{(k+1)} = \arg\min_{z} L_{\rho}(x^{(k+1)}, z, u^{(k)})$$
(5)

$$\mathbf{u}_{st}^{(k+1)} = \mathbf{u}_{st}^{(k)} + \rho(\mathbf{x}_s^{(k+1)} - \mathbf{z}_{st}^{(k+1)}), \ \mathbf{u}_{ts}^{(k+1)} = \mathbf{u}_{ts}^{(k)} + \rho(\mathbf{x}_t^{(k+1)} - \mathbf{z}_{ts}^{(k+1)}). \tag{6}$$

The advantage of this algorithm is that, in each iteration, the optimization problem can be decomposed into smaller subproblems: the updates of x requires solving problems in the form of  $\min_{\mathbf{x}_i} f_i(\mathbf{x}_i) + \|\mathbf{x}_i - \mathbf{t}\|^2$ , which has explicit solutions for a large range of  $f_i$ ; and the updates of z requires solving  $\min_{\mathbf{z}_{st},\mathbf{z}_{ts}} \|\mathbf{z}_{ts} - \mathbf{t}_1\|^2 + \|\mathbf{z}_{st} - \mathbf{t}_2\|^2 + \lambda \|\mathbf{z}_{ts} - \mathbf{z}_{st}\|$ , which has explicit solutions.

## 2.2 Proposed Method: A different way of splitting the objective function

In this section, we will propose another ADMM algorithm for solving (1), based on the reformulation as follows: We will divide the set of edges  $\mathcal{E}$  into  $\mathcal{E}_0$  and  $\mathcal{E}_1$ , such that the set  $\mathcal{E}_0$  does not contain two neighboring edges, and then solve the following optimization problem:

$$\underset{\{\mathbf{x}_{i}\}_{i \in \mathcal{V}}, \{\mathbf{z}_{st}\}_{(s,t) \in \mathcal{E}_{1}}}{\operatorname{arg min}} \left( \sum_{i \in \mathcal{V}} f_{i}(\mathbf{x}_{i}) + \lambda \sum_{(s,t) \in \mathcal{E}_{0}} \|\mathbf{x}_{s} - \mathbf{x}_{t}\| \right) + \lambda \sum_{(s,t) \in \mathcal{E}_{1}} \|\mathbf{z}_{st} - \mathbf{z}_{ts}\|,$$
(7)
$$\text{s.t. } \mathbf{x}_{s} = \mathbf{z}_{st} \text{ and } \mathbf{x}_{t} = \mathbf{z}_{ts} \text{ for all } (s,t) \in \mathcal{E}_{1}.$$

Since this formulation is different than (2), its associated ADMM routine is also different. Let  $\mathbf{u}_{st}$ ,  $\mathbf{u}_{ts}$  be the dual variables for  $\mathbf{x}_s - \mathbf{z}_{st}$  and  $\mathbf{x}_t - \mathbf{z}_{ts}$  respectively, then the augmented Lagrangian is

$$\hat{L}_{\rho}(x, z, u) = \sum_{i \in \mathcal{V}} f_i(\mathbf{x}_i) + \lambda \sum_{(s,t) \in \mathcal{E}_0} \|\mathbf{x}_s - \mathbf{x}_t\| + \sum_{(s,t) \in \mathcal{E}_1} \left(\lambda \|\mathbf{z}_{st} - \mathbf{z}_{ts}\| + \mathbf{u}_{st}^T (\mathbf{x}_s - \mathbf{z}_{st})\right)$$
(8)

+ 
$$\mathbf{u}_{ts}^{T}(\mathbf{x}_{t} - \mathbf{z}_{ts}) + \frac{\rho}{2} \|\mathbf{x}_{s} - \mathbf{z}_{st}\|^{2} + \frac{\rho}{2} \|\mathbf{x}_{t} - \mathbf{z}_{ts}\|^{2}$$
. (9)

and the update formula is

$$x^{(k+1)} = \arg\min_{x} \hat{L}_{\rho}(x, z^{(k)}, u^{(k)})$$
(10)

$$z^{(k+1)} = \arg\min_{z} \hat{L}_{\rho}(x^{(k+1)}, z, u^{(k)})$$
(11)

$$\mathbf{u}_{st}^{(k+1)} = \mathbf{u}_{st}^{(k)} + \rho(\mathbf{x}_s^{(k+1)} - \mathbf{z}_{st}^{(k+1)}), \ \mathbf{u}_{ts}^{(k+1)} = \mathbf{u}_{ts}^{(k)} + \rho(\mathbf{x}_t^{(k+1)} - \mathbf{z}_{ts}^{(k+1)}). \tag{12}$$

While the update formula for x (10) is similar to (4), it requires solving a slightly different problem due to the additional component  $\lambda \sum_{(s,t)\in\mathcal{E}_0} \|\mathbf{x}_s - \mathbf{x}_t\|$ . For any  $(s,t)\in\mathcal{E}_0$ , the ADMM procedure needs to solve

$$\underset{\mathbf{x}_s, \mathbf{x}_t \in \mathbb{R}^p}{\arg \min} f_s(\mathbf{x}_s) + f_t(\mathbf{x}_t) + \frac{\rho}{2} d_s(\mathbf{x}_s - \mathbf{t}_1)^2 + \frac{\rho}{2} d_t(\mathbf{x}_t - \mathbf{t}_2)^2 + \lambda \|\mathbf{x}_s - \mathbf{x}_t\|, \tag{13}$$

where  $\mathbf{t}_1, \mathbf{t}_2 \in \mathbb{R}^p$  and  $d_s$  denotes the degree of the vertex s in the graph  $(\mathcal{V}, \mathcal{E}_1)$ . For many choices of  $f_s$  and  $f_t$  (for example, square functions), this problem has an explicit solution.

Intuitively, we expect that the proposed algorithm would achieve a faster convergence rate than (2): (7) has fewer "dummy variables" in the form of  $\mathbf{z}_{st}$  ( $2|\mathcal{E}_1|$  instead of  $2|\mathcal{E}|$ ) and (8) has fewer dual parameters than (7). As a result, the Lagrangian in (8) contains fewer variables than (3) and we expect the algorithm to converge faster.

## 2.3 An equivalent formulation

In this section, we propose an equivalent form of the update formula (10)-(12), with a smaller computational cost per iteration. In particular, we consider the ADMM algorithm for solving

$$\underset{\{\mathbf{x}_{i}\}_{i \in \mathcal{V}}, \{\mathbf{z}_{st}\}_{(s,t) \in \mathcal{E}_{1}}}{\operatorname{arg min}} \left( \sum_{i \in \mathcal{V}} f_{i}(\mathbf{x}_{i}) + \lambda \sum_{(s,t) \in \mathcal{E}_{0}} \|\mathbf{x}_{s} - \mathbf{x}_{t}\| \right) + \lambda \sum_{(s,t) \in \mathcal{E}_{1}} \|\mathbf{z}_{st}\|, \tag{14}$$
s.t.  $\mathbf{z}_{st} = \mathbf{x}_{s} - \mathbf{x}_{t}$ .

For its Lagrangian

$$\tilde{L}_{\rho}(x, z, u) = \sum_{i \in \mathcal{V}} f_i(\mathbf{x}_i) + \lambda \sum_{(s,t) \in \mathcal{E}_0} \|\mathbf{x}_s - \mathbf{x}_t\| + \sum_{(s,t) \in \mathcal{E}_1} \left(\lambda \|\mathbf{z}_{st}\| + \mathbf{u}_{st}^T (\mathbf{z}_{st} - \mathbf{x}_s + \mathbf{x}_t)\right)$$
(15)

$$+\frac{\rho}{2}\|\mathbf{z}_{st} - \mathbf{x}_s + \mathbf{x}_t\|^2\Big),\tag{16}$$

the preconditioned ADMM algorithm can be written as

$$x^{(k+1)} = \arg\min_{x} \tilde{L}_{\rho}(x, z^{(k)}, u^{(k)}) + \frac{\rho}{2} \sum_{(s,t) \in \mathcal{E}_{1}} \|\mathbf{x}_{s} + \mathbf{x}_{t} - \mathbf{x}_{s}^{(k)} - \mathbf{x}_{t}^{(k)}\|^{2}$$
(17)

$$z^{(k+1)} = \arg\min_{z} \tilde{L}_{\rho}(x^{(k+1)}, z, u^{(k)})$$
(18)

$$u_{st}^{(k+1)} = u_{st}^{(k)} + \rho(z_{st}^{(k+1)} - z_{s}^{(k+1)} + z_{t}^{(k+1)}).$$
(19)

It is called preconditioned ADMM due to the additional component  $\frac{\rho}{2} \sum_{(s,t) \in \mathcal{E}_1} \|\mathbf{x}_s + \mathbf{x}_t - \mathbf{x}_s^{(k)} - \mathbf{x}_t^{(k)}\|^2$  in the update of x.

Compared with the update formula (10)-(12), the computational cost of (17)-(19) is smaller since there are no "dummy variables"  $\mathbf{z}_{ts}$  and its associated dual variables  $\mathbf{u}_{ts}$ . In addition, Lemma 2.1 shows that (17)-(19) is equivalent to the update formula (10)-(12). Its proof is deferred to Section 5.

**Lemma 2.1.** The update formula (17)-(19) with  $\rho = \rho_0$  is equivalent to the update formula (10)-(12) with  $\rho = 2\rho_0$ .

#### 2.4 Implementation and its computational cost per iteration

Based on the update formulas (17)-(19), the implementation of our proposed algorithm is described as Algorithm 1.

As this algorithm depends on the graph decomposition  $\mathcal{E}_0 \cup \mathcal{E}_1$ , in practice, we use a greedy algorithm to find  $\mathcal{E}_0$  as follows: First, label all edges in some arbitrary order and  $\mathcal{E}_0$  be an empty set. Second, cycle once through each edge and add it to  $\mathcal{E}_0$  if it is not neighboring any existing edges in  $\mathcal{E}_0$ . In fact,  $\mathcal{E}_0$  is called matching in graph theory and there are numerous algorithms for finding a matching within a graph.

To compare the computational cost per iteration Algorithm 1 and the network lasso, we investigate a commonly used special case that  $f_i(\mathbf{x}_i) = ||\mathbf{x}_i - \mathbf{y}_i||^2$ . Then step 1 in Algorithm 1 requires the following Lemma 2.2. We skip its proof since it is relatively straightforward and has been discussed in works such as (Hallac et al., 2015).

Algorithm 1 The implementation of the ADMM method in (17).

**Input:** Graph  $(\mathcal{V}, \mathcal{E})$  and its partition  $\mathcal{E} = \mathcal{E}_0 \cup \mathcal{E}_1$  ( $\mathcal{E}_0$  has not neighboring edges); loss functions  $\{f_i\}_{i\in\mathcal{V}}$ ; parameters  $\rho$  and  $\lambda$ .

Initialization: Initialize  $\{\mathbf{x}_i^{(0)}\}, \{\mathbf{z}_i^{(0)}\}_{i\in\mathcal{V}}\subset\mathbb{R}^p, \{\mathbf{u}_{st}^{(0)}\}_{(s,t)\in\mathcal{E}_1}\subset\mathbb{R}^p.$ 

**Loop:** Iterate Steps 1–4 until convergence:

1: For any  $(s,t) \in \mathcal{E}_0$ ,  $(\mathbf{x}_s^{(k+1)}, \mathbf{x}_t^{(k+1)}) = \arg\min_{\mathbf{x}_s, \mathbf{x}_t} f_s(\mathbf{x}_s) + f_t(\mathbf{x}_t) + \lambda \|\mathbf{x}_s - \mathbf{x}_t\| + \rho (d_s \|\mathbf{x}_s\|^2 + d_t \|\mathbf{x}_t\|^2) + \mathbf{x}_s^T \mathbf{t}_s^{(k)} + \mathbf{x}_t^T \mathbf{t}_t^{(k)}$ , where

$$\mathbf{t}_{i}^{(k)} = \sum_{j:(i,j) \in \mathcal{E}_{1}} [-(\mathbf{u}_{ij}^{(k)} + \rho \mathbf{z}_{ij}^{(k)}) - \rho(\mathbf{x}_{i}^{(k)} + \mathbf{x}_{j}^{(k)})] + \sum_{j:(j,i) \in \mathcal{E}_{1}} [(\mathbf{u}_{ji}^{(k)} + \rho \mathbf{z}_{ji}^{(k)}) - \rho(\mathbf{x}_{i}^{(k)} + \mathbf{x}_{j}^{(k)})]$$

2:For any  $i \in \mathcal{V}$  and does not belong to any edges in  $\mathcal{E}_0$ ,  $\mathbf{x}_i^{(k+1)} = \arg\min_{\mathbf{x}_i} f_i(\mathbf{x}_i) + \rho d_i ||\mathbf{x}_i||^2 + \mathbf{x}_i^T \mathbf{t}_i^{(k)}$ .

3: For any  $(s,t) \in \mathcal{E}_1$ ,  $\mathbf{z}_{st}^{(k+1)} = \arg\min_{\mathbf{z}_{st}} \frac{\rho}{2} \|\mathbf{z}_{st}\|^2 + \mathbf{z}_{st} (\mathbf{u}_{st}^{(k)} - \rho \mathbf{x}_s^{(k+1)} + \rho \mathbf{x}_t^{(k+1)}) + \lambda \|\mathbf{z}_{st}\| =$ threshold $(\mathbf{x}_s^{(k+1)} - \mathbf{x}_t^{(k+1)} - \frac{\mathbf{u}_{st}^{(k)}}{\rho}, \frac{\lambda}{\rho})$ .

4: For any  $(s,t) \in \mathcal{E}_1$ ,  $\mathbf{u}_{st}^{(k+1)} = \mathbf{u}_{st}^{(k)} + \rho(\mathbf{z}_{st}^{(k+1)} - \mathbf{x}_{s}^{(k+1)} + \mathbf{x}_{t}^{(k+1)})$ .

**Output:** The solution to (1),  $\hat{\mathbf{x}}_i = \lim_{k \to \infty} \mathbf{x}_i^{(k)}$  for all  $i \in \mathcal{V}$ .

Lemma 2.2. For any  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ ,

$$\arg \min_{\mathbf{x}, \mathbf{y} \in \mathbb{R}^p} c_1 \|\mathbf{x} - \mathbf{a}\|^2 + c_2 \|\mathbf{y} - \mathbf{b}\|^2 + \lambda \|\mathbf{x} - \mathbf{y}\|$$

$$= \begin{cases} \left(\frac{c_1 \mathbf{a} + c_2 \mathbf{b}}{c_1 + c_2}, \frac{c_1 \mathbf{a} + c_2 \mathbf{b}}{c_1 + c_2}\right), & \text{if } 2c_1 c_2 \|\mathbf{a} - \mathbf{b}\| \leq (c_1 + c_2)\lambda \\ \left(\mathbf{a} - \frac{\lambda}{2c_1} \frac{\mathbf{a} - \mathbf{b}}{\|\mathbf{a} - \mathbf{b}\|}, \mathbf{b} - \frac{\lambda}{2c_2} \frac{\mathbf{b} - \mathbf{a}}{\|\mathbf{b} - \mathbf{a}\|}\right), & \text{otherwise.} \end{cases}$$

Now let us investigate the computational complexity per iteration of Algorithm 1, by keep track of the multiplications of a scalar and a vector of  $\mathbb{R}^p$  (denoted as multiplications) and the additions of two vectors of  $\mathbb{R}^p$  (denoted as additions). In particular, the calculation of  $\mathbf{t}_s/(1+\rho_s)$  in the update of x requires  $2|\mathcal{E}_1|+n$  multiplications and  $2|\mathcal{E}_1|+n$  additions between, and step 1 requires an additional cost of at most  $2|\mathcal{E}_0|$  multiplications and  $2|\mathcal{E}_0|$  additions, and  $|\mathcal{E}_0|$  operations of finding the norm of a vector of length p and  $|\mathcal{E}_0|$  operations of comparing two scalars. Step 3 requires  $3|\mathcal{E}_1|$  additions,  $|\mathcal{E}_1|$  multiplications and  $|\mathcal{E}_1|$  comparisons. Step 4 requires  $3|\mathcal{E}_1|$  additions and  $|\mathcal{E}_1|$  multiplications.

Note that the network lasso algorithm is equivalent to the case where  $\mathcal{E}_0 = \emptyset$  and  $\mathcal{E}_1 = \mathcal{E}$ , we may compare the computational cost per iteration between Algorithm 1 and the network lasso algorithm, and it is clear that Algorithms 1 has a smaller computational cost per iteration compared to network lasso.

#### 2.5 Convergence Rate

Now let us investigate the theoretical convergence rate. First, we introduce a general theory on the local convergence of ADMM. Its proof is deferred to Section 5.

Theorem 2.3 (Local Convergence Rate of ADMM). Considering the problem of minimizing  $f_1(\mathbf{x}_1) + f_2(\mathbf{x}_2)$  subject to  $\mathbf{A}_1\mathbf{x}_1 + \mathbf{A}_2\mathbf{x}_2 = \mathbf{b}$ . Assuming that around the solution  $(\mathbf{x}^*, \mathbf{y}^*)$ ,  $\partial f_1(\mathbf{x}) = \mathbf{C}_1\mathbf{x} + \mathbf{c}_1$  and  $\partial f_2(\mathbf{x}) = \mathbf{C}_2\mathbf{x} + \mathbf{c}_2$ , then the local convergence rate of the ADMM algorithm is  $O(c(\rho)^k)$ , where k is the number of iterations and  $c(\rho)$  is the largest real components among all eigenvalue of

$$\frac{1}{2}[(\mathbf{I} - 2(\mathbf{I} + \rho \mathbf{A}_2 \mathbf{C}_2^{-1} \mathbf{A}_2^T)^{-1})(\mathbf{I} - 2(\mathbf{I} + \rho \mathbf{A}_1 \mathbf{C}_1^{-1} \mathbf{A}_1^T)^{-1}) + \mathbf{I}].$$

Considering that Algorithm 1 is obtained from solving (7), the convergence rate of Algorithm 1 follows from this theorem with

$$f_1(\{\mathbf{x}_i\}_{i\in\mathcal{V}}) = \sum_{i\in\mathcal{V}} f_i(\mathbf{x}_i) + \lambda \sum_{(s,t)\in\mathcal{E}_0} \|\mathbf{x}_s - \mathbf{x}_t\|, \ f_2(\{\mathbf{z}_{st}\}_{(s,t)\in\mathcal{E}_1}) = \lambda \sum_{(s,t)\in\mathcal{E}_1} \|\mathbf{z}_{st} - \mathbf{z}_{ts}\|.$$

That is,  $\mathbf{x}_1$  in Theorem 2.3 is replaced by  $\{\mathbf{x}_i\}_{i\in\mathcal{G}}$ ,  $\mathbf{x}_2$  in Theorem 2.3 is replaced by  $\{\mathbf{z}_{st}\}_{(s,t)\in\mathcal{E}_1}$ , and  $\mathbf{A}_1\mathbf{x}_1 + \mathbf{A}_2\mathbf{x}_2 = \mathbf{b}$  is replaced by  $\mathbf{x}_s = \mathbf{z}_{st}$  and  $\mathbf{x}_t = \mathbf{z}_{ts}$  for all  $(s,t)\in\mathcal{E}_1$ . Therefore, we have  $\mathbf{A}_1 \in \mathbb{R}^{np\times 2p|\mathcal{E}_1|}$ , defined such that  $\mathbf{A}_1(2i-1,s_i) = \mathbf{I}_{p\times p}$  and  $\mathbf{A}_1(2i,t_i) = \mathbf{I}_{p\times p}$  if  $(s_i,t_i)$  is the *i*-th edge in  $\mathcal{E}_1$ , and  $\mathbf{A}_2 = -\mathbf{I}_{2p|\mathcal{E}_1|\times 2p|\mathcal{E}_1|}$ . The matrix  $\mathbf{C}_1 \in \mathbb{R}^{pn\times pn}$  can be generated by the following three steps. First, the (i,i)-th  $p\times p$  block is given by

$$\mathbf{C}_1(i,i) = \operatorname{Hessian} f_i(\mathbf{x}_i^*)$$

Second, for  $(i, j) \in \mathcal{E}_0$  we let  $\mathbf{T}(i, j) = \frac{1}{\|\mathbf{x}_i^* - \mathbf{x}_j^*\|} \mathbf{I} - \frac{1}{\|\mathbf{x}_i^* - \mathbf{x}_j^*\|^3} (\mathbf{x}_i^* - \mathbf{x}_j^*) (\mathbf{x}_i^* - \mathbf{x}_j^*)^T$  if  $\mathbf{x}_i^* \neq \mathbf{x}_j^*$ , and  $\mathbf{T}(i, j) = \infty \mathbf{I}$  if  $\mathbf{x}_i^* = \mathbf{x}_j^*$ . Third, we update the (i, i), (i, j), (j, i), (j, j)-th  $p \times p$  blocks of  $\mathbf{C}_1$  by

$$\mathbf{C}_1(i,i) \leftarrow \mathbf{C}_1(i,i) + \mathbf{T}(i,j), \quad \mathbf{C}_1(j,j) \leftarrow \mathbf{C}_1(j,j) + \mathbf{T}(i,j),$$
  
 $\mathbf{C}_1(i,j) \leftarrow \mathbf{C}_1(i,j) - \mathbf{T}(i,j), \quad \mathbf{C}_1(j,i) \leftarrow \mathbf{C}_1(j,i) - \mathbf{T}(i,j).$ 

The matrix  $\mathbf{C}_2 \in \mathbb{R}^{2p|\mathcal{E}_1|\times 2p|\mathcal{E}_1|}$  is generated as follows: if the *i*-th edge in  $\mathcal{E}_1$ ,  $(s_i, t_i)$ , satisfies that  $\mathbf{x}_{s_i}^* \neq \mathbf{x}_{t_i}^*$ , then the (i, i)-th  $2p \times 2p$  block of  $\mathbf{C}_2$  is given by  $[\mathbf{T}(s_i, t_i), -\mathbf{T}(s_i, t_i); -\mathbf{T}(s_i, t_i), \mathbf{T}(s_i, t_i)]$ . The remaining  $2p \times 2p$  blocks are all zero matrices.

Note that the network lasso algorithm is equivalent to Algorithm 1 with  $\mathcal{E}_0 = \emptyset$  and  $\mathcal{E}_1 = \mathcal{E}$ , this result can also be applied to analyze the convergence rate of the network lasso algorithm.

While it is difficult to compare their convergence rates in general due to the complexities of  $\mathbf{A}_i$  and  $\mathbf{C}_i$ , we can calculate the convergence rate numerically for some specific examples. Here we assume that  $\mathcal{G}$  is the one-dimensional chain graph defined by  $\mathcal{V} = \{1, 2, \dots, n\}$  and  $\mathcal{E} = \{(1, 2), (2, 3), \dots, (n - 1, n)\}$ , and the partition such that  $\mathcal{E}_0 = \{(1, 2), (3, 4), \dots\}$  and  $\mathcal{E}_1 = \{(2, 3), (4, 5), \dots\}$ . We first compare the theoretical convergence rate (measured by  $c(\rho)$  in Theorem 2.3) of Algorithm 1 and the network lasso  $c(\rho)$  for the following three settings:

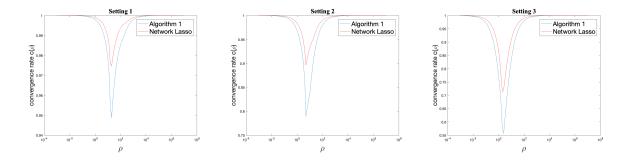


Figure 1: Comparison of the theoretical local convergence rates between Algorithm 1, the network lasso, and the standard lasso.

- 1.  $\mathbf{x}_{i}^{*} \in \mathbb{R}^{2}$ ,  $\hat{\mathbf{x}}_{i} \neq \hat{\mathbf{x}}_{i+1}$  when i = 50,  $\lambda = 1$ .
- 2.  $\mathbf{x}_{i}^{*} \in \mathbb{R}^{2}, \, \hat{\mathbf{x}}_{i} \neq \hat{\mathbf{x}}_{i+1} \text{ when } i = 10, 20, \dots, 90, \, \lambda = 1.$
- 3.  $\mathbf{x}_{i}^{*} \in \mathbb{R}^{2}$ ,  $\hat{\mathbf{x}}_{i} \neq \hat{\mathbf{x}}_{i+1}$  when  $i = 2, 4, \dots, 98, \lambda = 1$ .

The comparison of  $c(\rho)$  of Algorithm 1 and network lasso is visualized in Figure 1. We can see that Algorithm 1 consistently has a smaller  $c(\rho)$ , which implies a faster convergence rate.

## 2.6 Comparison with other works based on graph decomposition

Decomposing a graph into edges or paths is an idea that has been applied in existing works. However, we remark that our approach is different from previous works. For example, (Tansey and Scott, 2015) investigates the idea that decomposes the graph into a set of trails (this idea is also explored by Barbero and Sra (Barbero and Sra, 2014) for the two-dimensional grid graph), and then apply existing algorithms to solve each problem. In particular, it decomposes  $\mathcal{E}$  into K sets  $\mathcal{E}_1 \cup \cdots \cup \mathcal{E}_K$  such that for each  $1 \leq k \leq K$ ,  $\mathcal{E}_k$  is a trail. By writing the optimization problem as

$$\min_{\mathbf{x},\mathbf{z}} \sum_{i \in \mathcal{V}} f_i(\mathbf{x}_i) + \lambda \sum_{1 \le k \le K} \sum_{(s,t) \in \mathcal{E}_k} \|\mathbf{z}_{\mathcal{E}_k,s} - \mathbf{z}_{\mathcal{E}_k,t}\|,$$

s.t.  $\mathbf{z}_{\mathcal{E}_k,s} = \mathbf{x}_s$  for all  $1 \leq k \leq K$  and s in some edge of  $\mathcal{E}_k$ ,

then the ADMM algorithm can be used to update  $\mathbf{x}$  and  $\mathbf{z}$  alternatively. We remark that there are two main differences: first, their method only works well for the case where  $\mathbf{x}_i$  are scalars (i.e., p=1). In comparison, our method can handle the case where  $\mathbf{x}_i \in \mathbb{R}^p$  with p>1. Second and more importantly, the total variation penalty term in their method was not partitioned and it is addressed using the augmented variable  $\mathbf{z}$ ; while in our case, the total variation penalty term is partitioned and part of the  $\ell_1$  penalty was handled through the variable  $\mathbf{x}$ . In fact, the idea in this work can be combined with their idea for the case p=1 and the problem could be written as follows: first, we decompose  $\mathcal{E}$  into sets  $(\mathcal{E}_1^1 \cup \cdots \cup \mathcal{E}_{K_1}^1) \cup (\mathcal{E}_1^0 \cup \cdots \cup \mathcal{E}_{K_0}^0)$ , such that all  $\mathcal{E}_k^0$  and  $\mathcal{E}_k^1$  are trails, and the trails  $\mathcal{E}_1^0, \cdots, \mathcal{E}_{K_0}^0$  are disjoint. Then writing the optimization problem as

$$\min_{\mathbf{x}, \mathbf{z}} \left( \sum_{i \in \mathcal{V}} f_i(\mathbf{x}_i) + \lambda \sum_{1 \le k \le K_0} \sum_{(s,t) \in \mathcal{E}_k^0} \|\mathbf{x}_s - \mathbf{x}_t\| \right) + \lambda \sum_{1 \le k \le K_1} \sum_{(s,t) \in \mathcal{E}_k^1} \|\mathbf{z}_{\mathcal{E}_k^1, s} - \mathbf{z}_{\mathcal{E}_k^1, t}\|,$$
s.t.  $\mathbf{z}_{\mathcal{E}_k^1, s} = \mathbf{x}_s$  for all  $1 \le k \le K - 1$  and  $s$  in some edge of  $\mathcal{E}_k^1$ .

This formulation would give another algorithm for solving the problem in (Tansey and Scott, 2015), but we will leave it for possible future investigations.

## 3 Experiments

In this section, the proposed Algorithm 1 will be compared with network lasso for solving (1) under various scenarios. We measure their error at iteration k by the difference between its objective value and the optimal objective value. We remark that all ADMM algorithms require an augmented Lagrangian parameter  $\rho$ , and the algorithms would converge slowly when  $\rho$  is too large or too small. While there have been many works on the choice of  $\rho$ . For example, a simple varying penalty strategy based on residual balancing is suggested in Section 3.4.1 of (Boyd et al., 2011), and another choice based on the Barzilai-Borwein spectral method for gradient descent is proposed in there do not exists an optimal choice for settings (Xu et al., 2017). Considering that there is no general consensus on the optimal strategy of the choice of  $\rho$ , we will test the performance of the algorithms on a range of  $\rho$  in the following simulations, and we choose the range so that the optimal  $\rho$  is inside the

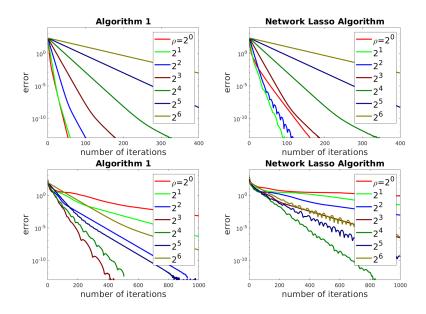


Figure 2: Comparison the convergence rates under the 1D chain graph setting with  $\lambda = 1$  (top row) and  $\lambda = 10$  (second row).

chosen range.

#### 3.1 Simulations

We first test our result on the one-dimensional chain graph. Following (Zhu, 2017), we use the model that

$$\mathbf{y}_{i}^{*} = \begin{cases} [1, 1], & \text{if } 1 \leq i \leq 11 \\ [-1, 1], & \text{if } 12 \leq i \leq 22 \\ [2, 2], & \text{if } 1 \leq 23 \leq 33 \\ [-1, -1], & \text{if } 34 \leq i \leq 44 \\ [0, 0], & \text{if } i \geq 45 \end{cases}.$$

Then we let  $\mathbf{y}_i = \mathbf{y}_i^* + N(\mathbf{0}, \mathbf{I}_{2\times 2})$  and  $\lambda = 1$  or 10, and compare Algorithm 1 will be compared with network lasso in Figure 2 with various choices of  $\rho$ . The figures indicate that for both choices of  $\lambda$ , Algorithm 1 always performs better with a good choice of  $\rho$ . In fact, if  $\rho$  is chosen to be the optimal values for both algorithms, Algorithm 1 converges twice as fast as network lasso.

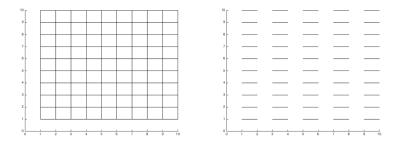


Figure 3: Visualization of a two-dimensional grid graph of size  $10 \times 10$  (left) and the corresponding  $\mathcal{E}_0$  for this graph (right).

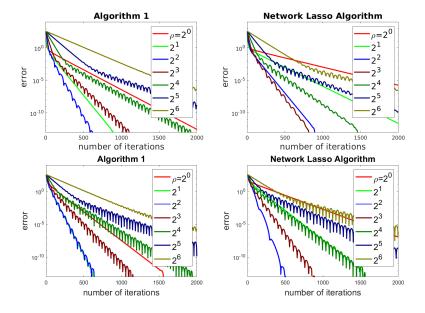


Figure 4: Comparison the convergence rates under the 2D grid graph setting with  $\lambda = 1$  (top row) and  $\lambda = 5$  (second row).

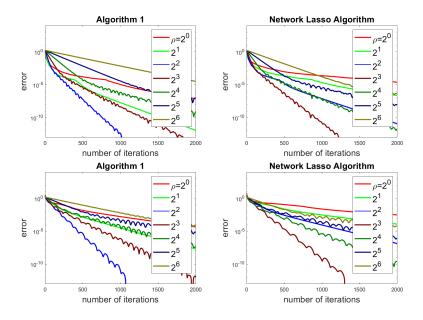


Figure 5: Comparison the convergence rates for the Chicago crime dataset with  $\lambda = 0.05$  (top row) and  $\lambda = 0.25$  (second row).

We also test the two-dimensional grid graph example. We generate data in the form of a 64 by 64 grid of points; the value is equal to [0,0,0] for points within a distance of 16 from the middle of the grid, and [0.4,0.7,1] for all other points, and add noise of  $N(\mathbf{0},\mathbf{I}_{3\times3})$  for all points. For this example, we have a natural choice of  $\mathcal{E}_0$  as visualized in the right Figure of Figure 3. The convergence rates are shown in Figure 4, which shows that Algorithm 1 has a comparable or faster convergence rate as the network lasso algorithm ADMM. Combining it with the fact that Algorithm 1 has a smaller computational complexity per iteration, this implies the numerical superiority of Algorithm 1.

#### 3.2 Real dataset

We also use an example of a graphical fused lasso problem with a reasonably large, geographically-defined underlying graph. The data comes from police reports made publicly available by the city of Chicago, from 2001 until the present (Chicago Police Department 2014) and is available as the supplementary files of (Arnold and Tibshirani, 2016). In this dataset, the vertices represent the census blocks and the edges represent the neighboring blocks, and more detailed explanation of this dataset is deferred to the supplementary file.

This graph has 2162 vertices, 6995 edges and by running the greedy algorithm as described in Section 2.2, we obtain  $\mathcal{E}_0$  with 1025 edges and  $\tilde{m}_0 = 797$ . The result of the experiment in Figure 5 shows when  $\lambda = 0.05$  or 0.25, Algorithm 1 converges much faster than the network lasso algorithm. Since Algorithm 1 has a smaller computational complexity per iteration as analyzed, this experiment shows that our algorithm numerically converges faster than network lasso.

## 4 Conclusion

This paper proposes a new ADMM algorithm for solving graphic fused lasso, based on a novel method of dividing the objective function into two components. Compared with the standard ADMM algorithm of network lasso, it has a similar complexity per iteration while usually converges within fewer iterations. As for future directions, it would be interesting to theoretically analyze its advantage, and explore other ways of dividing the objective function which could even future improve the performance of the ADMM algorithm for graph-fused lasso.

The idea of the proposed algorithm can also be applied to other problems such as trend filtering (Ramdas and Tibshirani, 2015), defined by  $\arg\min_{\mathbf{x}_i,i=1,\cdots,n} \sum_{i=1}^n f_i(\mathbf{x}_i) + \lambda \sum_{i=2}^{n-1} \|\mathbf{x}_{i-1} - 2\mathbf{x}_i + \mathbf{x}_{i+1}\|$ . We may divide the objective function into  $\sum_{i=1}^n f_i(\mathbf{x}_i) + \|\mathbf{x}_1 - 2\mathbf{x}_2 + \mathbf{x}_3\| + \|\mathbf{x}_4 - 2\mathbf{x}_5 + \mathbf{x}_6\| + \cdots$  and  $\|\mathbf{x}_2 - 2\mathbf{x}_3 + \mathbf{x}_4\| + \|\mathbf{x}_3 - 2\mathbf{x}_4 + \mathbf{x}_5\| + \|\mathbf{x}_5 - 2\mathbf{x}_6 + \mathbf{x}_7\| + \|\mathbf{x}_6 - 2\mathbf{x}_7 + \mathbf{x}_8\| + \cdots$ . The analysis of its performance and the comparison with standard algorithms would be another possible future direction.

## 5 Proofs

#### 5.1 Proof of Lemma 2.1

Let us consider the problem

$$\underset{\{\mathbf{x}_{i}\}_{i \in \mathcal{E}}, \{\mathbf{z}_{st}\}_{(s,t) \in \mathcal{E}_{1}}}{\operatorname{arg min}} \left( \sum_{i \in \mathcal{E}} f_{i}(\mathbf{x}_{i}) + \lambda \sum_{(s,t) \in \mathcal{E}_{0}} \|\mathbf{x}_{s} - \mathbf{x}_{t}\| \right) + \lambda \sum_{(s,t) \in \mathcal{E}_{1}} \|\mathbf{z}_{st}\|, \tag{20}$$
s.t.  $\mathbf{z}_{st} = \mathbf{x}_{s} - \mathbf{x}_{t}, \mathbf{z}_{ts} = \mathbf{x}_{s} + \mathbf{x}_{t}.$ 

and its associated Lagrangian of

$$\bar{L}_{\rho}(x, z, u) = \sum_{i \in \mathcal{V}} f_{i}(\mathbf{x}_{i}) + \lambda \sum_{(s, t) \in \mathcal{E}_{0}} \|\mathbf{x}_{s} - \mathbf{x}_{t}\| + \sum_{(s, t) \in \mathcal{E}_{1}} \left(\lambda \|\mathbf{z}_{st}\| + \mathbf{u}_{st}^{T}(\mathbf{z}_{st} - \mathbf{x}_{s} + \mathbf{x}_{t}) + \mathbf{u}_{ts}^{T}(\mathbf{z}_{ts} - \mathbf{x}_{s} - \mathbf{x}_{t}) + \frac{\rho}{2} \|\mathbf{z}_{st} - \mathbf{x}_{s} + \mathbf{x}_{t}\|^{2} + \frac{\rho}{2} \|\mathbf{z}_{ts} - \mathbf{x}_{s} - \mathbf{x}_{t}\|^{2}\right),$$
(21)

as well as the ADMM algorithm of

$$x^{(k+1)} = \arg\min_{x} \bar{L}_{\rho}(x, z^{(k)}, u^{(k)})$$
(23)

$$z^{(k+1)} = \arg\min_{z} \bar{L}_{\rho}(x^{(k+1)}, z, u^{(k)})$$
(24)

$$u_{st}^{(k+1)} = u_{st}^{(k)} + \rho(z_{st}^{(k+1)} - x_s^{(k+1)} + x_t^{(k+1)}), \ u_{ts}^{(k+1)} = u_{ts}^{(k)} + \rho(z_{ts}^{(k+1)} - x_s^{(k+1)} - x_t^{(k+1)}).$$
(25)

It can be verified that the update of (23)-(25) with  $\bar{L}_{\rho}$  is in fact identical to the update of (10)-(12) with  $\hat{L}_{2\rho}$ . In addition, using the fact that (20) is (14) with additional constraints  $\mathbf{z}_{ts} = \mathbf{x}_s + \mathbf{x}_t$ , and Lemma 5.1 implies the equivalence between (23)-(25) is identical to (17)-(19).

**Lemma 5.1.** The preconditioned ADMM procedure for solving  $\min_{\mathbf{x},\mathbf{y}} f(\mathbf{x}) + g(\mathbf{y})$  subject to  $\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} = \mathbf{c}$ 

$$\mathbf{x}^{(k+1)} = \arg\min_{\mathbf{x}} L(\mathbf{x}, \mathbf{y}^{(k)}, \mathbf{v}^{(k)}) + \frac{\rho}{2} (\mathbf{x} - \mathbf{x}^{(k)})^T \mathbf{C}_1^T \mathbf{C}_1 (\mathbf{x} - \mathbf{x}^{(k)}), \tag{26}$$

$$\mathbf{y}^{(k+1)} = \underset{\mathbf{y}}{\operatorname{arg min}} L(\mathbf{x}^{(k+1)}, \mathbf{y}, \mathbf{v}^{(k)}) + \frac{\rho}{2} (\mathbf{y} - \mathbf{y}^{(k)})^T \mathbf{C}_2^T \mathbf{C}_2 (\mathbf{y} - \mathbf{y}^{(k)}), \tag{27}$$

$$\mathbf{v}^{(k+1)} = \mathbf{v}^{(k)} + \rho(\mathbf{A}\mathbf{x}^{(k+1)} + \mathbf{B}\mathbf{y}^{(k+1)} - \mathbf{c}), \tag{28}$$

where  $L(\mathbf{x}, \mathbf{y}, \mathbf{v}) = f(\mathbf{x}) + g(\mathbf{y}) + \mathbf{v}^T (\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} - \mathbf{c}) + \frac{\rho}{2} ||\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} - \mathbf{c}||^2$ , is equivalent to the standard ADMM procedure applied to the augmented problem

$$\min_{\mathbf{x},\mathbf{y}} f(\mathbf{x}) + g(\mathbf{y}), \text{ s.t. } [\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y}, \mathbf{C}_1\mathbf{x}, \mathbf{C}_2\mathbf{y}] = [\mathbf{c}, \mathbf{z}, \mathbf{w}].$$

#### 5.2 Proof of Lemma 5.1

Proof of Lemma 5.1. Applying the standard ADMM routine to optimize  $(\mathbf{x}, \mathbf{w})$  and  $(\mathbf{y}, \mathbf{z})$  alternatively, the update formula for the augmented ADMM is

$$\mathbf{x}^{(k+1)} = \arg\min_{\mathbf{x}} L(\mathbf{x}, \mathbf{y}^{(k)}, \mathbf{v}^{(k)}) + \frac{\rho}{2} \|\mathbf{C}_{1}^{0.5} \mathbf{x} - \mathbf{z}^{(k)}\|^{2} + \mathbf{v}_{1}^{(k)T} (\mathbf{C}_{1}^{0.5} \mathbf{x} - \mathbf{z}^{(k)}), \tag{29}$$

$$\mathbf{w}^{(k+1)} = C_2^{0.5} \mathbf{y}^{(k)} + \frac{1}{\rho} \mathbf{v}_2^{(k)},\tag{30}$$

$$\mathbf{y}^{(k+1)} = \arg\min_{\mathbf{y}} L(\mathbf{x}^{(k+1)}, \mathbf{y}, \mathbf{v}^{(k)}) + \frac{\rho}{2} \|\mathbf{C}_{2}^{0.5} \mathbf{y} - \mathbf{w}^{(k+1)}\|^{2} + \mathbf{v}_{2}^{(k)T} (\mathbf{C}_{2}^{0.5} \mathbf{y} - \mathbf{w}^{(k+1)}), \quad (31)$$

$$\mathbf{v}^{(k+1)} = \mathbf{v}^{(k)} + \rho(\mathbf{A}\mathbf{x}^{(k+1)} + \mathbf{B}\mathbf{y}^{(k+1)} - \mathbf{c}), \tag{32}$$

$$\mathbf{z}^{(k+1)} = C_1^{0.5} \mathbf{x}^{(k+1)} + \frac{1}{\rho} \mathbf{v}_1^{(k)}$$
(33)

$$\mathbf{v}_{1}^{(k+1)} = \mathbf{v}_{1}^{(k)} + \rho(\mathbf{C}_{1}^{0.5}\mathbf{x}^{(k+1)} - \mathbf{z}^{(k+1)}), \ \mathbf{v}_{2}^{(k+1)} = \mathbf{v}_{2}^{(k)} + \rho(\mathbf{C}_{2}^{0.5}\mathbf{y}^{(k+1)} - \mathbf{w}^{(k+1)}).$$
(34)

Note that by plugging the definition of  $\mathbf{z}^{(k+1)}$  in (32) to the definition of  $\mathbf{v}^{(k+1)}$  (34),  $\mathbf{v}_1^{(k+1)} = 0$  for all k. So (32) implies that  $\mathbf{z}^{(k+1)} = C_1^{0.5} \mathbf{x}^{(k+1)}$  and (29) is equivalent to (26).

Plugging in the definition of  $\mathbf{w}^{(k+1)}$  to (31), we obtain the equivalence between (31) and (27).

#### 5.3 Proof of Theorem 2.3

Proof of Theorem 2.3. By calculation, the ADMM algorithm is equivalent to the Douglas-Rachford splitting method applied to

$$\max_{\mathbf{z}} -\mathbf{b}^T \mathbf{z} - f_1^* (-\mathbf{A}_1^T \mathbf{z}) - f_2^* (-\mathbf{A}_2^T \mathbf{z})$$

with two parts given by  $f = f_2^*(-\mathbf{A}_2^T\mathbf{z})$  and  $g = \mathbf{b}^T\mathbf{z} + f_1^*(-\mathbf{A}_1^T\mathbf{z})$  respectively, and the Douglas-Rachford splitting method is an iterative method that minimizing  $f(\mathbf{x}) + g(\mathbf{x})$  with update formula

$$\mathbf{x}^{(k+1)} = \frac{1}{2} [(\mathbf{I} - 2\operatorname{prox}_{\rho f})(\mathbf{I} - 2\operatorname{prox}_{\rho g}) + \mathbf{I}](\mathbf{x}^{(k)}).$$

Note that locally around the optimal solution we have

$$\operatorname{prox}_{\rho f} = (\mathbf{I} + \rho \partial f)^{-1}, \ \partial f(\mathbf{x}) = \mathbf{A}_2 \mathbf{C}_2^{-1} \mathbf{A}_2^T \mathbf{x} + \mathbf{c}_1,$$
$$\operatorname{prox}_{\rho g} = (\mathbf{I} + \rho \partial g)^{-1}, \ \partial g(\mathbf{x}) = \mathbf{A}_1 \mathbf{C}_1^{-1} \mathbf{A}_1^T \mathbf{x} + \mathbf{c}_2$$

for some  $c_1$  and  $c_2$ , each iteration of the algorithm is a linear operator in the form of

$$\frac{1}{2}[(\mathbf{I} - 2\operatorname{prox}_{\rho f})(\mathbf{I} - 2\operatorname{prox}_{\rho g}) + \mathbf{I}](\mathbf{x})$$

$$= \frac{1}{2}[(\mathbf{I} - 2(\mathbf{I} + \rho \mathbf{A}_2 \mathbf{C}_2^{-1} \mathbf{A}_2^T)^{-1})(\mathbf{I} - 2(\mathbf{I} + \rho \mathbf{A}_1 \mathbf{C}_1^{-1} \mathbf{A}_1^T)^{-1}) + \mathbf{I}](\mathbf{x}) + \mathbf{c}_0.$$

As a result, the algorithm converges in the order of

$$\left(\frac{1}{2}[(\mathbf{I} - 2(\mathbf{I} + \rho \mathbf{A}_2 \mathbf{C}_2^{-1} \mathbf{A}_2^T)^{-1})(\mathbf{I} - 2(\mathbf{I} + \rho \mathbf{A}_1 \mathbf{C}_1^{-1} \mathbf{A}_1^T)^{-1}) + \mathbf{I}]\right)^k,$$

where k is the number of iterations. The theorem is then proved.

## References

Arnold, T. B. and R. J. Tibshirani (2016). Efficient implementations of the generalized lasso dual path algorithm. *Journal of Computational and Graphical Statistics* 25(1), 1–27.

Barbero, A. and S. Sra (2014). Modular proximal optimization for multidimensional total-variation regularization.

Bleakley, K. and J.-P. Vert (2011, June). The group fused Lasso for multiple change-point detection. working paper or preprint.

Boyd, S., N. Parikh, E. Chu, B. Peleato, and J. Eckstein (2011, January). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.* 3(1), 1–122.

Chambolle, A. and J. Darbon (2009). On total variation minimization and surface evolution using parametric maximum flows. *International Journal of Computer Vision* 84(3).

- Chen, X., Q. Lin, S. Kim, J. G. Carbonell, and E. P. Xing (2012, 06). Smoothing proximal gradient method for general structured sparse regression. *Ann. Appl. Stat.* 6(2), 719–752.
- Condat, L. (2013, Nov). A direct algorithm for 1-d total variation denoising. *IEEE Signal Processing Letters* 20(11), 1054–1057.
- Davies, P. L. and A. Kovac (2001, 02). Local extremes, runs, strings and multiresolution.

  Ann. Statist. 29(1), 1–65.
- Eckstein, J. and D. P. Bertsekas (1992, Apr). On the douglas—rachford splitting method and the proximal point algorithm for maximal monotone operators. *Mathematical Programming* 55(1), 293–318.
- Friedman, J., T. Hastie, H. Hfling, and R. Tibshirani (2007, 12). Pathwise coordinate optimization. *Ann. Appl. Stat.* 1(2), 302–332.
- Hallac, D., J. Leskovec, and S. Boyd (2015, aug). Network Lasso: Clustering and Optimization in Large Graphs. *KDD*: proceedings. International Conference on Knowledge Discovery & Data Mining 2015, 387–396.
- Johnson, N. A. (2013). A dynamic programming algorithm for the fused lasso and 1 0-segmentation. *Journal of Computational and Graphical Statistics* 22(2), 246–260.
- Kolmogorov, V., T. Pock, and M. Rolinek (2016). Total variation on a tree. SIAM Journal on Imaging Sciences 9(2), 605–636.
- Kovac, A. and A. D. A. C. Smith (2011). Nonparametric regression on a graph. *Journal of Computational and Graphical Statistics* 20(2), 432–447.
- Landrieu, L. and G. Obozinski (2017). Cut pursuit: Fast algorithms to learn piecewise constant functions on general weighted graphs. SIAM Journal on Imaging Sciences 10(4), 1724–1766.
- Lin, X., M. Pham, and A. Ruszczynski (2014). Alternating linearization for structured regularization problems. 15, 3447–3481.

- Liu, J., L. Yuan, and J. Ye (2010). An efficient algorithm for a class of fused lasso problems. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '10, New York, NY, USA, pp. 323–332. ACM.
- Ramdas, A. and R. J. Tibshirani (2015, August). Fast and Flexible ADMM Algorithms for Trend Filtering. *Journal of Computational and Graphical Statistics* 25(3), 839–858.
- Rudin, L. I., S. Osher, and E. Fatemi (1992). Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena* 60(1), 259 268.
- Tansey, W. and J. G. Scott (2015, may). A Fast and Flexible Algorithm for the Graph-Fused Lasso.
- Tibshirani, R., M. Saunders, S. Rosset, J. Zhu, and K. Knight (2005). Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society Series B*, 91–108.
- Wahlberg, B., S. Boyd, M. Annergren, and Y. Wang (2012). An admm algorithm for a class of total variation regularized estimation problems. In *Preprints of the 16th IFAC Symposium on System Identification*, pp. 83–88. QC 20121112.
- Xu, Z., M. A. T. Figueiredo, and T. Goldstein (2017). Adaptive admm with spectral penalty parameter selection. In *AISTATS*.
- Ye, G.-B. and X. Xie (2011). Split bregman method for large scale fused lasso. Computational Statistics & Data Analysis 55(4), 1552 1569.
- Yu, D., J.-H. Won, T. Lee, J. Lim, and S. Yoon (2015). High-dimensional fused lasso regression using majorization-minimization and parallel processing. *Journal of Computational and Graphical Statistics* 24(1), 121–153.
- Zhu, Y. (2017). An augmented admm algorithm with application to the generalized lasso problem. *Journal of Computational and Graphical Statistics* 26(1), 195–204.