
Optimization of Graph Total Variation via Active-Set-based Combinatorial Reconditioning

Zhenzhang Ye
TU Munich
zhenzhang.ye@tum.de

Thomas Möllenhoff
TU Munich
thomas.moellenhoff@tum.de

Tao Wu
TU Munich
tao.wu@tum.de

Daniel Cremers
TU Munich
cremers@tum.de

Abstract

Structured convex optimization on weighted graphs finds numerous applications in machine learning and computer vision. In this work, we propose a novel adaptive preconditioning strategy for proximal algorithms on this problem class. Our preconditioner is driven by a sharp analysis of the local linear convergence rate depending on the “active set” at the current iterate. We show that nested-forest decomposition of the inactive edges yields a guaranteed local linear convergence rate. Further, we propose a practical greedy heuristic which realizes such nested decompositions and show in several numerical experiments that our reconditioning strategy, when applied to proximal gradient or primal-dual hybrid gradient algorithm, achieves competitive performances. Our results suggest that local convergence analysis can serve as a guideline for selecting variable metrics in proximal algorithms.

1 Introduction

Preconditioning, as a way of transforming a difficult linear system into one that is easier to solve, enjoys a rich and successful history. Recently, *proximal algorithms* [18, 53, 16] have received a surge of popularity in solving structured non-smooth convex optimization problems appearing across many fields in science and engineering. Unlike in the case of linear systems, putting forward a satisfactory theory and implementation of preconditioning in the general non-smooth setting remains a largely unsolved challenge [54, 26, 27, 40, 9, 21, 28, 16, 7, 45].

Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics (AISTATS) 2020, Palermo, Italy. PMLR: Volume 108. Copyright 2020 by the author(s).

This is mainly due to two obstacles:

(i) The non-linear dynamics of proximal algorithms, as well as the geometry of the non-smooth energy are more involved than in the quadratic case. A precise characterization of the convergence behavior, which could guide the proper choice of metric (preconditioner), is challenging.

(ii) In cases where the proper choice of metric is clear, non-diagonal preconditioners typically make the proximal operators in the algorithm much more expensive to evaluate. While favorably reducing the number of outer iterations, each inner iteration could be even of similar complexity as the original problem [40].

In this vein, numerous efforts have been devoted to a better understanding of the dynamics of proximal algorithms (see, e.g., [49, 24]), and exploring scenarios where non-diagonally scaled proximal mappings are still efficient to evaluate [22, 6, 7].

In this paper we take a novel perspective, circumventing issue (i) by resorting to the local convergence analysis. This does not yield provable guarantees on the global iteration complexity. Nevertheless, we show empirically that our preconditioners guided by the local analysis yield an improvement long before the local linear convergence regime is entered (see Fig. 3).

To overcome difficulty (ii) we restrict ourselves to structured convex problems on weighted graphs, where metrics based on tree decompositions are amenable to efficient proximal evaluation thanks to recent message-passing algorithms [36]. Specifically, given an undirected weighted graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \omega)$, whose edges are weighted by a function $\omega : \mathcal{E} \rightarrow \mathbb{R}_{>0}$, we consider the structured convex optimization on \mathcal{G} :

$$\min_{u \in \mathbb{R}^{\mathcal{V}}} G(u) + \text{TV}_{\mathcal{G}}(u), \quad (1)$$

where $\text{TV}_{\mathcal{G}}$ is the graph total variation

$$\text{TV}_{\mathcal{G}}(u) = \sum_{e=(i,j) \in \mathcal{E}} \omega_e |u_i - u_j|. \quad (2)$$

The function $G : \mathbb{R}^{\mathcal{V}} \rightarrow \mathbb{R} \cup \{+\infty\}$ is assumed to be proper, lower semi-continuous and convex.

We define the vertex-to-edge map $K : \mathbb{R}^{\mathcal{V}} \rightarrow \mathbb{R}^{\mathcal{E}}$ by

$$K = \text{diag}(\omega)\nabla,$$

where ∇ is the (transposed) incidence matrix of \mathcal{G} , i.e.,

$$(\nabla u)_e = u_i - u_j, \quad \forall e = (i, j) \in \mathcal{E},$$

with arbitrarily fixed orientation. With this notation we can succinctly write $\text{TV}_{\mathcal{G}}(u) = \|Ku\|_1$.

Problems of form (1) are, for example, relevant in image processing and computer vision [25, 44, 14, 48], unsupervised and transductive learning [30, 31, 10, 23], collaborative filtering [8] and clustering [23].

For separable convex $G(u) = \sum_{i \in \mathcal{V}} g_i(u_i)$, problem (1) can be efficiently solved (up to machine precision) in polynomial time by parametric max-flow methods [12, 32]. To handle non-separable but differentiable G , the authors in [62] propose a (primal) proximal gradient iteration, reducing (1) to a sequence of separable problems which are solved by parametric max-flow. For problems on regular grids, several authors proposed a splitting into *chains*, leading to 1D total variation subproblems which can be solved efficiently [19, 3, 36]. An active-set method for submodular minimization (which includes the graph total variation as a special case) was proposed in [38], which is different from the active-set strategy pursued here. Landrieu et al. recently proposed a fast method for graph total variation [39, 56] by assuming that the solution is piecewise constant and refining that partition by solving a sequence of max-flow problems. Closely related to the present approach are projected Newton methods [58], which have also been applied to the total variation [2].

In contrast, the main focus of this paper is to advance the understanding of preconditioning in proximal algorithms. To solve the problem class (1), we consider two types of algorithms:

(1) (Dual) proximal gradient (PG). Assume G^* is C^2 such that $l_{G^*}I \preceq \nabla^2 G^*(\cdot) \preceq L_{G^*}I$ for some constants $l_{G^*}, L_{G^*} > 0$. Based on the (Fenchel) dual formulation of (1), written

$$\min_{p \in \mathbb{R}^{\mathcal{E}}} G^*(-K^\top p) + \delta\{\|p\|_\infty \leq 1\}, \quad (3)$$

one can apply the proximal (or projected) gradient:

$$p^{k+1} = \arg \min_{p \in \mathbb{R}^{\mathcal{E}}} -\langle K\nabla G^*(-K^\top p^k), p \rangle + \delta\{\|p\|_\infty \leq 1\} + \frac{t}{2}\|p - p^k\|_{T_k}^2. \quad (4)$$

Here $T_k \in \mathbb{R}^{|\mathcal{E}| \times |\mathcal{E}|}$ is a symmetric positive definite matrix which induces a scaled norm $\|\cdot\|_{T_k}$ defined by $\|u\|_{T_k}^2 = \langle u, u \rangle_{T_k} = u^\top T_k u$.

(2) Primal-dual hybrid gradient (PDHG). Another equivalent formulation of (1) is the following convex-concave saddle-point problem:

$$\min_{u \in \mathbb{R}^{\mathcal{V}}} \max_{p \in \mathbb{R}^{\mathcal{E}}} \langle Ku, p \rangle + G(u) - \delta\{\|p\|_\infty \leq 1\}, \quad (5)$$

to which one can apply the primal-dual hybrid gradient (PDHG) algorithm:

$$\begin{aligned} u^{k+1} &= \arg \min_{u \in \mathbb{R}^{\mathcal{V}}} G(u) + \langle p^k, Ku \rangle + \frac{s}{2}\|u - u^k\|^2, \quad (6) \\ p^{k+1} &= \arg \min_{p \in \mathbb{R}^{\mathcal{E}}} -\langle K(2u^{k+1} - u^k), p \rangle \\ &\quad + \delta\{\|p\|_\infty \leq 1\} + \frac{t}{2}\|p - p^k\|_{T_k}^2. \quad (7) \end{aligned}$$

1.1 Related work on preconditioning

The (vanilla) PG and PDHG (typically with $T_k \equiv I$), as special instances of proximal algorithms, are widely applied in convex optimization – we refer to [18, 53, 16] for up-to-date surveys which contain relevant historical accounts and interconnection of algorithms. Acceleration of these algorithms is of significant research as well as practical interests. To this end, momentum-based acceleration techniques, which are traced back to the seminal works by Nesterov [47] and Polyak [55], were recently developed for PG [5, 52] and PDHG [15] and achieved impressive performances [16].

In contrast to momentum methods, preconditioning techniques for proximal algorithms are less developed and understood, as previously discussed in the introduction. To clarify further, in the context of proximal methods there are roughly two separate streams of ideas referred to as preconditioning.

In the first one, the aim is to make the individual update steps in the algorithm easier while retaining a convergent method [9, 14]. While making each iteration faster, the effect on the overall complexity is unclear.

The second line of works, aims at improving the theoretical convergence rate and thereby reducing the number of outer iterations, see [26, 27, 28]. However, these works make very restrictive assumptions on the problem class and do not apply to our setting. A consensus among these works is to minimize the (*finite*) condition number $\kappa(T^{-1/2}K)$, which is defined by

$$\kappa(\cdot) = \frac{\sigma_{\max}(\cdot)}{\sigma_{\min>0}(\cdot)}, \quad (8)$$

as a reasonable heuristic in practice. This was (approximately) pursued for general problems in [54, 21,

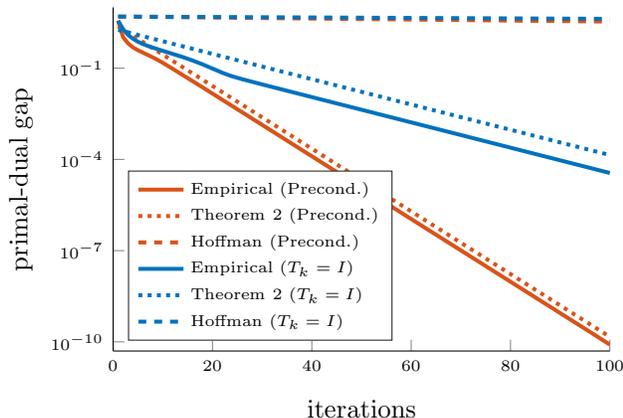


Figure 1: Local vs global analysis of the linear convergence of the PG iteration (4). The local linear rate sharply matches the observed convergence behaviour, while the global rate based on Hoffman’s bound is not informative. We guide the construction of our preconditioner based on the local convergence theory.

20]. In particular, forest-structured preconditioners for $K = \text{diag}(\omega)\nabla$ which are provably optimal in terms of $\kappa(T^{-1/2}K)$ were proposed recently in [45].

2 Local convergence analysis

While the condition number $\kappa(T^{-1/2}K)$ has proven to be reasonable heuristic in practice, a more quantified connection between the convergence rate and the preconditioner T would be desirable.

For problems of form (1), global linear convergence of PG (4) can be established using Hoffman’s bound [33, 35, 46, 34]. However, the linear rate obtained from that bound is mainly of theoretical interest, as it does not really inform us about the practical performance of the method but rather gives a (weak) upper bound. Secondly, Hoffman’s bound is an inherently combinatorial expression that is very challenging to compute even for small problem instances.

Instead, we aim to choose the preconditioner to improve the local convergence behaviour of the method. It turns out that for a wide range of *partly smooth* functions the local dynamics of the PG, PDHG and accelerated variants thereof are well understood, see [41, 42, 43]. This will serve as a basis for our theory.

In Fig. 1 we show the linear rate predicted by Hoffman’s bound to the local rate on a small 4×3 grid graph for which Hoffman’s bound is still tractable to compute. As discussed above, the global rate by Hoffman’s bound is not informative. The local analysis we present in Theorem 2 below (which proceeds similar to [41]) is sharp, matches the empirical performance

and will be the guide of our preconditioners.

Next we establish the local linear convergence of (4). Our strategy is to prove that, locally, iteration (4) reduces to gradient descent on a modified unconstrained problem. Therefore, the local linear rate is inherited from the one of gradient descent.

Lemma 1. *Let h be C^2 with $l_h I \preceq \nabla^2 h(\cdot) \preceq L_h I$ for some constants $l_h, L_h > 0$. Then the gradient descent on $\min_x h(Ax + b)$ with step size $1/t = 2/(L_h \sigma_{\max}(A)^2 + l_h \sigma_{\min>0}(A)^2)$ satisfies*

$$\|x^{k+1} - x^*\| \leq \frac{\varphi - 1}{\varphi + 1} \|x^k - x^*\|, \quad (9)$$

with $\varphi = \kappa(A)^2 \cdot \kappa(h)$, $\kappa(h) := L_h/l_h$.

Proof. See the supplementary material. \square

The analysis in Theorem 2 below hinges on finite identification of the *active set* define as

$$\mathcal{A}(p) = \{e \in \mathcal{E} : |p_e| = 1\}. \quad (10)$$

The associated projection matrix is defined as

$$(P_{\mathcal{A}}p)_e = \begin{cases} p_e & \text{if } e \in \mathcal{A}, \\ 0 & \text{if } e \notin \mathcal{A}. \end{cases} \quad (11)$$

Correspondingly, let $\mathcal{I}(p) := \mathcal{E} \setminus \mathcal{A}(p)$ be the *inactive set* and $P_{\mathcal{I}} := I - P_{\mathcal{A}}$.

Theorem 2. *Suppose that (4) generates a sequence $\{p^k\}$ which converges to a minimizer $p^* \in \mathbb{R}^{\mathcal{E}}$ of (3). Under the assumptions that*

$$(A1) \text{ For each } e \in \mathcal{E}, (K\nabla G^*(-K^\top p^*))_e = 0 \Rightarrow |p_e^*| < 1;$$

$$(A2) \text{ For each } k \in \mathbb{N}, \underline{t}I \preceq T_k \preceq \bar{t}I \text{ with fixed } \underline{t}, \bar{t} > 0;$$

$$(A3) T_k \text{ depends on } p^k \text{ only through } \mathcal{A}(p^k);$$

there exists $\bar{k} \in \mathbb{N}$ such that for all $k \geq \bar{k}$:

(i) *Finite identification, i.e.,*

$$\mathcal{A}(p^k) = \mathcal{A}(p^*) \equiv \mathcal{A}^*, T_k \equiv T. \quad (12)$$

(ii) *Local linear convergence, i.e.,*

$$\|p^k - p^*\|_T \leq \left(\frac{\varphi - 1}{\varphi + 1}\right)^{k - \bar{k}} \|p^{\bar{k}} - p^*\|_T, \quad (13)$$

with

$$\varphi = \kappa(\Pi_{U(\mathcal{A}^*)} T^{-1/2} K)^2 \cdot \kappa(G^*), \quad (14)$$

and $\Pi_{U(\mathcal{A}^*)}$ the orthogonal projection onto the subspace $U(\mathcal{A}^*) := \ker(P_{\mathcal{A}^*} T^{-1/2})$.

Proof. (i) Finite identification of the active set follows by invoking [11, Corollary 3.6]. The strict complementary condition at p^* required by that corollary is (A1). Further, the corollary requires

$$\text{dist}(0, \nabla J(p^k) + N(p^k)) \rightarrow 0, \quad (15)$$

where $J = G^* \circ (-K^\top)$ and

$$N(\bar{p}) = \left\{ p \in \mathbb{R}^{\mathcal{E}} : p_e = 0 \text{ if } e \notin \mathcal{A}(\bar{p}), \right. \\ \left. \text{sgn}(\bar{p}_e) \cdot p_e \geq 0 \text{ if } e \in \mathcal{A}(\bar{p}) \right\}, \quad (16)$$

denotes the normal cone at \bar{p} . From the optimality conditions of (4) it follows

$$tT_k(p^k - p^{k+1}) - (\nabla J(p^k) - \nabla J(p^{k+1})) \\ \in \nabla J(p^{k+1}) + N(p^{k+1}). \quad (17)$$

Then we have

$$\text{dist}(0, \nabla J(p^{k+1}) + N(p^{k+1})) \\ \leq \|tT_k(p^k - p^{k+1}) - (\nabla J(p^k) - \nabla J(p^{k+1}))\| \\ \leq (t\|T_k\| + L_{G^*} \lambda_{\max}(K^\top K)) \|p^k - p^{k+1}\|. \quad (18)$$

Convergence of $\{p^k\}$ to p^* implies $\|p^k - p^{k+1}\| \rightarrow 0$ and (15) follows by (A2). Since the active set is constant for $k \geq \bar{k}$ we have by (A3) that $T_k \equiv T$.

(ii). Assume in the following that $k \geq \bar{k}$. Since $T_k = T$ due to (i), (4) is equivalent the projected gradient descent applied to

$$\min_{q \in \mathbb{R}^{\mathcal{E}}} \tilde{J}(q) \quad \text{s.t. } \|T^{-1/2}q\|_\infty \leq 1, \quad (19)$$

under the change of variable $p = T^{-1/2}q$, $\tilde{J} = J \circ T^{-1/2}$. The iteration in q is given by

$$q^{k+1} = \arg \min_{\|T^{-1/2}q\|_\infty \leq 1} \langle \nabla \tilde{J}(q^k), q \rangle + \frac{t}{2} \|q - q^k\|^2, \quad (20)$$

whose optimality condition reads

$$t(q^k - q^{k+1}) \in T^{-1/2}N(T^{-1/2}q^{k+1}) + \nabla \tilde{J}(q^k). \quad (21)$$

From (i) we know that $P_{\mathcal{A}^*}p^{k+1} = P_{\mathcal{A}^*}p^k$, which yields

$$P_{\mathcal{A}^*}T^{-1/2}q^{k+1} = P_{\mathcal{A}^*}T^{-1/2}q^k, \\ \Rightarrow q^{k+1} - q^k \in \ker(P_{\mathcal{A}^*}T^{-1/2}) = U(\mathcal{A}^*). \quad (22)$$

In addition, in view of (16) we have

$$T^{-1/2}N(T^{-1/2}q^k) \subset U(\mathcal{A}^*)^\perp. \quad (23)$$

Thus, applying $\Pi_{U(\mathcal{A}^*)}$ on both sides of (21) yields an *equivalent* characterization:

$$0 = \Pi_{U(\mathcal{A}^*)} \nabla \tilde{J}(q^k) + t(q^{k+1} - q^k). \quad (24)$$

Indeed, this is the gradient descent on \tilde{J} restricted to $U(\mathcal{A}^*)$, which we rewrite as

$$q^{k+1} = q^k + t^{-1} \Pi_{U(\mathcal{A}^*)} T^{-1/2} K^\top \nabla G^* (-K^\top T^{-1/2} q^k) \\ = q^k + t^{-1} \Pi_{U(\mathcal{A}^*)} T^{-1/2} K^\top \nabla G^* (-K^\top T^{-1/2} \\ (\Pi_{U(\mathcal{A}^*)} q^k + \Pi_{U(\mathcal{A}^*)^\perp} q^{\bar{k}})). \quad (25)$$

Hence (20) is equivalent to gradient descent on the function $G^* \circ (A \cdot + b)$ with $A = -K^\top T^{-1/2} \Pi_{U(\mathcal{A}^*)}$, $b = -K^\top T^{-1/2} \Pi_{U(\mathcal{A}^*)^\perp} p^{\bar{k}}$. Using Lemma 6 yields the linear convergence in $\{q^k\}$. As $\|q^k\| = \|T^{1/2}p^k\| = \|p^k\|_T$, we achieve the linear convergence, with respect to the T -norm, of the original sequence $\{p^k\}$. \square

Corollary 3. *Let φ be given as in (14). Locally (i.e., for $k \geq \bar{k}$), with fixed $T \equiv T_{\bar{k}}$ we have $\|p^k - p^*\| \leq \varepsilon$ whenever*

$$k \geq \bar{k} + \frac{\varphi + 1}{2} \log \left(\frac{\|p^{\bar{k}} - p^*\| \sqrt{\kappa(T)}}{\varepsilon} \right). \quad (26)$$

We remark that there are bounds in literature on \bar{k} , see [42, Prop. 3.6] or the recent works [51, 50]. Analyzing which choice of variable metric T_k lead to fast identification of \mathcal{A}^* is beyond the scope of this work.

3 Combinatorial preconditioner

Suggested by the local convergence analysis and Corollary 3 from the previous section, an ideal preconditioner T ought to minimize the condition number $\kappa(\Pi_{U(\mathcal{A}^*)} T^{-1/2} K)$ once the active set \mathcal{A}^* is identified. In practice, however, computationally amenable choices of T are rather constrained due to a generic *trade-off* between convergence speed of (outer) iterations and per-iteration cost, i.e., the T -scaled proximal evaluation in (4) or (7). A dense matrix T , in general, will render inner iterations very expensive, as in the case of proximal Newton method [40]. For this reason, many authors consider diagonal preconditioners [54, 26, 27] or (diagonal + low-rank) preconditioners [6, 7] to keep the inner iterations fast and tractable.

Towards yet better balance of this trade-off, a recent paper [45] makes use of fast TV solver on trees [19, 36] and proposes a class of block diagonal preconditioners via graph partitioning (aiming at optimizing $\kappa(T^{-1/2}K)$ heuristically, however). There the optimal condition number $\kappa(T^{-1/2}K)$ is achieved by matroid partitioning. As a remark, combinatorial preconditioners for solving linear systems involving graph Laplacians date back to the early work by Vaidya in 1990's [61]; refer to [59] for a more detailed survey.

In this section, we construct combinatorial preconditioners which are more faithful, compared to the ones

from [45], to the (local) convergence analysis. In a nutshell, given the current active/inactive sets of edges, we partition the graph into *inactively nested forests* in the sense of (37), so that the resulting preconditioner yields a guaranteed (local) convergence rate, which is made precise in Theorem 5.

To construct our preconditioner, let the edge set \mathcal{E} be partitioned into L mutually disjoint subsets, i.e., $\mathcal{E} = \bigsqcup_{l=1}^L \mathcal{E}_l$, such that each subgraph $\mathcal{G}_l = (\mathcal{V}, \mathcal{E}_l, \omega|_{\mathcal{E}_l})$ is a *forest*. Correspondingly, we define P_l as the canonical projection from $\mathbb{R}^{\mathcal{E}}$ to $\mathbb{R}^{\mathcal{E}_l}$, i.e., $P_l p = p|_{\mathcal{E}_l}$ for any $p \in \mathbb{R}^{\mathcal{E}}$. Thus, the matrix K can be decomposed into submatrices $\{K_l\}_{l=1}^L$ where each $K_l = P_l K \in \mathbb{R}^{|\mathcal{E}_l| \times |\mathcal{V}|}$. Analogously, let $\nabla_l = P_l \nabla$. Note that each ∇_l^\top (or K_l^\top) has full column rank, and hence

$$T_l := K_l K_l^\top, \quad \forall l \in \{1, \dots, L\}, \quad (27)$$

is symmetric positive definite.

We then define our preconditioner as

$$T := \sum_{l=1}^L P_l^\top T_l P_l. \quad (28)$$

In view of Theorem 2, we analyze in the following the condition number of the following matrix:

$$\begin{aligned} \Pi_{\mathcal{I}} &:= K^\top T^{-1/2} \Pi_{U(\mathcal{A})} T^{-1/2} K \\ &= K^\top T^{-1/2} (I - T^{-1/2} P_{\mathcal{A}} (T^{-1/2} P_{\mathcal{A}})^\dagger) T^{-1/2} K. \end{aligned} \quad (29)$$

As a preparatory result, the following lemma decomposes $\Pi_{\mathcal{I}}$ into orthogonal projections onto subspaces.

Lemma 4. *Given $\mathcal{E} = \mathcal{A} \sqcup \mathcal{I}$, let \mathcal{G} be partitioned into L nonempty forests $\{\mathcal{G}_l\}_{l=1}^L$. Then the matrix defined in (29) can be characterized as*

$$\Pi_{\mathcal{I}} = \sum_{l=1}^L \Pi_{\mathcal{I},l}, \quad (30)$$

where each $\Pi_{\mathcal{I},l}$ is the orthogonal projection onto the linear subspace $\mathcal{S}_{\mathcal{I},l}$ defined by

$$\mathcal{S}_{\mathcal{I},l} := \text{span}\{\nabla_e^\top : e \in \mathcal{I} \cap \mathcal{E}_l\}. \quad (31)$$

Proof. (i) We show the identity (30) with $I_l := P_l I P_l^\top$, $P_{\mathcal{A},l} := P_l P_{\mathcal{A}} P_l^\top$, $P_{\mathcal{I},l} := P_l P_{\mathcal{I}} P_l^\top$, and

$$\begin{aligned} \Pi_{\mathcal{I},l} &:= K_l^\top T_l^{-1/2} (I_l - (T_l^{-1/2} P_{\mathcal{A},l}) (T_l^{-1/2} P_{\mathcal{A},l})^\dagger) \\ &\quad T_l^{-1/2} K_l. \end{aligned} \quad (32)$$

Note that

$$\begin{aligned} K^\top T^{-1/2} &= \sum_{l=1}^L K^\top P_l^\top T_l^{-1/2} P_l \\ &= \sum_{l=1}^L K_l^\top T_l^{-1/2} P_l, \end{aligned} \quad (33)$$

$$\begin{aligned} T^{-1/2} P_{\mathcal{A}} &= \left(\sum_{l=1}^L P_l^\top T_l^{-1/2} P_l \right) \left(\sum_{l'=1}^L P_{l'}^\top P_{\mathcal{A},l'} P_{l'} \right) \\ &= \sum_{l=1}^L P_l^\top T_l^{-1/2} P_{\mathcal{A},l} P_l, \end{aligned} \quad (34)$$

$$(T^{-1/2} P_{\mathcal{A}})^\dagger = \sum_{l=1}^L P_l^\top (T_l^{-1/2} P_{\mathcal{A},l})^\dagger P_l. \quad (35)$$

By plugging (33)–(35) into (29), we accomplish (i).

(ii) We show each $\Pi_{\mathcal{I},l}$ is the orthogonal projection onto $\mathcal{S}_{\mathcal{I},l}$. First, it is easy to see $\Pi_{\mathcal{I},l}$ is symmetric and $\Pi_{\mathcal{I},l}^2 = \Pi_{\mathcal{I},l}$, and hence an orthogonal projection. Secondly, note that $\text{rank } \Pi_{\mathcal{I},l} = |\mathcal{I} \cap \mathcal{E}_l| = \text{rank } K_l^\top P_{\mathcal{I},l}$. Furthermore, we have the following equation:

$$\begin{aligned} \Pi_{\mathcal{I},l} K_l^\top P_{\mathcal{I},l} &= K_l^\top P_{\mathcal{I},l} - K_l^\top T_l^{-1/2} \\ &\quad (T_l^{-1/2} P_{\mathcal{A},l}) (T_l^{-1/2} P_{\mathcal{A},l})^\dagger T_l^{1/2} P_{\mathcal{I},l} \\ &= K_l^\top P_{\mathcal{I},l}, \end{aligned} \quad (36)$$

which completes step (ii). \square

Theorem 5. *Given $\mathcal{E} = \mathcal{A} \sqcup \mathcal{I}$, let \mathcal{G} be partitioned into L nonempty, inactively nested forests $\{\mathcal{G}_l\}_{l=1}^L$ in the sense that*

$$\mathcal{S}_{\mathcal{I},1} = \dots = \mathcal{S}_{\mathcal{I},\hat{l}} \supseteq \mathcal{S}_{\mathcal{I},\hat{l}+1} \supseteq \dots \supseteq \mathcal{S}_{\mathcal{I},L} \supseteq \{0\}, \quad (37)$$

with the subspaces defined in (31). Then we have $\lambda_{\min>0}(\Pi_{\mathcal{I}}) = \hat{l}$ and the (local) convergence rate in Theorem 2 is $\varphi = (L/\hat{l}) \cdot \kappa(G^*)$.

Proof. By Lemma 4, we have $\lambda_{\max}(\Pi_{\mathcal{I}}) \leq \sum_{l=1}^L \lambda_{\max}(\Pi_{\mathcal{I},l}) \leq L$. In fact, the equality holds since $\Pi_{\mathcal{I}} v = Lv$ for some nonzero $v \in \mathcal{S}_{\mathcal{I},L}$. On the other hand, for any $v \in \text{ran } \Pi_{\mathcal{I}}$, we have $\langle v, \Pi_{\mathcal{I}} v \rangle \geq \sum_{l=1}^{\hat{l}} \langle v, \Pi_{\mathcal{I},l} v \rangle = \hat{l} \|v\|^2$. The equality holds for some nonzero $v \in \text{ran } \Pi_{\mathcal{I}} \cap (\mathcal{S}_{\mathcal{I},\hat{l}+1})^\perp$. This yields $\lambda_{\min>0}(\Pi_{\mathcal{I}}) = \hat{l}$. \square

4 Implementation

In this section, we specify how to construct our preconditioner and apply active-set-based reconditioning to both PG and PDHG algorithms. We only trigger reconditioning every n iterations. When reconditioning

is performed at iteration k , a greedy heuristic is used for constructing the preconditioner T_k ; see Section 4.1. Then, using the separability of $\|\cdot\|_\infty$, we perform the updates across the subgraphs $\{\mathcal{G}_l\}_{l=1}^L$:

$$p^{k+1}|_{\mathcal{E}_l} = \arg \min_{p \in \mathbb{R}^{\mathcal{E}_l}} -\langle K\bar{u}^k|_{\mathcal{E}_l}, p \rangle + \delta\{\|p\|_\infty \leq 1\} + \frac{t}{2}\|p - p^k\|_{T_{k,l}}^2, \quad (38)$$

where \bar{u}^k is defined as:

$$\bar{u}^k = \begin{cases} \nabla G^*(-K^\top p^k), & \text{for PG,} \\ 2u^{k+1} - u^k, & \text{for PDHG.} \end{cases} \quad (39)$$

The proximal evaluation required by (38) is detailed in Section 4.2, which invokes the message-passing algorithm on trees. The overall complexity of the reconditioned algorithm is discussed in Section 4.3.

4.1 Constructing preconditioner

Following Theorem 5 we aim to find a preconditioner T_k which minimizes the condition number $\varphi = (L/\hat{l}) \cdot \kappa(G^*)$, and hence the local linear convergence rate. Theoretically, optimal T_k can be found in polynomial time by Matroid partitioning as in [45]. The computation time is prohibitively large for the graphs in practical problems, however. Here we present a greedy heuristic to find inactive nested forests.

Given an input graph \mathcal{G} we partition the graph based on the active set at the current dual variable p^k . We assign to each edge $e \in \mathcal{E}$ an additional weight $\rho_e = 1 - |1 - |p_e^k||$. Then, a minimum spanning forest according to that weight is generated using Kruskal's algorithm [37]. This spanning forest is then subtracted from current graph and will be added to the set $\{\mathcal{G}_l\}_{l=1}^L$. We perform this generation and subtraction iteratively until no edges remain in the original graph.

The partitioning weight is introduced for two reasons: Firstly, we found it unstable to determine the active set $\mathcal{A}(p^k)$ numerically according to a threshold; Secondly, computing the preconditioner T_k is quite expensive for large graphs. This strategy could extend the suitable duration of current preconditioner since a potential active edge often has a larger partitioning weight.

4.2 Backward solver

The introduction of the proposed preconditioner T_k makes the backward update (38) more expensive. Here we describe how to solve it efficiently, following the approach in [45]. Combining the linear and the quadratic term, (38) can be re-written as:

$$p^{k+1}|_{\mathcal{E}_l} = \arg \min_{\|p\|_\infty \leq 1} \frac{1}{2}\|K_l^\top p + f_l\|^2, \quad (40)$$

where $f_l = -K_l^\top p^k|_{\mathcal{E}_l} - \bar{u}^k/t$. The (Fenchel) dual problem of (40) is given by

$$v_l = \arg \min_{u \in \mathbb{R}^{\mathcal{V}}} \frac{1}{2}\|u - f_l\|^2 + \|K_l u\|_1, \quad (41)$$

which is simply a weighted total variation problem on the individual trees in the forest \mathcal{G}_l . We solve the problem (41) using the message-passing algorithm introduced in [36]. To retrieve $p^{k+1}|_{\mathcal{E}_l}$ from v_l one can use the optimality condition:

$$K_l^\top p^{k+1}|_{\mathcal{E}_l} = v_l - f_l. \quad (42)$$

4.3 Discussion on complexity

For non-preconditioned proximal gradient, the complexity of each iteration is $\mathcal{O}(|\mathcal{E}|)$. For the preconditioned variant it is $\mathcal{O}(\sum_{t=1}^{\mathcal{T}} |\mathcal{E}_t| \log(|\mathcal{E}_t|))$ where \mathcal{T} is the total number of trees using the aforementioned message-passing algorithm [36]. The preconditioned update can still be parallelized to some extent, as the message-passing can run for each tree in parallel.

Construction of the preconditioner T_k based on the greedy inactive nested forest strategy with Prim's or Kruskal's algorithm is $\mathcal{O}(|\mathcal{E}|^2 \log(|\mathcal{E}|)/|\mathcal{V}|)$ [17].

After entering the local linear convergence phase, the overall iteration complexity is $\mathcal{O}(\varphi \log(1/\varepsilon))$ to find an ε -accurate solution (see Corollary 3). While each iteration of the preconditioned algorithm is slightly more costly (by roughly a factor of $\log(|\mathcal{E}|)$), the condition number φ is drastically reduced. For regular grids we have that $\varphi \in \mathcal{O}(|\mathcal{V}|)$ (cf. [45, Theorem 4]) in the non-preconditioned case. The proposed preconditioner improves this to a *constant* $\varphi \in \mathcal{O}(1)$, independent of problem size at the expense of a slightly more expensive dual update step (up to a logarithmic factor).

5 Applications

In the following experiments we compare four preconditioning strategies: non-preconditioned $T_k = I$, diagonally scaled $T_k = \text{diag}(KK^\top)$, nested (linear) forest from [45] and the

5.1 Numerical validation on synthetic data

As a first numerical example, we consider the fused Lasso [60] (also called ROF model in imaging [57])

$$\min_{u \in \mathbb{R}^{\mathcal{V}}} \frac{1}{2}\|u - f\|^2 + \|Ku\|_1. \quad (43)$$

We solve (43) on random graphs with fixed $|\mathcal{V}| = 512$ using proximal gradient (PG) with f chosen uniformly

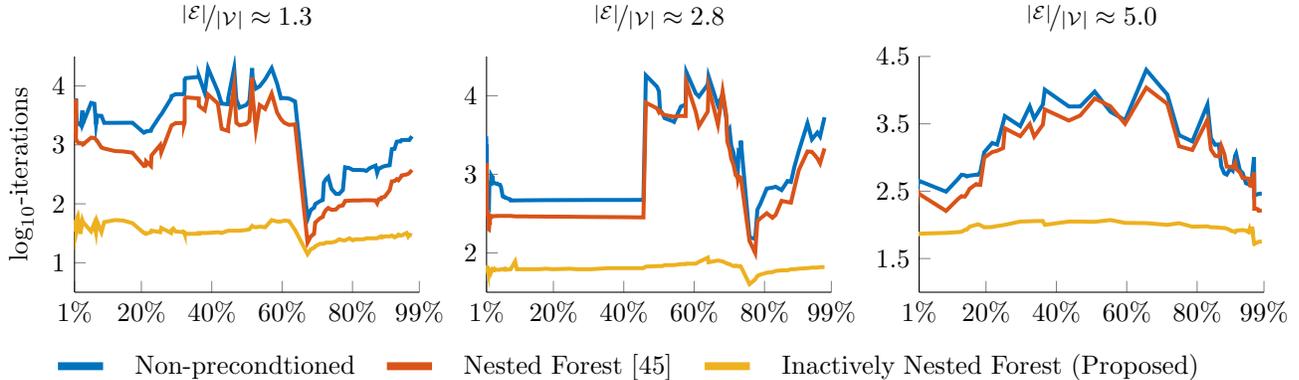


Figure 2: We show \log_{10} -iterations required by PG to reach a primal-dual gap smaller than 10^{-10} over percentage of active edges at the optimal solution for random graphs with varying edge-to-vertex ratio. The reconditioning strategy requires several orders of magnitude less iterations than no preconditioner and the preconditioner [45].

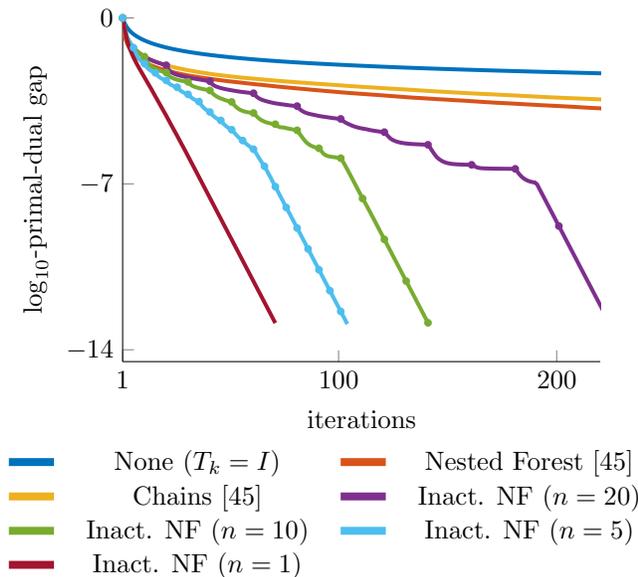


Figure 3: We show \log_{10} -primal dual gap vs iterations for PG (4) with various choices of T_k . The non-preconditioned choice $T_k = I$ performs the worst, followed by the preconditioners proposed in [45]. We indicate reconditioning by a dot and carry it out every n iterations for $n \in \{20, 10, 5, 1\}$. Smaller n leads to an increasingly improved performance.

random in $[0, 1]$. We consider two factors: edge-to-vertex ratio and percentage of active edges at the optimal solution. For the proposed reconditioning we set the frequency to $n = 1$.

The results are shown in Fig. 2. For reasonable amounts of active edges at the solution (30% – 80%) the proposed preconditioning strategy requires orders of magnitude less iterations to reach a primal-dual gap

under 10^{-10} . Moreover, it is shown that we require the fewest iterations across all scenarios.

In Fig. 3 we show \log_{10} -primal-dual gap over iterations for PG applied to (43) on a 100×100 grid graph with different choices of T_k and moderate regularization strength (30% of active edges at the optimal solution). The proposed preconditioner outperforms vanilla PG ($T_k = I$) and the recent (fixed) preconditioners proposed in [45]. Reconditioning more often leads to faster convergence, but as recomputing the preconditioner is expensive there is a trade-off between reducing the number of iterations and fast updates. In practice, a choice of the reconditioning frequency n between 5 and 30 leads to the best performance.

5.2 Fused Lasso on real-world graphs

To consider a more realistic scenario, we solve the model (43) on real-world graphs from a popular graph-cut benchmark considered in [29]. Furthermore, instead of using standard PG we used the accelerated FISTA variant [13, 1, 42] with overrelaxation parameter $\beta_k = (k - 1)/(k + 2)$. Reconditioning takes place at every 30 iterations. We discard the momentum for one iteration after reconditioning, which improved the stability. In Table 1, we show the running time and number of iterations of FISTA with non-preconditioned, diagonal preconditioner, nested forest [45], linear forest [45] and the proposed inactively nested forest.

Our preconditioner outperforms the other methods in all cases on number of iterations, despite a rather large choice of $n = 30$. However, the linear forest from [45] perform better with respect to the running time on 5 out of 10 datasets. The two datasets with grid structures leads to chain partition on which message-passing is much faster than on trees. The sizes of last

Instance name	$\frac{ A_* }{ \mathcal{E} }$	None		Diagonal		Nest. Forest		Lin. Forest		Inact. NF	
		it[10 ³]	time[s]	it[10 ³]	time[s]						
rmf-long	0.02	–	–	19	473	12	2539	18	233.3	1.9	474.9
rmf-wide	0.19	–	–	62	665	27	2274	43	213.1	0.19	18.54
horse	0.02	–	–	–	–	2.9	340.8	37	355.6	0.73	155.3
alue	0.03	–	–	–	–	4.5	117.2	100	270.9	0.71	155.3
lux	0.01	–	–	–	–	13	1254	–	–	0.40	54.34
punch	0.01	488	968	–	–	14.9	1445	203	872.1	0.34	62.66
BVZ*	0.35	27	74.0	22	730	1.12	434.8	0.57	6.59	0.49	261
manga*	0.05	–	–	–	–	38	48591	5.3	230	1.41	4609
KZ2	0.5	419	3042	1.6	159	0.43	614.1	0.6	56.0	0.42	965.9
ferro	0.09	9.25	186	5.83	639.6	0.36	430.3	0.93	86.23	0.27	609.3

Table 1: We show the number of iterations and running time to reach a relative primal dual gap less than 10^{-10} on (43) on real-world graphs. FISTA with various choices of T_k is used to solve these problems. “–” means the algorithm failed to reach the tolerance within 5×10^5 iterations. “*” means that graph has a grid structure.

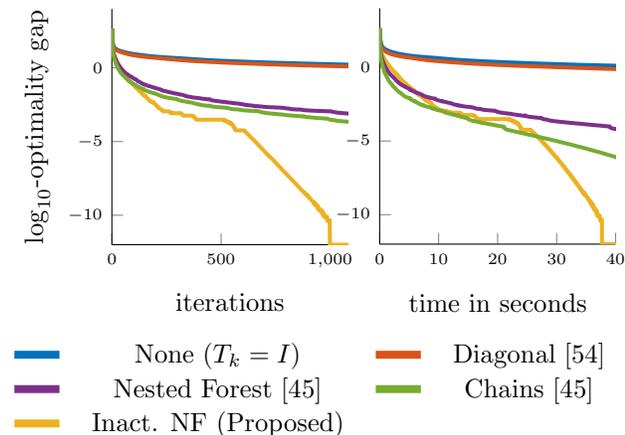


Figure 4: \log_{10} -optimality gap over iterations (left) and time (right) for PDHG with various preconditioners applied to a TV deconvolution problem.

two graphs are huge ($|\mathcal{V}| \approx 250,000$, $|\mathcal{E}| \approx 600,000$) and therefore partitioning is quite expensive. To summarize, the proposed preconditioning strategy consistently improves the number of iterations, but to ensure a shorter overall running time, an efficient implementation or improved strategy on reconstructing the tree decomposition might be required.

5.3 Linear inverse problems

In this image processing experiment we consider a TV deconvolution problem on a regular 2D grid of size 116×87 . The data term is given by $G(u) = \frac{1}{2} \|Au - f\|^2$, where the forward model A is a convolution with motion blur kernel with radius 3. We construct f by applying the forward model and adding Gaussian noise. The overall problem is solved using PDHG. The primal update is a quadratic problem and we use a few itera-

tions of (warm started) conjugate gradient. Considering the size of the problem, we set the reconditioning frequency to $n = 5$ for the proposed approach.

In Fig. 4 we show the \log_{10} -optimality gap over iterations and time for various choices of preconditioners. The diagonal preconditioner is the one from [54] with $\alpha = 1$. The forest preconditioners perform comparably when the accuracy is lower. Once the local convergence regime is entered, the proposed algorithm achieves linear convergence rate. Especially for high accuracies, the proposed inactively nested forest reconditioning strategy outperforms the other approaches with respect to overall running time and iterations.

6 Discussion and conclusion

We presented an efficient reconditioning strategy for proximal algorithms on graphs. By relying on a sharp analysis of the local linear convergence rate we proposed an edge partitioning of the graph into forests which provably boosts the linear convergence rate. The scaled dual updates are still efficiently computable thanks to a message-passing algorithm on trees.

While one is tempted to commit to a super-linearly convergent solver once the optimal active set is identified (as e.g., mentioned in [41, 42, 43, 51]), it is unfortunately difficult to verify in practice whether the current active set is the optimal. Furthermore, as observed in the numerical experiments, the adaptive preconditioning strategy practically improves the convergence also *before* the local linear convergence regime is entered. The result suggests that local convergence analysis can serve as a practical guideline for constructing preconditioners for proximal algorithms.

References

- [1] H. Attouch, J. Peypouquet, and P. Redont. Fast convergence of an inertial gradient-like system with vanishing viscosity. *arXiv:1507.04782*, 2015.
- [2] A. Barbero and S. Sra. Fast Newton-type methods for total variation regularization. In *Proceedings of the 32nd International Conference on Machine Learning, ICML*, 2011.
- [3] A. Barbero and S. Sra. Modular proximal optimization for multidimensional total-variation regularization. *arXiv:1411.0589*, 2014.
- [4] H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer Science & Business Media, 2011.
- [5] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.*, 2:183–202, 2009.
- [6] S. Becker and M. J. Fadili. A quasi-Newton proximal splitting method. In *Proceedings of the 26th International Conference on Neural Information Processing Systems, NIPS*, 2012.
- [7] S. Becker, J. Fadili, and P. Ochs. On quasi-Newton forward-backward splitting: Proximal calculus and convergence. *arXiv:1801.08691*, 2018.
- [8] K. Benzi, V. Kalofolias, X. Bresson, and P. Vandergheynst. Song recommendation with non-negative matrix factorization and graph total variation. In *Proceedings of the 42nd IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP*, 2016.
- [9] K. Bredies and H. Sun. Preconditioned Douglas-Rachford splitting methods for convex-concave saddle-point problems. *SIAM J. Numer. Anal.*, 53:421–444, 2015.
- [10] X. Bresson, T. Laurent, D. Uminsky, and J. Von Brecht. Multiclass total variation clustering. In *Proceedings of the 27th International Conference on Neural Information Processing Systems, NIPS*, 2013.
- [11] J. V. Burke and J. J. Moré. On the identification of active constraints. *SIAM J. Numer. Anal.*, 25(5):1197–1211, 1988.
- [12] A. Chambolle and J. Darbon. On total variation minimization and surface evolution using parametric maximum flows. *Int. J. Comput. Vis.*, 84:288–307, 2009.
- [13] A. Chambolle and C. Dossal. On the convergence of the iterates of the “fast iterative shrinkage/thresholding algorithm”. *J. Optim. Theory Appl.*, 166:968–982, 2015.
- [14] A. Chambolle and T. Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *J. Math. Imaging Vis.*, 40:120–145, 2011.
- [15] A. Chambolle and T. Pock. On the ergodic convergence rates of a first-order primal-dual algorithm. *Math. Program.*, 159:253–287, 2016.
- [16] A. Chambolle and T. Pock. An introduction to continuous optimization for imaging. *Acta Numer.*, 25:161–319, 2016.
- [17] D. Cheriton and R. E. Tarjan. Finding minimum spanning trees. *SIAM J. Comput.*, 5(4):724–742, 1976.
- [18] P. L. Combettes and J.-C. Pesquet. Proximal splitting methods in signal processing. In *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, pages 185–212. Springer, 2011.
- [19] L. Condat. A direct algorithm for 1D total variation denoising. *IEEE Signal Process. Lett.*, 20:1054–1057, 2013.
- [20] S. Diamond and S. Boyd. Stochastic matrix-free equilibration. *J. Optim. Theory Appl.*, 172:436–454, 2017.
- [21] C. Fougner and S. Boyd. Parameter selection and pre-conditioning for a graph form solver. *arXiv:1503.08366*, 2015.
- [22] M. P. Friedlander and G. Goh. Efficient evaluation of scaled proximal operators. *Electron. Trans. Numer. Anal.*, 46:1–22, 2017.
- [23] C. Garcia-Cardona, E. Merkurjev, A. L. Bertozzi, A. Flenner, and A. G. Percus. Multiclass data segmentation using diffuse interface methods on graphs. *IEEE Trans. Pattern Anal. Mach. Intell.*, 36(8):1600–1613, 2014.
- [24] G. Garrigos, L. Rosasco, and S. Villa. Convergence of the forward-backward algorithm: Beyond the worst case with the help of geometry. *arXiv:1703.09477*, 2017.
- [25] G. Gilboa and S. Osher. Nonlocal operators with applications to image processing. *Multiscale Modeling & Simulation*, 7(3):1005–1028, 2008.
- [26] P. Giselsson and S. Boyd. Diagonal scaling in Douglas-Rachford splitting and ADMM. In *Proceedings of the 53rd IEEE Conference on Decision and Control, CDC*, 2014.
- [27] P. Giselsson and S. Boyd. Preconditioning in fast dual gradient methods. In *Proceedings of the 53rd IEEE Conference on Decision and Control, CDC*, 2014.

- [28] P. Giselsson and S. Boyd. Metric selection in fast dual forward–backward splitting. *Automatica*, 62: 1–10, 2015.
- [29] A. Goldberg, S. Hed, H. Kaplan, R. Tarjan, and R. Werneck. Maximum flows by incremental breadth-first search. *European Symposium on Algorithms, ALGO ESA*, 2011.
- [30] M. Hein and S. Setzer. Beyond spectral clustering – tight relaxations of balanced graph cuts. In *Proceedings of the 25th International Conference on Neural Information Processing Systems, NIPS*, 2011.
- [31] M. Hein, S. Setzer, L. Jost, and S. S. Rangapuram. The total variation on hypergraphs – learning on hypergraphs revisited. In *Proceedings of the 27th International Conference on Neural Information Processing Systems, NIPS*, 2013.
- [32] D. S. Hochbaum. An efficient algorithm for image segmentation, Markov random fields and related problems. *Journal of the ACM (JACM)*, 48(4): 686–701, 2001.
- [33] A. J. Hoffman. On approximate solutions of systems of linear inequalities. *J. Res. Natl. Bur. Standards*, 49:263–265, 1952.
- [34] H. Karimi, J. Nutini, and M. Schmidt. Linear convergence of gradient and proximal-gradient methods under the Polyak-Lojasiewicz condition. In *Proceedings of the Joint European Conference on Machine Learning and Knowledge Discovery in Databases ECML PKDD*, 2016.
- [35] D. Klatté and G. Thiere. Error bounds for solutions of linear equations and inequalities. *Zeitschrift für Operations Research*, 41(2):191–214, 1995.
- [36] V. Kolmogorov, T. Pock, and M. Rolinek. Total variation on a tree. *SIAM J. Imaging Sci.*, 9:605–636, 2016.
- [37] J. B. Kruskal. On the shortest spanning subtree of a graph and the traveling salesman problem. *Proceedings of the American Mathematical Society*, 1956.
- [38] K. Kumar and F. Bach. Active-set methods for submodular minimization problems. *J. Mach. Learn. Res.*, 18(132):1–31, 2017.
- [39] L. Landrieu and G. Obozinski. Cut pursuit: Fast algorithms to learn piecewise constant functions on general weighted graphs. *SIAM J. on Imaging Sci.*, 10(4):1724–1766, 2017.
- [40] J. D. Lee, Y. Sun, and M. A. Saunders. Proximal Newton-type methods for minimizing composite functions. *SIAM J. Optim.*, 24:1420–1443, 2014.
- [41] J. Liang, J. Fadili, and G. Peyré. Local linear convergence of forward–backward under partial smoothness. In *Proceedings of the 28th International Conference on Neural Information Processing Systems, NIPS*, 2014.
- [42] J. Liang, J. Fadili, and G. Peyré. Activity identification and local linear convergence of forward–backward-type methods. *SIAM J. Optim.*, 27: 408–437, 2017.
- [43] J. Liang, J. Fadili, and G. Peyré. Local linear convergence analysis of primal–dual splitting methods. *Optimization*, 67(6):821–853, 2018.
- [44] Y. Lou, X. Zhang, S. Osher, and A. Bertozzi. Image recovery via nonlocal operators. *Journal of Scientific Computing*, 42(2):185–197, 2010.
- [45] T. Möllenhoff, Z. Ye, T. Wu, and D. Cremers. Combinatorial preconditioners for proximal algorithms on graphs. In *Proceedings of the 21st International Conference on Artificial Intelligence and Statistics, AISTATS*, 2018.
- [46] I. Necoara, Y. Nesterov, and F. Glineur. Linear convergence of first order methods for non-strongly convex optimization. *arXiv:1504.06298*, 2015.
- [47] Y. Nesterov. A method for solving the convex programming problem with convergence rate $\mathcal{O}(1/k^2)$. *Soviet Mathematics Doklady*, 269:543–547, 1983.
- [48] R. A. Newcombe, S. J. Lovegrove, and A. J. Davison. DTAM: Dense tracking and mapping in real-time. In *Proceedings of the 13th International Conference on Computer Vision, ICCV*, 2011.
- [49] R. Nishihara, L. Lessard, B. Recht, A. Packard, and M. Jordan. A general analysis of the convergence of ADMM. In *Proceedings of the 32nd International Conference on Machine Learning, ICML*, 2015.
- [50] J. Nutini, I. Laradji, and M. Schmidt. Let’s make block coordinate descent go fast: Faster greedy rules, message-passing, active-set complexity, and superlinear convergence. *arXiv:1712.08859*, 2017.
- [51] J. Nutini, M. Schmidt, and W. Hare. ”Active-set complexity” of proximal gradient: How long does it take to find the sparsity pattern? *arXiv:1712.03577*, 2017.
- [52] P. Ochs, Y. Chen, T. Brox, and T. Pock. iPiano: Inertial proximal algorithm for nonconvex optimization. *SIAM J. Imaging Sci.*, 7:1388–1419, 2014.
- [53] N. Parikh and S. Boyd. Proximal algorithms. *Foundations and Trends in Optimization*, 1:123–231, 2013.

- [54] T. Pock and A. Chambolle. Diagonal preconditioning for first order primal-dual algorithms in convex optimization. In *Proceedings of the 13th International Conference on Computer Vision, ICCV*, 2011.
- [55] B. T. Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4:1–17, 1964.
- [56] H. Raguét and L. Landrieu. Cut-pursuit algorithm for regularizing nonsmooth functionals with graph total variation. In *Proceedings of the 35th International Conference on Machine Learning, ICML*, 2018.
- [57] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena*, 60(1):259–268, 1992.
- [58] M. Schmidt, D. Kim, and S. Sra. Projected Newton-type methods in machine learning. *Optimization for Machine Learning*, 2012.
- [59] D. A. Spielman. Algorithms, graph theory, and linear equations in Laplacian matrices. In *Proceedings of the International Congress of Mathematicians, ICM*, 2010.
- [60] R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight. Sparsity and smoothness via the fused lasso. *J. R. Stat. Soc.*, 67(1):91–108, 2005.
- [61] P. M. Vaidya. Solving linear equations with symmetric diagonally dominant matrices by constructing good preconditioners. (A talk based on the manuscript was presented at the IMA Workshop on Graph Theory and Sparse Matrix Computation), 1991.
- [62] B. Xin, Y. Kawahara, Y. Wang, and W. Gao. Efficient generalized fused Lasso and its application to the diagnosis of Alzheimer’s disease. In *Proceedings of the 28th AAAI Conference on Artificial Intelligence*, 2014.

Optimization of Graph Total Variation via Active-Set-based Combinatorial Reconditioning

— Supplementary Material —

Zhenzhang Ye
TU Munich
zhenzhang.ye@tum.de

Thomas Möllenhoff
TU Munich
thomas.moellenhoff@tum.de

Tao Wu
TU Munich
tao.wu@tum.de

Daniel Cremers
TU Munich
cremers@tum.de

Lemma 6. *Let h be C^2 with $l_h I \preceq \nabla^2 h(\cdot) \preceq L_h I$ for some constants $l_h, L_h > 0$. Then the gradient descent on $\min_x h(Ax + b)$ with step size $1/t = 2/(L_h \sigma_{\max}(A)^2 + l_h \sigma_{\min>0}(A)^2)$ satisfies*

$$\|x^{k+1} - x^*\| \leq \frac{\varphi - 1}{\varphi + 1} \|x^k - x^*\|, \quad (44)$$

with $\varphi = \kappa(A)^2 \cdot \kappa(h)$, $\kappa(h) := L_h/l_h$.

Proof of Lemma 6. Clearly $t > (L_h \cdot \lambda_{\max}(A^\top A))/2$, so classical theory (e.g. [4]) guarantees $x^k \rightarrow x^*$. First note that

$$x^{k+1} - x^k = -\frac{1}{t} A^\top \nabla h(Ax^k + b) \in \text{ran } A^\top, \quad (45)$$

and consequently $x^k - x^0 \in (\ker A)^\perp$, also $x^* - x^k \in (\ker A)^\perp$. Inserting the gradient step for x^{k+1} yields

$$\|x^{k+1} - x^*\| = \|x^k - x^* - \frac{1}{t} A^\top (\nabla h(Ax^k + b) - \nabla h(Ax^* + b))\|. \quad (46)$$

From the mean value theorem it follows

$$\nabla h(Ax^k + b) - \nabla h(Ax^* + b) = M(Ax^k - Ax^*), \quad (47)$$

with $M = \int_0^1 \nabla^2 h(Ax^* + b + \alpha(Ax^k - Ax^*)) d\alpha$. Since $l_h I \preceq \nabla^2 h(\cdot) \preceq L_h I$ we have $l_h I \preceq M \preceq L_h I$. This yields due to $x^k - x^* \in (\ker A)^\perp$ that

$$\|x^{k+1} - x^*\| = \|(I - \frac{1}{t} A^\top M A)(x^k - x^*)\| \leq \max\{|1 - l_h \sigma_{\min>0}(A)^2/t|, |1 - L_h \sigma_{\max}(A)^2/t|\} \cdot \|x^k - x^*\|. \quad (48)$$

The choice $t = (l_h \sigma_{\min>0}(A)^2 + L_h \sigma_{\max}(A)^2)/2$ minimizes the above rate and yields the desired result. \square