

Multi-stage Convex Relaxation for Non-convex Optimization

Tong Zhang
Statistics Department
Rutgers University, NJ
tzhang@stat.rutgers.edu

Abstract

We consider learning formulations with non-convex objective functions that often occur in practical applications. There are two approaches to this problem:

- Heuristic methods such as gradient descent that only find a local minimum. A drawback of this approach is the lack of theoretical guarantee showing that the local minimum gives a good solution.
- Convex relaxation such as L_1 -regularization that solves the problem under some conditions. However it often leads to a sub-optimal solution in reality.

This paper tries to remedy the above gap between theory and practice. In particular, we investigate a multi-stage convex relaxation scheme for solving problems with non-convex objective functions. The general approach is outlined first. Then for learning formulations with sparse regularization, we analyze the behavior of a resulting multi-stage relaxation scheme. Under appropriate conditions, we show that the local solution obtained by this procedure is superior to the global solution of the standard L_1 convex relaxation for learning sparse targets.

1 Introduction

We consider the general regularization framework for machine learning, where a loss function is minimized, subject to a regularization condition on the model parameter. For many natural machine learning problems, either the loss function or the regularization condition can be non-convex. For example, the loss function is non-convex for classification problems, and the regularization condition is non-convex in problems with sparse parameters.

A major difficulty with nonconvex formulations is that the global optimal solution cannot be efficiently computed, and the behavior of a local solution is hard to analyze. In practice, convex relaxation (such as support vector machine for classification or L_1 regularization for sparse learning) has been adopted to remedy the problem. The choice of convex formulation makes the solution unique and efficient to compute. Moreover, the solution is easy to analyze theoretically. That is, it can be shown that the solution of the convex formulation approximately solves the original problem under appropriate assumptions. However, for many practical problems, such simple convex relaxation schemes can be sub-optimal.

Because of the above gap between practice and theory, it is important to study direct solutions of non-convex optimization problems beyond the standard convex relaxation. Our goal is to design a numerical procedure that leads to a *reproducible solution* which is better than the standard convex

relaxation solution. In order to achieve this, we present a general framework of multi-stage convex relaxation, which iteratively refine the convex relaxation formulation to give better solutions. The method is derived from concave duality, and involves solving a sequence of convex relaxation problems, leading to better and better approximations to the original nonconvex formulation. Since each stage is a convex optimization problem, the approach is computationally efficient. Although the method only leads to a local optimal solution for the original nonconvex problem, this local solution is a refinement of the global solution for the initial convex relaxation. Therefore intuitively one expects that the local solution is better than the standard one-stage convex relaxation. In order to prove this observation more rigorously, we consider least squares regression with nonconvex sparse regularization terms, for which we can analyze the effectiveness of the multi-stage convex relaxation. It is shown that under appropriate assumptions, the (local) solution computed by the multi-stage convex relaxation method using nonconvex regularization achieves better parameter estimation performance than the standard convex relaxation with L_1 regularization.

The paper is organized into two parts. The first part, from Section 2 to Section 4, presents the general framework for multi-stage convex relaxation, with potential application examples. The second part, presented in Section 5, considers the special case of sparse learning, where we derive theoretical results showing that under appropriate conditions, it is beneficial to use multi-stage convex relaxation with nonconvex regularization as opposed to the standard convex L_1 regularization.

2 Non-convex Formulations in Machine Learning

The combination of regularization and risk minimization is essential in modern machine learning. We shall first motivate this class of learning algorithms from supervised learning as follows. Consider a set of input vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in R^d$, with corresponding desired output variables y_1, \dots, y_n . The task of supervised learning is to estimate the functional relationship $y \approx f(\mathbf{x})$ between the input \mathbf{x} and the output variable y from the training examples $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$. The quality of prediction is often measured through a loss function $\phi(f(\mathbf{x}), y)$. We assume that $\phi(f, y)$ is convex in f throughout the paper. In this paper, we consider linear prediction model $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$. As in boosting or kernel methods, nonlinearity can be introduced by including nonlinear features in \mathbf{x} .

For linear models, we are mainly interested in the scenario that $d \gg n$. That is, there are many more features than the number of samples. In this case, an unconstrained empirical risk minimization is inadequate because the solution overfits the data. The standard remedy for this problem is to impose a constraint on \mathbf{w} to obtain a *regularized* problem. This leads to the following regularized empirical risk minimization method:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in R^d} \left[\sum_{i=1}^n \phi(\mathbf{w}^\top \mathbf{x}_i, y_i) + \lambda g(\mathbf{w}) \right], \quad (1)$$

where $\lambda > 0$ is an appropriately chosen regularization condition.

Loss Function

Examples of loss function $\phi(\mathbf{w}^\top \mathbf{x}, y)$ in (1) include least squares for regression: $\phi(\mathbf{w}^\top \mathbf{x}, y) = (\mathbf{w}^\top \mathbf{x} - y)^2$, and 0-1 binary classification error: $\phi(\mathbf{w}^\top \mathbf{x}, y) = I(\mathbf{w}^\top \mathbf{x} y \leq 0)$, where $y \in \{\pm 1\}$ are the class labels, and $I(\cdot)$ is the set indicator function. The latter is nonconvex. In practice, for computationally reasons, a convex relaxation such as the SVM loss $\phi(\mathbf{w}^\top \mathbf{x}, y) = \max(0, 1 - \mathbf{w}^\top \mathbf{x} y)$,

is often used to substitute the classification error loss. Such a convex loss is often referred to as a surrogate loss function, and the resulting method becomes a convex relaxation method for solving binary classification. This class of methods have been theoretically analyzed in [2, 13]. While asymptotically, convex surrogate methods are consistent (that is, they can be used to obtain Bayes optimal classifiers when the sample size approaches infinity), for finite data, these methods can be more sensitive to outliers. In order to alleviate the effect of outliers, one may consider the smoothed classification error loss function $\phi(\mathbf{w}^\top \mathbf{x}, y) = \min(\alpha, \max(0, 1 - \mathbf{w}^\top \mathbf{x}y))$ ($\alpha \geq 1$). This loss function is bounded, and thus more robust to outliers than SVMs under finite sample size; moreover, it is piece-wise differentiable, and thus computationally feasible. For comparison purpose, the three loss functions are plotted in Figure 1.

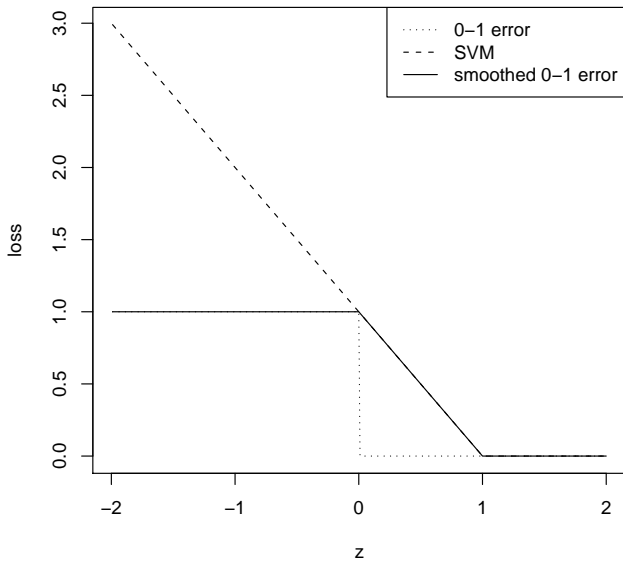


Figure 1: Loss Functions: classification error versus smoothed classification error ($\alpha = 1$) and SVM

Regularization Condition

Some examples of regularization conditions in (1) include squared regularization $g(\mathbf{w}) = \mathbf{w}^\top \mathbf{w}$, and 1-norm regularization $g(\mathbf{w}) = \|\mathbf{w}\|_1$. The former can be generalized to kernel methods, and the latter leads to sparse solutions. Sparsity is an important regularization condition, which corresponds to the (non-convex) L_0 regularization, defined as $\|\mathbf{w}\|_0 = |\{j : \mathbf{w}_j \neq 0\}| = k$. That is, the measure of complexity is the number of none zero coefficients. If we know the sparsity parameter k for the target vector, then a good learning method is L_0 regularization:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{w}^\top \mathbf{x}_i, y_i) \quad \text{subject to } \|\mathbf{w}\|_0 \leq k, \tag{2}$$

which applies the standard empirical risk minimization formulation to learning L_0 constrained sparse targets.

If k is not known, then one may regard k as a tuning parameter, which can be selected through cross-validation. This method is often referred to as *subset selection* in the literature. Sparse learning is an essential topic in machine learning, which has attracted considerable interests recently. It can be shown that the solution of the L_0 regularization problem in (2) achieves good prediction accuracy if the target function can be approximated by a sparse $\bar{\mathbf{w}}$. However, a fundamental difficulty with this method is the computational cost, because the number of subsets of $\{1, \dots, d\}$ of cardinality k (corresponding to the nonzero components of \mathbf{w}) is exponential in k .

Due to the computational difficult, in practice, it is necessary to replace (2) by some easier to solve formulations in (1). Specifically, L_0 regularization is equivalent to (1) by choosing the regularization function as $g(\mathbf{w}) = \|\mathbf{w}\|_0$. However, this function is discontinuous. For computational reasons, it is helpful to consider a continuous approximation with $g(\mathbf{w}) = \|\mathbf{w}\|_p^p$, where $p > 0$. If $p \geq 1$, the resulting formulation is convex. In particular, by choosing the closest approximation with $p = 1$, one obtain *Lasso*, which is the standard convex relaxation formulation for sparse learning. With $p \in (0, 1)$, the L_p regularizer $\|\mathbf{w}\|_p^p$ is non-convex but continuous.

General Regularized Learning Formulation

Supervised learning can be solved using general empirical risk minimization formulation in (1). Both ϕ and g can be non-convex in application problems. The traditional approach is to use convex relaxation to approximate it, leading to a single stage convex formulation. In this paper, we try to extend the idea by looking at a more general multi-stage convex relaxation method, which leads to more accurate approximations.

We consider an optimization formulation more general than (1) as follows:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} R(\mathbf{w}),$$

$$R(\mathbf{w}) = R_0(\mathbf{w}) + \sum_{k=1}^K R_k(\mathbf{w}), \quad (3)$$

where $R(\mathbf{w})$ is the general form of a regularized objective function. Moreover, for convenience, we assume that $R_0(\mathbf{w})$ is convex in \mathbf{w} , and each $R_k(\mathbf{w})$ is non-convex. In the proposed work, we shall employ convex/concave duality to derive convex relaxations of (3) that can be efficiently solved.

Related to (3), one may also consider the constrained formulation

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} R_0(\mathbf{w}) \quad \text{subject to} \quad \sum_{k=1}^K R_k(\mathbf{w}) \leq A, \quad (4)$$

where A is a constant. One may also mix (3) and (4).

For illustration, we consider the following examples which will be used in our later discussion.

- Smoothed classification error loss: formulation (1) with convex regularization $g(\mathbf{w})$ and non-convex loss function (with $\alpha \geq 1$)

$$\phi(\mathbf{w}^\top \mathbf{x}, y) = \min(\alpha, \max(0, 1 - \mathbf{w}^\top \mathbf{x}y)).$$

This corresponds to $R_0(\mathbf{w}) = \lambda g(\mathbf{w})$, and $R_k(\mathbf{w}) = \phi(\hat{\mathbf{w}}^\top \mathbf{x}_k, y_k)$ for $k = 1, \dots, n$ in (3).

- L_p regularization ($0 \leq p \leq 1$): formulation (1) with nonconvex regularization $g(\mathbf{w}) = \|\mathbf{w}\|_p^p$ and a loss function $\phi(\cdot, \cdot)$ that is convex in \mathbf{w} . This corresponds to $R_0(\mathbf{w}) = n^{-1} \sum_{i=1}^n \phi(\mathbf{w}^\top \mathbf{x}_i, y_i)$, and $R_k(\mathbf{w}) = \lambda |\mathbf{w}_k|^p$ for $k = 1, \dots, d$ in (3).
- Smoothed L_p regularization (with parameters $\alpha > 0$ and $0 \leq p \leq 1$): formulation (1) with nonconvex regularization $g(\mathbf{w}) = \sum_k [(\alpha + |\mathbf{w}_k|)^p - \alpha^p] / (p\alpha^{p-1})$, and a loss function $\phi(\cdot, \cdot)$ that is convex in \mathbf{w} . This corresponds to $R_0(\mathbf{w}) = n^{-1} \sum_{i=1}^n \phi(\mathbf{w}^\top \mathbf{x}_i, y_i)$, and $R_k(\mathbf{w}) = \lambda [(\alpha + |\mathbf{w}_k|)^p - \alpha^p] / (p\alpha^{p-1})$ for $k = 1, \dots, d$ in (3). The main difference between standard L_p and smoothed L_p is at $|\mathbf{w}_k| = 0$, where the smoothed L_p regularization is differentiable, with derivative 1. This difference is theoretically important as discussed in Section 5.1.
- Smoothed log regularization (with parameter $\alpha > 0$): formulation (1) with nonconvex regularization $g(\mathbf{w}) = \sum_k \alpha \ln(\alpha + |\mathbf{w}_k|)$, and a loss function $\phi(\cdot, \cdot)$ that is convex in \mathbf{w} . This corresponds to $R_0(\mathbf{w}) = n^{-1} \sum_{i=1}^n \phi(\mathbf{w}^\top \mathbf{x}_i, y_i)$, and $R_k(\mathbf{w}) = \lambda \alpha \ln(\alpha + |\mathbf{w}_k|)$ for $k = 1, \dots, d$ in (3). Similar to the smoothed L_p regularization, the smoothed log-loss has derivative 1 at $|\mathbf{w}_k| = 0$.
- Capped- L_1 regularization (with parameter $\alpha > 0$): formulation (1) with nonconvex regularization $g(\mathbf{w}) = \sum_{j=1}^d \min(|\mathbf{w}_j|, \alpha)$, and a loss function $\phi(\cdot, \cdot)$ that is convex in \mathbf{w} . This corresponds to $R_0(\mathbf{w}) = n^{-1} \sum_{i=1}^n \phi(\mathbf{w}^\top \mathbf{x}_i, y_i)$, and $R_k(\mathbf{w}) = \lambda \min(|\mathbf{w}_k|, \alpha)$ for $k = 1, \dots, d$ in (3). The capped- L_1 regularization is a good approximation to L_0 because as $\alpha \rightarrow 0$, $\sum_k \min(|\mathbf{w}_k|, \alpha) / \alpha \rightarrow \|\mathbf{w}\|_0$. Therefore when $\alpha \rightarrow 0$, this regularization condition is equivalent to the sparse L_0 regularization upto a rescaling of λ .

3 Concave Duality

In the following discussion, we consider a single nonconvex component $R_k(\mathbf{w})$ in (3), which we shall rewrite using concave duality. Let $\mathbf{h}_k(\mathbf{w}) : R^d \rightarrow \Omega_k \subset R^{d_k}$ be a vector function with range Ω_k . It may not be a one-to-one map. However, we assume that there exists a function \bar{R}_k defined on Ω_k so that we can express $R_k(\mathbf{w})$ as

$$R_k(\mathbf{w}) = \bar{R}_k(\mathbf{h}_k(\mathbf{w})).$$

Assume that we can find \mathbf{h}_k so that the function $\bar{R}_k(\mathbf{u}_k)$ is concave on $\mathbf{u}_k \in \Omega_k$. Under this assumption, we can rewrite the regularization function $R_k(\mathbf{w})$ as:

$$R_k(\mathbf{w}) = \inf_{\mathbf{v}_k \in R^{d_k}} \left[\mathbf{v}_k^\top \mathbf{h}_k(\mathbf{w}) + R_k^*(\mathbf{v}_k) \right] \quad (5)$$

using concave duality [10]. In this case, $R_k^*(\mathbf{v}_k)$ is the concave dual of $\bar{R}_k(\mathbf{u}_k)$ given below

$$R_k^*(\mathbf{v}_k) = \inf_{\mathbf{u}_k \in \Omega_k} \left[-\mathbf{v}_k^\top \mathbf{u}_k + \bar{R}_k(\mathbf{u}_k) \right].$$

Moreover, it is well-known that the minimum of the right hand side of (5) is achieved at

$$\hat{\mathbf{v}}_k = \nabla_{\mathbf{u}} \bar{R}_k(\mathbf{u})|_{\mathbf{u}=\mathbf{h}_k(\mathbf{w})}. \quad (6)$$

This is a very general framework. For illustration, we include some of the example non-convex conditions discussed in the introduction.

Smoothed classification error

We consider a loss term of the form $R_k(\mathbf{w}) = \min(\alpha, \max(0, 1 - \mathbf{w}^\top \mathbf{x}_k y_k))$ for $k = 1, \dots, n$ (with $\alpha \geq 1$), and relax it to the SVM loss $\mathbf{h}_k(\mathbf{w}) = \max(0, 1 - \mathbf{w}^\top \mathbf{x}_k y_k)$. In this case, each \mathbf{u}_k is a scalar in the range $\Omega_k = [0, \infty)$, and $\bar{R}_k(\mathbf{u}_k) = \min(\alpha, \mathbf{u}_k)$. We have $R_k^*(\mathbf{v}_k) = \alpha(1 - \mathbf{v}_k)I(\mathbf{v}_k \in [0, 1])$. The solution in (6) is given by $\hat{\mathbf{v}}_k = I(\mathbf{w}^\top \mathbf{x}_k y_k \geq 1 - \alpha)$ for $k = 1, \dots, n$.

L_p regularization

We consider a regularization term $R_k(\mathbf{w}) = \lambda|\mathbf{w}_k|^p$ ($k = 1, \dots, d$) for some $p \in (0, 1)$. Given any $q > p$, (5) holds with $\mathbf{u}_k = \mathbf{h}_k(\mathbf{w}) = |\mathbf{w}_k|^q \in [0, \infty)$, and $\bar{R}_k(\mathbf{u}_k) = \lambda|\mathbf{u}_k|^{p/q}$, where $\mathbf{u}_k \in \Omega_k = [0, \infty)$. The dual is $R_k^*(\mathbf{v}_k) = \lambda c(p, q)(\mathbf{v}_k/\lambda)^{p/(p-q)}$, defined on the domain $\mathbf{v}_k \geq 0$, where $c(p, q) = (q-p)p^{p/(q-p)}q^{q/(p-q)}$. The solution in (6) is given by $\hat{\mathbf{v}}_k = \lambda(p/q)|\mathbf{w}_k|^{p-q}$.

Smoothed L_p regularization

We consider a regularization term $R_k(\mathbf{w}) = \lambda[(\alpha + |\mathbf{w}_k|)^p - \alpha^p]/(p\alpha^{p-1})$ ($k = 1, \dots, d$) for some $p \in (0, 1)$ and $\alpha > 0$. Given any $q > p$, (5) holds with $\mathbf{u}_k = \mathbf{h}_k(\mathbf{w}) = (\alpha + |\mathbf{w}_k|)^q \in [\alpha^q, \infty)$, and $\bar{R}_k(\mathbf{u}_k) = \lambda[\mathbf{u}_k^{p/q} - \alpha^p]/(p\alpha^{p-1})$, where $\mathbf{u}_k \in \Omega_k = [0, \infty)$. The dual is $R_k^*(\mathbf{v}_k) = \lambda c(p, q)(\mathbf{v}_k/\lambda)^{p/(p-q)} - \lambda\alpha^p/(p(\alpha^{p-1}))$, defined on the domain $\mathbf{v}_k \geq 0$, where $c(p, q) = (q-p)p^{p/(q-p)}q^{q/(p-q)}$. The solution in (6) is given by $\hat{\mathbf{v}}_k = \lambda/(q\alpha^{p-1})(\alpha + |\mathbf{w}_k|)^{p-q}$.

An alternative is to relax smoothed L_p regularization ($p \in (0, 1)$) directly to L_q regularization for $q \geq 1$ (one usually takes either $q = 1$ or $q = 2$). In this case, $\mathbf{u}_k = \mathbf{h}_k(\mathbf{w}) = |\mathbf{w}_k|^q \in [0, \infty)$, and $\bar{R}_k(\mathbf{u}_k) = \lambda[(\alpha + \mathbf{u}_k^{1/q})^p - \alpha^p]/(p\alpha^{p-1})$. Although it is not difficult to verify that $\bar{R}_k(\mathbf{u}_k)$ is concave, we do not have a simple closed form for $R_k^*(\mathbf{v}_k)$. However, it is easy to check that the solution in (6) is given by $\hat{\mathbf{v}}_k = \lambda/(q\alpha^{p-1})(\alpha + |\mathbf{w}_k|)^{p-1}|\mathbf{w}_k|^{1-q}$.

Smoothed log regularization

We consider a regularization term $R_k(\mathbf{w}) = \lambda\alpha \ln(\alpha + |\mathbf{w}_k|)$ for some $\alpha > 0$. Given any $q > 0$, (5) holds with $\mathbf{u}_k = \mathbf{h}_k(\mathbf{w}) = (\alpha + |\mathbf{w}_k|)^q \in [\alpha^q, \infty)$, and $\bar{R}_k(\mathbf{u}_k) = \lambda(\alpha/q) \ln(\mathbf{u}_k)$, where $\mathbf{u}_k \in \Omega_k = [0, \infty)$. The dual is $R_k^*(\mathbf{v}_k) = \lambda(\alpha/q)[- \ln \mathbf{v}_k - 1]$, defined on the domain $\mathbf{v}_k \geq 0$. The solution in (6) is given by $\hat{\mathbf{v}}_k = \lambda(\alpha/q)(\alpha + |\mathbf{w}_k|)^{-q}$.

Similar to smoothed L_p , we may relax directly to L_q , with $\mathbf{u}_k = \mathbf{h}_k(\mathbf{w}) = |\mathbf{w}_k|^q \in [0, \infty)$. $\bar{R}_k(\mathbf{u}_k) = \lambda\alpha \ln(\alpha + \mathbf{u}_k^{1/q})$, where The solution in (6) is given by $\hat{\mathbf{v}}_k = \lambda(\alpha/q)(\alpha + |\mathbf{w}_k|)^{-1}|\mathbf{w}_k|^{1-q}$.

Capped- L_1 regularization

We consider a regularization term $R_k(\mathbf{w}) = \lambda \min(|\mathbf{w}_k|, \alpha)$ ($k = 1, \dots, d$) for some $\alpha > 0$. In this case, (5) holds with $\mathbf{u}_k = \mathbf{h}_k(\mathbf{w}) = |\mathbf{w}_k| \in [0, \infty)$, and $\bar{R}_k(\mathbf{u}_k) = \lambda \min(\mathbf{u}_k, \alpha)$, where $\mathbf{u}_k \in \Omega_k = [0, \infty)$. The dual is $R_k^*(\mathbf{v}_k) = \lambda\alpha(1 - \mathbf{v}_k/\lambda)I(\mathbf{v}_k \in [0, \lambda])$ defined $[0, \infty)$, where $I(\cdot)$ is the set indicator function. The solution in (6) is given by $\hat{\mathbf{v}}_k = \lambda I(|\mathbf{w}_k| \leq \alpha)$.

4 Multi-stage Convex Relaxation

Using concave duality given in the previous section, we can derive a general convex relaxation based procedure for solving the penalized formulation (3) or the constrained formulation (4).

4.1 Penalized formulation

Let $h_k(\mathbf{w})$ be a convex relaxation of $R_k(\mathbf{w})$ that dominates $R_k(\mathbf{w})$ (for example, it can be the smallest convex upperbound (i.e., the inf over all convex upperbounds) of $R_k(\mathbf{w})$). A simple convex relaxation of (3) becomes

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in R^d} \left[R_0(\mathbf{w}) + \sum_{k=1}^K \mathbf{h}_k(\mathbf{w})^\top \mathbf{v}_k \right]. \quad (7)$$

This simple relaxation yields a solution that is different from the solution of (3). However, if each \mathbf{h}_k satisfies the condition of Section 3, then it is possible to write $R_k(\mathbf{w})$ using (5). Now, with this new representation, we can rewrite (3) as

$$[\hat{\mathbf{w}}, \hat{\mathbf{v}}] = \arg \min_{\mathbf{w}, \{\mathbf{v}_k\}} \left[R_0(\mathbf{w}) + \sum_{k=1}^K (\mathbf{h}_k(\mathbf{w})^\top \mathbf{v}_k + R_k^*(\mathbf{v}_k)) \right]. \quad (8)$$

This is clearly equivalent to (1) because of (5). If we can find a good approximation of $\hat{\mathbf{v}} = \{\hat{\mathbf{v}}_k\}$ that improves upon the initial value of $\hat{\mathbf{v}}_k = \mathbf{1}$, then the above formulation can lead to a refined convex problem in \mathbf{w} that is a better convex relaxation than (7).

Our numerical procedure exploits the above fact, which tries to improve the estimation of \mathbf{v}_k over the initial choice of $\mathbf{v}_k = \mathbf{1}$ in (7) using an iterative algorithm. This can be done using an alternating optimization procedure, which repeatedly applies the following two steps:

- First we optimize \mathbf{w} with \mathbf{v} fixed: this is a convex problem in \mathbf{w} with appropriately chosen $\mathbf{h}(\mathbf{w})$.
- Second we optimize \mathbf{v} with \mathbf{w} fixed: although non-convex, it has a closed form solution that is given by (6).

Initialize $\hat{\mathbf{v}} = \mathbf{1}$
Repeat the following two steps until convergence:

- Let
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[R_0(\mathbf{w}) + \sum_{k=1}^K \mathbf{h}_k(\mathbf{w})^\top \hat{\mathbf{v}}_k \right]. \quad (9)$$
- Let $\hat{\mathbf{v}}_k = \nabla_{\mathbf{u}} \bar{R}_k(\mathbf{u})|_{\mathbf{u}=\mathbf{h}_k(\hat{\mathbf{w}})}$ ($k = 1, \dots, K$)

Figure 2: Multi-stage Convex Relaxation Method

The general procedure for solving (8) is presented in Figure 2. It can be regarded as a generalization of CCCP (concave-convex programming) [11], which takes $\mathbf{h}(\mathbf{w}) = \mathbf{w}$. By repeatedly refining the parameter \mathbf{v} , we can potentially obtain better and better convex relaxation, leading to a solution superior to that of the initial convex relaxation. Since at each step the procedure decreases the objective function in (8), its convergence to a local minimum is not difficult to prove. In fact, in order to achieve convergence, one only needs to approximately minimize (9) and reasonably decrease the objective value at each step. We skip the detailed analysis here, because in the general case, a

local solution is not necessarily a good solution, and there are other approaches (such as gradient descent) that can compute a local solution. In order to demonstrate the effectiveness of multi-stage convex relaxation, we shall include a more careful analysis for the special case of sparse regularization in Section 5.1. Our theory shows that the local solution of multi-stage relaxation with a nonconvex sparse regularizer is superior to the convex L_1 regularization solution (under appropriate conditions).

4.2 Constrained formulation

The multi-stage convex relaxation idea can also be used to solve the constrained formulation (4). The one-stage convex relaxation of (4), given fixed relaxation parameter \mathbf{v}_k , becomes

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in R^d} R_0(\mathbf{w}) \quad \text{subject to} \quad \sum_{k=1}^K \mathbf{h}_k(\mathbf{w})^\top \mathbf{v}_k \leq A - \sum_{k=1}^K R_k^*(\mathbf{v}_k).$$

Now by optimizing \mathbf{v} in addition to \mathbf{w} , we obtain the following algorithm:

- Initialize $\hat{\mathbf{v}} = \mathbf{1}$
- Repeat the following two steps until convergence:

– Let

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} R_0(\mathbf{w}) \quad \text{subject to} \quad \sum_{k=1}^K \mathbf{h}_k(\mathbf{w})^\top \hat{\mathbf{v}}_k \leq A - \sum_{k=1}^K R_k^*(\hat{\mathbf{v}}_k).$$

– Let $\hat{\mathbf{v}}_k = \nabla_{\mathbf{u}} \bar{R}_k(\mathbf{u})|_{\mathbf{u}=\mathbf{h}_k(\hat{\mathbf{w}})}$ ($k = 1, \dots, K$)

If an optimization problem includes both nonconvex penalization and nonconvex constraints, then one may use the above algorithm with Figure 2. The constrained formulation is not the focus of this paper. We include it here only to show that it can be handled under the general framework.

4.3 Some Examples of Multi-stage Convex Relaxation Methods

The multi-stage convex relaxation method can be used with examples in Section 3 to obtain concrete algorithms for various formulations. We describe some examples here.

Smoothed classification loss

In this case, the optimization problem is

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[\sum_{i=1}^n \min(\alpha, \max(0, 1 - \mathbf{w}^\top \mathbf{x}_i y_i)) + \lambda g(\mathbf{w}) \right],$$

where we assume that $g(\mathbf{w})$ is a convex regularization condition such as $g(\mathbf{w}) = \lambda \|\mathbf{w}\|_2^2$.

From Section 3, the multi-stage convex relaxation solves the weighted SVM formulation

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[\sum_{i=1}^n \hat{\mathbf{v}}_i \max(0, 1 - \mathbf{w}^\top \mathbf{x}_i y_i) + \lambda g(\mathbf{w}) \right],$$

where the relaxation parameter \mathbf{v} is updated as

$$\hat{\mathbf{v}}_i = I(\hat{\mathbf{w}}^\top \mathbf{x}_i y_i \geq 1 - \alpha) \quad (i = 1, \dots, n).$$

Intuitively, the mis-classified points $\hat{\mathbf{w}}^\top \mathbf{x}_i y_i < 1 - \alpha$ are considered as outliers, and ignored.

L_p and smoothed L_p regularization

We consider the following optimization formulation for some $\alpha \geq 0$ and $p \in (0, 1]$:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[R_0(\mathbf{w}) + \lambda \sum_{j=1}^d (\alpha + |\mathbf{w}_j|)^p \right],$$

where we assume that $R_0(\mathbf{w})$ is a convex function of \mathbf{w} .

From Section 3, the multi-stage convex relaxation becomes a weighted L_q regularization formulation for $q \geq 1$:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[R_0(\mathbf{w}) + \sum_{j=1}^d \hat{\mathbf{v}}_j |\mathbf{w}_j|^q \right],$$

where the relaxation parameter \mathbf{v} is updated as

$$\hat{\mathbf{v}}_j = \lambda(p/q)(\alpha + |\hat{\mathbf{w}}_j|)^{p-1} |\hat{\mathbf{w}}_j|^{1-q} \quad (j = 1, \dots, d).$$

The typical choices of q are $q = 1$ or $q = 2$. That is, we relax L_p regularization to L_1 or L_2 regularization.

A similar formula holds for smoothed log regularization. The resulting procedure is the same as the one empirically studied in [6]. Finally, we note that the two stage version of L_p regularization, relaxed to L_q with $q = 1$, is referred to Adaptive-Lasso [16].

Capped L_1 regularization

In capped L_1 regularization, we consider the optimization problem

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[R_0(\mathbf{w}) + \lambda \sum_{j=1}^d \min(\alpha, |\mathbf{w}_j|) \right],$$

where we assume that $R_0(\mathbf{w})$ is a convex function of \mathbf{w} .

From Section 3, the multi-stage convex relaxation becomes a weighted L_1 regularization formulation:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[R_0(\mathbf{w}) + \sum_{j=1}^d \hat{\mathbf{v}}_j |\mathbf{w}_j| \right],$$

where the relaxation parameter \mathbf{v} is updated as

$$\hat{\mathbf{v}}_j = \lambda I(|\hat{\mathbf{w}}_j| \leq \alpha) \quad (j = 1, \dots, d).$$

This method has an intuitive interpretation: in order to achieve sparsity, the standard L_1 regularization not only shrinks small coefficients to zero, but also shrinks large coefficients. This causes a bias. The capped- L_1 formulation removes the bias by adaptively adjusting the relaxation parameter $\hat{\mathbf{v}}_j$ so that if $|\hat{\mathbf{w}}_j|$ is large, then we do not penalize the corresponding variable j .

Sparse Eigenvalue Problem

We use a simple example to illustrate that the multi-stage convex relaxation idea does not only apply to formulations with convex risks. Consider the sparse eigenvalue problem, where we are interested in finding the largest eigenvalue of a positive semi-definite matrix A . One formulation is

$$\hat{\mathbf{w}} = \arg \max_{\|\mathbf{w}\|_2 \leq 1} \left[\mathbf{w}^\top A \mathbf{w} - \lambda \sum_{j=1}^d (\alpha + |\mathbf{w}_j|)^p \right],$$

with parameter $p \in (0, 1)$ and a small parameter $\alpha > 0$ to encourage sparsity. If $\lambda = 0$, then it is the standard eigenvalue problem without sparsity constraints. Although the standard eigenvalue problem is not convex in \mathbf{w} , it has a convex relaxation to a semi-definite programming problem, and thus can be efficiently solved. For convenience, we think of the standard eigenvalue problem as “convex” for the purpose of this paper. The multi-stage convex relaxation becomes:

$$\hat{\mathbf{w}} = \arg \max_{\|\mathbf{w}\|_2 \leq 1} \left[\mathbf{w}^\top A \mathbf{w} - \sum_{j=1}^d \mathbf{v}_j \mathbf{w}_j^2 \right],$$

which is a standard eigenvalue problem. The relaxation parameter is updated as

$$\hat{\mathbf{v}}_j = \lambda(p/2)(\alpha + |\hat{\mathbf{w}}_j|)^{p-1} |\hat{\mathbf{w}}_j|^{-1} \quad (j = 1, \dots, d).$$

Matrix Regularization

Our final example is multi-task learning with matrix regularization, also considered in [1]. In this case, \mathbf{w} is not a vector, but a matrix, with columns (tasks) \mathbf{w}^ℓ . We solve a problem of the following form:

$$\mathbf{w} = \arg \min_{\mathbf{w}} \left[\sum_{\ell=1}^m R^\ell(\mathbf{w}^\ell) + \lambda \text{tr}((\alpha I + \mathbf{w} \mathbf{w}^\top)^{p/2}) \right].$$

In the above formulation, R^ℓ is the risk function for task ℓ . The matrix regularization used here is the counterpart of L_p regularization for vectors. It encourages low-rank if $p < 2$. In particular, the case of $p = 1$ is often called trace norm (or nuclear norm). It is convex and frequently used in the literature. The parameter $\alpha > 0$ gives some smoothness, similar to the vector case.

The case of $p < 1$ gives better low-rank approximation, similar to the vector regularization case. Again, this problem can be solved with multi-stage convex relaxation method. In this case, the relaxation parameter \mathbf{v} is positive semi-definite matrix. And we relax the regularization term to $\mathbf{h}(\mathbf{w}) = (\alpha I + \mathbf{w} \mathbf{w}^\top)$ as a matrix, and thus, the relaxed regularization term becomes $\text{tr}(\mathbf{v}(\alpha I + \mathbf{w} \mathbf{w}^\top))$. This regularization decouples the problems as the following problem, which allows us to solve each task ℓ separately:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[R^\ell(\mathbf{w}^\ell) + (\mathbf{w}^\ell)^\top \hat{\mathbf{v}} \mathbf{w}^\ell \right] \quad (\ell = 1, 2, \dots, m).$$

This is a key advantage of the method. Similar to the vector case, we have the following update formula for the relaxation parameter:

$$\hat{\mathbf{v}} = \lambda(p/2)(\alpha I + \hat{\mathbf{w}} \hat{\mathbf{w}}^\top)^{(p-2)/2}.$$

5 Multi-stage Convex Relaxation for Sparse Regularization

In recent years, convex relaxation has become a major theme in machine learning. The essential argument for convexity is that global optimal solutions can be computed. Moreover, in some cases, strong theoretical results can be obtained for convex relaxation solutions showing that they approximately solve the original nonconvex problem. A specific example is sparse learning (2), where various recent work showed that L_1 relaxation approximately solves the original sparse learning problem (2).

The purpose of this paper is to show that although convex relaxation has been successfully applied in machine learning, it is not necessarily the optimal approach for solving nonconvex problems. The multi-stage convex relaxation method tries to obtain better approximations of the original nonconvex problem by refining the convex relaxation formulation. Since the local solution found by the algorithm is the global solution of a refined convex relaxation formulation, it should be closer to the desired solution than that of the standard one-stage convex relaxation method. Although this high level intuition is appealing, it is still necessary to present a more rigorous theoretical result which can precisely demonstrate the advantage of the multi-stage approach over the standard single stage method. Unless we can develop a theory to show the effectiveness of the multi-stage procedure in Figure 2, our proposal is yet another local minimum finding scheme that may potentially get stuck at a bad local solution.

We shall point out that the framework presented in this work is quite general, and for many problems, even the behavior of one-stage convex relaxation is not well-understood. In order to obtain some strong theoretical results that can demonstrate our points, we consider the special case of sparse learning. This is because this problem has been well-studied in recent years, and the behavior of convex relaxation (L_1 regularization) is well-understood.

5.1 Theory of Sparse Regularization

For a non-convex but smooth regularization condition such as capped- L_1 or smoothed- L_p with $p \in (0, 1)$, standard numerical techniques such as gradient descent leads to a local minimum solution. Unfortunately, it is difficult to find the global optimum, and it is also difficult to analyze the quality of the local minimum. Although in practice, such a local minimum solution may outperform the Lasso solution, the lack of theoretical (and practical) performance guarantee prevents the more widespread applications of such algorithms. As a matter of fact, results with non-convex regularization are difficult to reproduce because different numerical optimization procedures can lead to different local minima. Therefore the quality of the solution heavily depend on the numerical procedure used.

The situation is very difficult for a convex relaxation formulation such as L_1 -regularization (Lasso). The global optimum can be easily computed using standard convex programming techniques. It is known that in practice, 1-norm regularization often leads to sparse solutions (although often suboptimal). Moreover, its performance has been theoretically analyzed recently. For example, it is known from the compressed sensing literature that under certain conditions, the solution of L_1 relaxation may be equivalent to L_0 regularization asymptotically (e.g. [5]). If the target is truly sparse, then it was shown in [15] that under some restrictive conditions referred to as *irrepresentable conditions*, 1-norm regularization solves the feature selection problem. The prediction performance of this method has been considered in [9, 14, 3, 4].

In spite of its success, L_1 -regularization often leads to suboptimal solutions because it is not a good approximation to L_0 regularization. Statistically, this means that even though it converges

to the true sparse target when $n \rightarrow \infty$ (consistency), the rate of convergence can be suboptimal. The only way to fix this problem is to employ a non-convex regularization condition that is closer to L_0 regularization. In the following, we formally prove a result for multi-stage convex relaxation with non-convex sparse regularization that is superior to the Lasso result. In essence, we establish a performance guarantee for non-convex formulations when they are solved by using the multi-stage convex relaxation approach which is more sophisticated than the standard one-stage convex relaxation.

In supervised learning, we observe a set of input vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in R^d$, with corresponding desired output variables y_1, \dots, y_n . In general, we may assume that there exists a target $\bar{\mathbf{w}} \in R^d$ such that

$$y_i = \bar{\mathbf{w}}^\top \mathbf{x}_i + \epsilon_i \quad (i = 1, \dots, n),$$

where ϵ_i are zero-mean independent random noises (but not necessarily identically distributed). Moreover, we assume that the target vector $\bar{\mathbf{w}}$ is sparse. That is, there exists $\bar{k} = \|\bar{\mathbf{w}}\|_0$ is small. This is the standard statistical model for sparse learning.

Let \mathbf{y} denote the vector of $[y_i]$ and X be the $n \times d$ matrix with each row a vector \mathbf{x}_i . We are interested in recovering $\bar{\mathbf{w}}$ from noisy observations using the following sparse regression method:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left[\frac{1}{n} \|X\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \sum_{j=1}^d g(|\mathbf{w}_j|) \right], \quad (10)$$

where $g(|\mathbf{w}_j|)$ is a regularization function. Here we require that $g'(u)$ is a non-negative which means we penalize larger u more significantly. Moreover, we assume $g'(u)$ is a non-increasing function when $u > 0$, which means that $[g(|\mathbf{w}_1|), \dots, g(|\mathbf{w}_d|)]$ is concave with respect to $\mathbf{h}(\mathbf{w}) = [|\mathbf{w}_1|, \dots, |\mathbf{w}_d|]$. It follows that (10) can be solved using the multi-stage convex relaxation algorithm in Figure 3, which we will analyze.

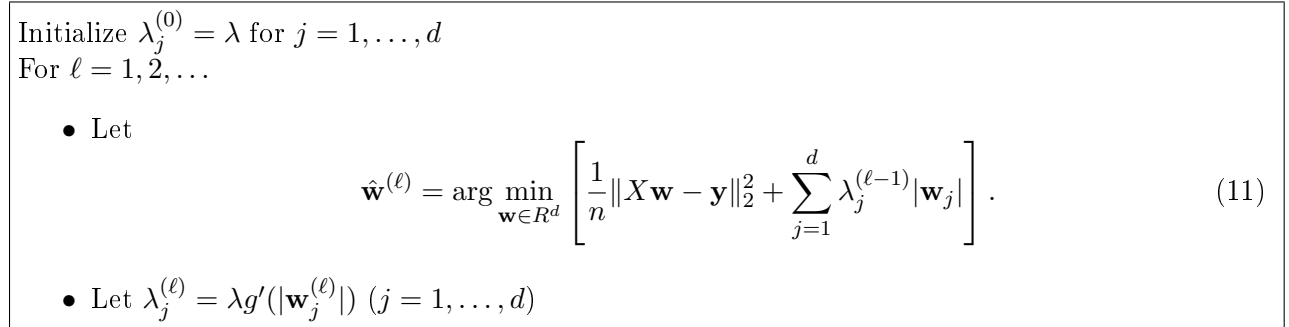


Figure 3: Multi-stage Convex Relaxation for Sparse Regularization

For convenience, we consider fixed design only, where X is fixed and the randomness is with respect to \mathbf{y} only. We require some technical conditions for our analysis. First we assume sub-Gaussian noise as follows.

Assumption 5.1 *Assume that $\{\epsilon_i\}_{i=1, \dots, n}$ are independent (but not necessarily identically distributed) sub-Gaussians: there exists $\sigma \geq 0$ such that $\forall i$ and $\forall t \in R$,*

$$\mathbf{E}_{\epsilon_i} e^{t\epsilon_i} \leq e^{\sigma^2 t^2 / 2}.$$

Both Gaussian and bounded random variables are sub-Gaussian using the above definition. For example, if a random variable $\xi \in [a, b]$, then $\mathbf{E}_\xi e^{t(\xi - \mathbf{E}\xi)} \leq e^{(b-a)^2 t^2 / 8}$. If a random variable is Gaussian: $\xi \sim N(0, \sigma^2)$, then $\mathbf{E}_\xi e^{t\xi} \leq e^{\sigma^2 t^2 / 2}$.

We also introduce the concept of sparse eigenvalue, which is standard in the analysis of L_1 regularization.

Definition 5.1 Given k , define

$$\begin{aligned}\rho_+(k) &= \sup \left\{ \frac{1}{n} \|X\mathbf{w}\|_2^2 / \|\mathbf{w}\|_2^2 : \|\mathbf{w}\|_0 \leq k \right\} \\ \rho_-(k) &= \inf \left\{ \frac{1}{n} \|X\mathbf{w}\|_2^2 / \|\mathbf{w}\|_2^2 : \|\mathbf{w}\|_0 \leq k \right\}\end{aligned}$$

Our main result is stated as follows. The proof is in the appendix.

Theorem 5.1 Let Assumption 5.1 hold. Assume also that the target $\bar{\mathbf{w}}$ is sparse, with $\mathbf{E}y_i = \bar{\mathbf{w}}^\top \mathbf{x}_i$, and $\bar{k} = \|\bar{\mathbf{w}}\|_0$. Choose λ such that

$$\lambda \geq 20\sigma \sqrt{2\rho_+(1) \ln(2d/\eta)/n}.$$

Assume that $g'(z) \geq 0$ is a non-increasing function such that $g'(z) = 1$ when $z \leq 0$. Moreover, we require that $g'(\theta) \geq 0.9$ with $\theta = 9\lambda/\rho_-(2\bar{k} + \ell)$. Assume that $\rho_+(\ell)/\rho_-(2\bar{k} + 2\ell) \leq 1 + 0.5\ell/\bar{k}$ for some $\ell \geq 2\bar{k}$, then

$$\begin{aligned}\|\Delta \hat{\mathbf{w}}^{(\ell)}\|_2 &\leq \frac{17}{\rho_-(2\bar{k} + \ell)} \left[2\sigma \sqrt{\rho_+(\bar{k})} \left(\sqrt{7.4\bar{k}/n} + \sqrt{2.7 \ln(2/\eta)/n} \right) + \lambda \left(\sum_{j: \bar{\mathbf{w}}_j \neq 0} g'(|\bar{\mathbf{w}}_j| - \theta)^2 \right)^{1/2} \right] \\ &\quad + 0.7^\ell \frac{10}{\rho_-(2\bar{k} + \ell)} \sqrt{\bar{k}} \lambda.\end{aligned}$$

Note that the theorem allows the situation $d \gg n$, which is what we are interested in. This is the first general analysis of multi-stage convex relaxation for high dimensional sparse learning, although some simpler results for low dimensional two-stage procedures were obtained in [16, 17].

The condition $\rho_+(\ell)/\rho_-(2\bar{k} + 2\ell) \leq 1 + 0.5\ell/\bar{k}$ requires the eigenvalue ratio $\rho_+(\ell)/\rho_-(\ell)$ to grow sub-linearly in ℓ . Such a condition, referred to as *sparse eigenvalue condition*, is also needed in the standard analysis of L_1 regularization [12, 14]. It is related but weaker than the *restricted isometry property* (RIP) in compressive sensing [5]. The theorem yields important insights into the behavior of multi-stage convex relaxation. Since for Lasso, $g'(|\mathbf{w}_j|) \equiv 1$, we obtain the following bound

$$\|\hat{\mathbf{w}}_{L_1} - \bar{\mathbf{w}}\|_2 = O(\sqrt{\bar{k}}\lambda),$$

where $\hat{\mathbf{w}}_{L_1}$ is the solution of the standard L_1 regularization. This bound is tight for Lasso, in the sense that the right hand side cannot be improved except for the constant — this can be easily verified with an orthogonal design matrix. It is known that in order for Lasso to be effective, one has to pick λ no smaller than the order $\sigma\sqrt{\ln d/n}$. Therefore, the parameter estimation error of the standard Lasso is of the order $\sigma\sqrt{\bar{k} \ln d/n}$, which cannot be improved.

In comparison, if we consider an appropriate regularization condition $g(|\mathbf{w}_j|)$ that is concave in $|\mathbf{w}_j|$. Since $g'(|\mathbf{w}_j|) \approx 0$ when $|\mathbf{w}_j|$ is large, the bound in Theorem 5.1 can be significantly better

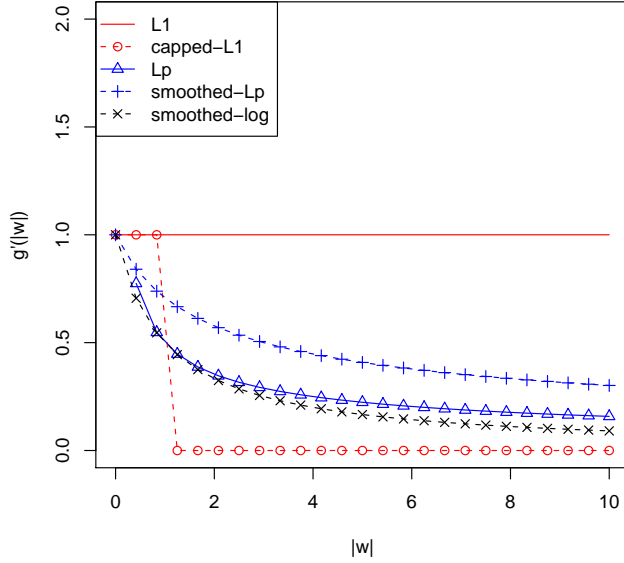


Figure 4: Derivative $g'(|\mathbf{w}_j|)$ of some sparse regularizers

when most non-zero coefficients of $\bar{\mathbf{w}}$ are relatively large in magnitude. For example, consider the capped- L_1 regularizer $g(|\mathbf{w}_j|) = \min(\alpha, |\mathbf{w}_j|)$ with $\alpha \geq \theta$; in the extreme case where $\min_j \|\mathbf{w}_j\| > \alpha + \theta$ (which can be achieved when all nonzero components of $\bar{\mathbf{w}}$ are larger than $O(\sigma\sqrt{\ln d/n})$), we obtain the better bound

$$\|\hat{\mathbf{w}}^{(\ell)} - \bar{\mathbf{w}}\|_2 = O(\sqrt{\bar{k}/n} + \sqrt{\ln(1/\eta)/n})$$

for the multi-stage procedure as $\ell \rightarrow \infty$. This is superior to the standard one-stage L_1 regularization bound $\|\hat{\mathbf{w}}_{L_1} - \bar{\mathbf{w}}\|_2 = O(\sqrt{\bar{k} \ln(d/\eta)/n})$, which is tight for Lasso. The difference can be significant when d is large.

Generally speaking, with a regularization condition $g(|\mathbf{w}_j|)$ that is concave in $|\mathbf{w}_j|$, the dependency on λ is through $g'(|\bar{\mathbf{w}}_j|)$ which decreases as $|\bar{\mathbf{w}}_j|$ increases. This removes the bias of the Lasso and leads to improved performance. Specifically, if $\bar{\mathbf{w}}_j$ is large, then $g'(|\bar{\mathbf{w}}_j|) \approx 0$. In comparison, the Lasso bias is due to the fact that $g'(|\bar{\mathbf{w}}_j|) \equiv 1$. For illustration, the derivative $g'(\cdot)$ of some sparse regularizers are plotted in Figure 4.

Note that our theorem only applies to regularizers with finite derivative at zero. That is, $g'(0) < \infty$. The result doesn't apply to L_p regularization with $p < 1$ because $g'(0) = \infty$. Although a weaker result can be obtained for such regularizers, we do not include it here. We only include an intuitive example below to illustrate why the condition $g'(0) < \infty$ is necessary for stronger results presented in the paper. Observe that the multi-stage convex relaxation method only computes a local minimum, and the regularization update rule is given by $\lambda_j^{(\ell-1)} = g'(\hat{\mathbf{w}}_j^{(\ell-1)})$. If $g'(0) = \infty$, then $\lambda_j^{(\ell-1)} = \infty$ when $\hat{\mathbf{w}}_j^{(\ell-1)} = 0$. This means that if a feature accidentally becomes zero in some stage, it will always remain zero. This is why only weaker results can be obtained for L_p regularizers ($p < 1$): we need to further assume that $\hat{\mathbf{w}}_j^{(\ell)}$ never becomes close to zero when $\bar{\mathbf{w}}_j \neq 0$. A toy

Stage ℓ	coefficients	$\ \hat{\mathbf{w}}^{(\ell)} - \bar{\mathbf{w}}\ _2$
multi-stage capped- L_1		
1	[6.0, 0.0, 4.7, 4.8, 3.9, 0.6, 0.7, 1.2, 0.0, ...]	4.4
2	[7.7, 0.4, 5.7, 6.3, 5.7, 0.0, 0.0, 0.2, 0.0, ...]	1.6
3	[7.8, 1.2, 5.7, 6.6, 5.7, 0.0, 0.0, 0.0, 0.0, ...]	0.98
4	[7.8, 1.2, 5.7, 6.6, 5.7, 0.0, 0.0, 0.0, 0.0, ...]	0.98
multi-stage $L_{0.5}$		
1	[6.0, 0.0, 4.7, 4.8, 3.9, 0.6, 0.7, 1.2, 0.0, ...]	4.4
2	[7.3, 0.0, 5.4, 5.9, 5.3, 0.0, 0.3, 0.3, 0.0, 0.0, ...]	2.4
3	[7.5, 0.0, 5.6, 6.1, 5.7, 0.0, 0.1, 0.0, 0.0, 0.0, ...]	2.2
4	[7.5, 0.0, 5.6, 6.2, 5.7, 0.0, 0.1, 0.0, 0.0, 0.0, ...]	2.1
target $\bar{\mathbf{w}}$	[8.2, 1.7, 5.4, 6.9, 5.7, 0.0, 0.0, 0.0, 0.0, ...]	

Table 1: An Illustrative Example for Multi-stage Sparse Regularization

example is presented in Table 1 to demonstrate this point. The example is a simulated regression problem with $d = 500$ variables and $n = 100$ training data. The first five variables of the target $\bar{\mathbf{w}}$ are non-zeros, and the remaining variables are zeros. For both capped- L_1 and L_p regularizers, the first stage is the standard L_1 regularization, which misses the correct feature #2 and wrongly selects some incorrect ones. For capped- L_1 regularization, in the second stage, because most correct features are identified, the corresponding “bias” is reduced by not penalizing the corresponding variables. This leads to improved performance. Since the correct feature #2 shows up in stage 2, we are able to identify it and further improve the convex relaxation in stage 3. After stage 3, the procedure stabilizes because it computes exactly the same relaxation. For L_p regularization, since feature #2 becomes zero in stage 1, it will remain zero thereafter because $\lambda_2^{(\ell)} = \infty$ when $\ell \geq 1$. In order to remedy this problem, one has to use a regularizer with $g'(0) < \infty$ such as the smoothed L_p regularizer.

5.2 Empirical Study

Although this paper focuses on the development of the general multi-stage convex relaxation framework as well as its theoretical understanding (in particular the major result given in Theorem 5.1), we include two simple numerical examples to verify our theory. More comprehensive empirical comparisons can be found in other related work such as [6, 16, 17].

In order to avoid cluttering, we only present results with capped- L_1 and L_p ($p = 0.5$) regularization methods. Note that based on Theorem 5.1, we may tune α in capped- L_1 by using a formula $\alpha = \alpha_0 \lambda$ where λ is the regularization parameter. We choose $\alpha_0 = 10$ and $\alpha_0 = 100$.

In the first experiment, we generate an $n \times d$ random matrix with its column j corresponding to $[\mathbf{x}_{1,j}, \dots, \mathbf{x}_{n,j}]$, and each element of the matrix is an independent standard Gaussian $N(0, 1)$. We then normalize its columns so that $\sum_{i=1}^n \mathbf{x}_{i,j}^2 = n$. A truly sparse target $\bar{\beta}$, is generated with k nonzero elements that are uniformly distributed from $[-10, 10]$. The observation $\mathbf{y}_i = \bar{\beta}^\top \mathbf{x}_i + \epsilon_i$, where each $\epsilon_i \sim N(0, \sigma^2)$. In this experiment, we take $n = 50, d = 200, k = 5, \sigma = 1$, and repeat the experiment 100 times. The average training error and 2-norm parameter estimation error are reported in Figure 5. We compare the performance of multi-stage methods with different regularization parameter λ . As expected, the training error for the multi-stage algorithms are smaller

than that of L_1 , due to the smaller bias. Moreover, substantially smaller parameter estimation error is achieved by the multi-stage procedures, which is consistent with Theorem 5.1. This can be regarded as an empirical verification of the theoretical result.

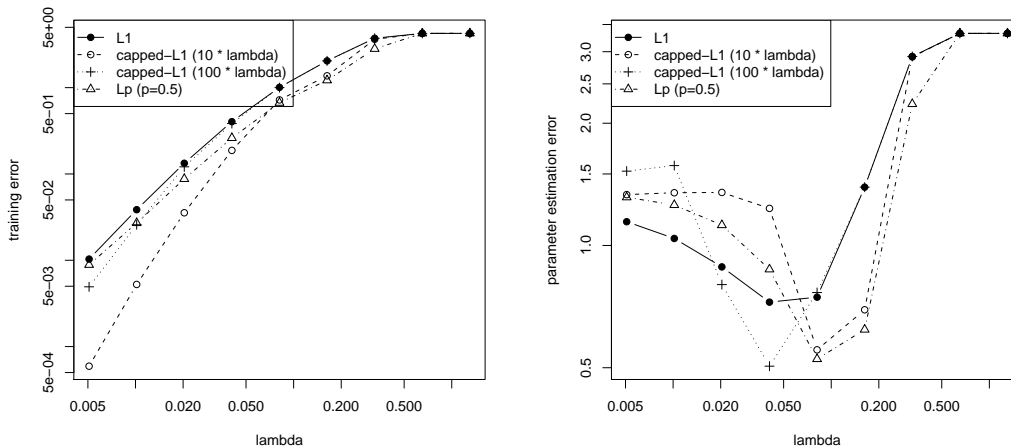


Figure 5: Performance of multi-stage convex relaxation on simulation data. Left: average training squared error versus λ ; Right: parameter estimation error versus λ .

In the second experiment, we use the *Boston Housing* data to illustrate the effectiveness of multi-stage convex relaxation. This data set contains 506 census tracts of Boston from the 1970 census, available from the *UCI Machine Learning Database Repository*: <http://archive.ics.uci.edu/ml/>. Each census tract is a data-point, with 13 features (we add a constant offset on e as the 14th feature), and the desired output is the housing price. In this example, we randomly partition the data into 20 training plus 486 test points. We perform the experiments 100 times, and report training and test squared error versus the regularization parameter λ for different q . The results are plotted in Figure 6. In this case, $L_{0.5}$ is not effective, while capped- L_1 regularization with $\alpha = 100\lambda$ is slightly better than Lasso. Note that this dataset contains only a small number ($d = 14$) features, which is not the case where we can expect significant benefit from the multi-stage approach (most of other UCI data similarly contain only small number of features). In order to illustrate the advantage of the multi-stage method more clearly, we also report results on a modified Boston Housing data, where we append 20 random features (similar to the simulation experiments) to the original Boston Housing data, and rerun the experiments. The results are shown in Figure 7. As expected from Theorem 5.1 and the discussion thereafter, since d becomes large, the multi-stage convex relaxation approach with capped- L_1 regularization and $L_{0.5}$ regularization perform significantly better than the standard Lasso.

6 Discussion

Many machine learning applications require solving nonconvex optimization problems. There are two approaches to this problem:

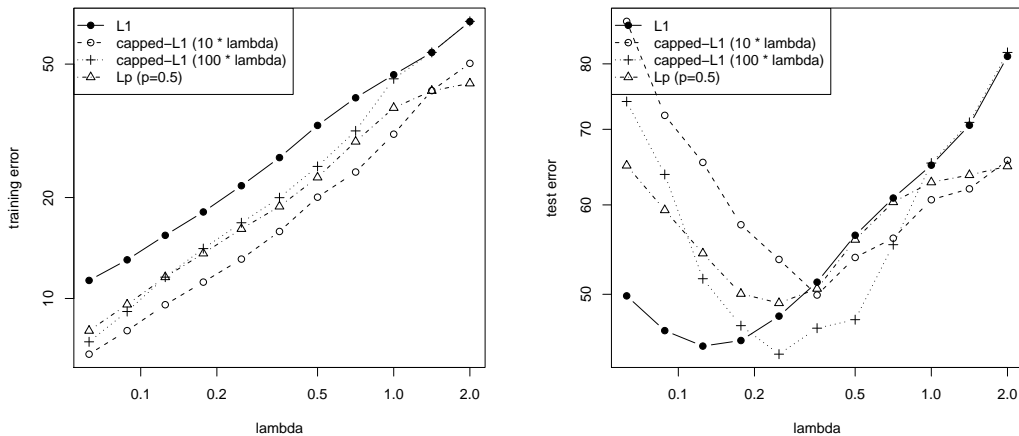


Figure 6: Performance of multi-stage convex relaxation on the original Boston Housing data. Left: average training squared error versus λ ; Right: test squared error versus λ .

- Heuristic methods such as gradient descent that only find a local minimum. A drawback of this approach is the lack of theoretical guarantee showing that the local minimum gives a good solution.
- Convex relaxation such as L_1 -regularization that solves the problem under some conditions. However it often leads to a sub-optimal solution in reality.

The goal of this paper is to remedy the above gap between theory and practice. In particular, we investigated a multi-stage convex relaxation scheme for solving problems with non-convex objective functions. The general algorithmic technique is presented first, which can be applied to a wide range of problems. The intuition is to refine convex relaxation iteratively by using solutions obtained from earlier stages. This leads to better and better convex relaxation formulations, and thus better and better solutions.

Although the scheme only finds a local minimum, the above argument indicates that the local minimum it finds should be closer to the original nonconvex problem than the standard convex relaxation solution. In order to prove the effectiveness of this approach theoretically, we considered the sparse learning problem where the behavior of convex relaxation (Lasso) has been well studied in recent years. We showed that under appropriate conditions, the local solution from the multi-stage convex relaxation algorithm is superior to the global solution of the standard L_1 convex relaxation for learning sparse targets. Experiments confirmed the effectiveness of this method.

Finally we shall mention that our theory only shows that nonconvex regularization behaves better than Lasso under appropriate sparse eigenvalue conditions. When such conditions hold, multi-stage convex relaxation is superior. On the other hand, when such conditions fail, neither Lasso nor (the local solution of) multi-stage convex relaxation can be shown to work well. In such case, some features will become highly correlated, and local solutions of non-convex formulations may become unstable. In order to improve stability, it may be helpful to employ ensemble methods such as bagging. Since our analysis doesn't yield any insights in this scenario, further theoretical

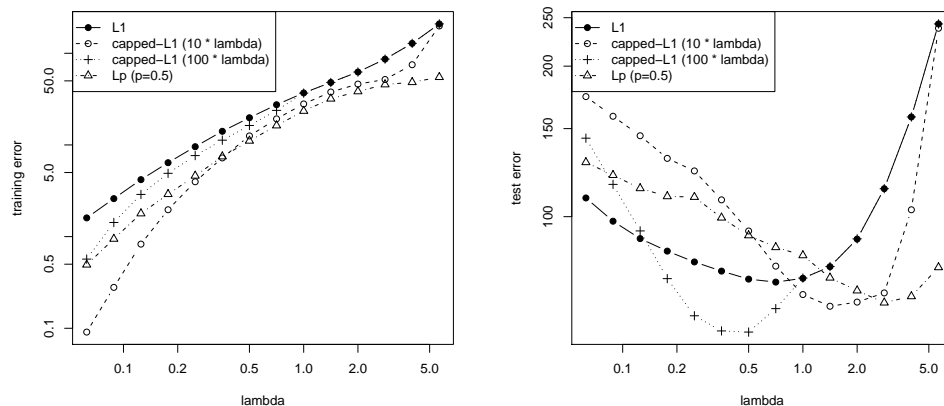


Figure 7: Performance of multi-stage convex relaxation on the modified Boston Housing data. Left: average training squared error versus λ ; Right: test squared error versus λ .

investigation is necessary.

References

- [1] Andreas Argyriou, Charles A. Micchelli, Massimiliano Pontil, and Yiming Ying. A spectral regularization framework for multi-task structure learning. In *NIPS'07*, 2008.
- [2] P.L. Bartlett, M.I. Jordan, and J.D. McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- [3] Peter Bickel, Yaacov Ritov, and Alexandre Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, 2008. to appear.
- [4] Florentina Bunea, Alexandre Tsybakov, and Marten H. Wegkamp. Sparsity oracle inequalities for the Lasso. *Electronic Journal of Statistics*, 1:169–194, 2007.
- [5] Emmanuel J. Candes and Terence Tao. Decoding by linear programming. *IEEE Trans. on Information Theory*, 51:4203–4215, 2005.
- [6] Emmanuel J. Candes, Michael B. Wakin, and Stephen P. Boyd. Enhancing sparsity by reweighted l_1 minimization. *Journal of Fourier Analysis and Applications*, 14(5):877–905, 2008.
- [7] David L. Donoho, Michael Elad, and Vladimir N. Temlyakov. Stable recovery of sparse over-complete representations in the presence of noise. *IEEE Trans. Info. Theory*, 52(1):6–18, 2006.
- [8] Junzhou Huang, Tong Zhang, and Dimitris Metaxas. Learning with structured sparsity. Technical report, Rutgers University, January 2009. A short version appears in ICML'09. Available from <http://arxiv.org/abs/0903.3002>.
- [9] Vladimir Koltchinskii. Sparsity in penalized empirical risk minimization. *Annales de l'Institut Henri Poincaré*, 2008.

- [10] R. Tyrrell Rockafellar. *Convex analysis*. Princeton University Press, Princeton, NJ, 1970.
- [11] Alan L. Yuille and Anand Rangarajan. The concave-convex procedure. *Neural Computation*, 15:915–936, 2003.
- [12] Cun-Hui Zhang and Jian Huang. The sparsity and bias of the lasso selection in high-dimensional linear regression. *Annals of Statistics*, 36(4):1567–1594, 2008.
- [13] Tong Zhang. Statistical behavior and consistency of classification methods based on convex risk minimization. *The Annals of Statistics*, 32:56–85, 2004. with discussion.
- [14] Tong Zhang. Some sharp performance bounds for least squares regression with L_1 regularization. *The Annals of Statistics*, 2009. to appear.
- [15] Peng Zhao and Bin Yu. On model selection consistency of Lasso. *Journal of Machine Learning Research*, 7:2541–2567, 2006.
- [16] Hui Zou. The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101:1418–1429, 2006.
- [17] Hui Zou and Runze Li. One-step sparse estimates in nonconcave penalized likelihood models. *Annals of Statistics*, 36(4):1509–1533, 2008.

A Proof of Theorem 5.1

The analysis is an adaptation of [14]. We first introduce some definitions. Consider the positive semi-definite matrix $A = n^{-1}X^\top X \in \mathbb{R}^{d \times d}$. Given $\ell, k \geq 1$ such that $\ell + k \leq d$. Let I, J be disjoint subsets of $\{1, \dots, d\}$ with k and ℓ elements respectively. Let $A_{I,J} \in \mathbb{R}^{k \times \ell}$ be the restriction of A to indices I, J , $A_{I,I} \in \mathbb{R}^{k \times k}$ be the restriction of A to indices I on the left and I on the right. Similarly we define restriction \mathbf{w}_I of a vector $\mathbf{w} \in \mathbb{R}^d$ on I ; and for convenience, we allow either $\mathbf{w}_I \in \mathbb{R}^k$ or $\mathbf{w}_I \in \mathbb{R}^d$ (where components not in I are zeros) depending on the context.

We also need the following quantity in our analysis:

$$\pi(k, \ell) = \sup_{\mathbf{v} \in \mathbb{R}^k, \mathbf{u} \in \mathbb{R}^{\ell, I, J}} \frac{\mathbf{v}^\top A_{I, J} \mathbf{u} \|\mathbf{v}\|_2}{\mathbf{v}^\top A_{I, I} \mathbf{v} \|\mathbf{u}\|_\infty}.$$

The following two lemmas are taken from [14]. We skip the proof.

Lemma A.1 *The following inequality holds:*

$$\pi(k, \ell) \leq \frac{\ell^{1/2}}{2} \sqrt{\rho_+(\ell) / \rho_-(k + \ell) - 1},$$

Lemma A.2 *Consider $k, \ell > 0$ and $G \subset \{1, \dots, d\}$ such that $|G^c| = k$. Given any $\mathbf{w} \in \mathbb{R}^d$. Let J be the indices of the ℓ largest components of \mathbf{w}_G (in absolute values), and $I = G^c \cup J$. Then*

$$\max(0, \mathbf{w}_I^\top A \mathbf{w}) \geq \rho_-(k + \ell) (\|\mathbf{w}_I\|_2 - \pi(k + \ell, \ell) \|\mathbf{w}_G\|_1 / \ell) \|\mathbf{w}_I\|_2.$$

The following lemma gives bounds for sub-Gaussian noise which are needed in our analysis. The proof can be found in [8]. Again we skip the derivation.

Lemma A.3 Define $\hat{\epsilon} = \frac{1}{n} \sum_{i=1}^n (\bar{\mathbf{w}}^\top \mathbf{x}_i - y_i) \mathbf{x}_i$. Under the conditions of Assumption 5.1, with probability larger than $1 - \eta$:

$$\|\hat{\epsilon}\|_\infty^2 \leq 2\sigma^2 \rho_+(1) \ln(2d/\eta)/n.$$

Moreover, for any fixed F , with probability larger than $1 - \eta$:

$$\|\hat{\epsilon}_F\|_2^2 \leq \rho_+(|F|)\sigma^2[7.4|F| + 2.7 \ln(2/\eta)]/n.$$

Lemma A.4 Consider $\bar{\mathbf{w}}$ such that $\{j : \bar{\mathbf{w}}_j \neq 0\} \subset F$ and $F \cap G = \emptyset$. Let $\hat{\mathbf{w}} = \hat{\mathbf{w}}^{(\ell)}$ be the solution of (11), and let $\Delta \hat{\mathbf{w}} = \hat{\mathbf{w}} - \bar{\mathbf{w}}$. Let $\lambda_G = \min_{j \in G} \lambda_j^{(\ell-1)}$ and $\lambda_0 = \max_j \lambda_j^{(\ell-1)}$. Then

$$\sum_{j \in G} |\hat{\mathbf{w}}_j| \leq \frac{2\|\hat{\epsilon}\|_\infty}{\lambda_G - 2\|\hat{\epsilon}\|_\infty} \sum_{j \notin F \cup G} |\hat{\mathbf{w}}_j| + \frac{2\|\hat{\epsilon}\|_\infty + \lambda_0}{\lambda_G - 2\|\hat{\epsilon}\|_\infty} \sum_{j \in F} |\Delta \hat{\mathbf{w}}_j|.$$

Proof For simplicity, let $\lambda_j = \lambda_j^{(\ell-1)}$. The first order equation implies that

$$\frac{1}{n} \sum_{i=1}^n 2(\mathbf{x}_i^\top \mathbf{w} - y_i) \mathbf{x}_{i,j} + \lambda_j \text{sgn}(\mathbf{w}_j) = 0,$$

where $\text{sgn}(\mathbf{w}_j) = 1$ when $\mathbf{w}_j > 0$, $\text{sgn}(\mathbf{w}_j) = -1$ when $\mathbf{w}_j < 0$, and $\text{sgn}(\mathbf{w}_j) \in [-1, 1]$ when $\mathbf{w}_j = 0$. This implies that for all $\mathbf{v} \in \mathbb{R}^d$, we have

$$2\mathbf{v}^\top A \Delta \hat{\mathbf{w}} \leq -2\mathbf{v}^\top \hat{\epsilon} - \sum_{j=1}^d \lambda_j \mathbf{v}_j \text{sgn}(\hat{\mathbf{w}}_j). \quad (12)$$

Now, let $\mathbf{v} = \Delta \hat{\mathbf{w}}$ in (12), we obtain

$$\begin{aligned} 0 &\leq 2\Delta \hat{\mathbf{w}}^\top A \Delta \hat{\mathbf{w}} \leq 2|\Delta \hat{\mathbf{w}}^\top \hat{\epsilon}| - \sum_{j=1}^d \lambda_j \Delta \hat{\mathbf{w}}_j \text{sgn}(\hat{\mathbf{w}}_j) \\ &\leq 2\|\Delta \hat{\mathbf{w}}\|_1 \|\hat{\epsilon}\|_\infty - \sum_{j \in F} \lambda_j \Delta \hat{\mathbf{w}}_j \text{sgn}(\hat{\mathbf{w}}_j) - \sum_{j \notin F} \lambda_j \Delta \hat{\mathbf{w}}_j \text{sgn}(\hat{\mathbf{w}}_j) \\ &\leq 2\|\Delta \hat{\mathbf{w}}\|_1 \|\hat{\epsilon}\|_\infty + \sum_{j \in F} \lambda_j |\Delta \hat{\mathbf{w}}_j| - \sum_{j \notin F} \lambda_j |\hat{\mathbf{w}}_j| \\ &\leq \sum_{j \in G} (2\|\hat{\epsilon}\|_\infty - \lambda_G) |\hat{\mathbf{w}}_j| + \sum_{j \notin G \cup F} 2\|\hat{\epsilon}\|_\infty |\hat{\mathbf{w}}_j| + \sum_{j \in F} (2\|\hat{\epsilon}\|_\infty + \lambda_0) |\Delta \hat{\mathbf{w}}_j|. \end{aligned}$$

By rearranging the above inequality, we obtain the desired bound. ■

Lemma A.5 Using the notations of Lemma A.4, and let J be the indices of the largest ℓ coefficients (in absolute value) of $\hat{\mathbf{w}}_G$. Let $I = G^c \cup J$ and $k = |G^c|$. If $(\lambda_0 + 2\|\hat{\epsilon}\|_\infty)/(\lambda_G - 2\|\hat{\epsilon}\|_\infty) \leq 3$, then

$$\|\Delta \hat{\mathbf{w}}\|_2 \leq (1 + (3k/\ell)^{0.5}) \|\Delta \hat{\mathbf{w}}_I\|_2.$$

Proof Using $(\lambda_0 + 2\|\hat{\epsilon}\|_\infty)/(\lambda_G - 2\|\hat{\epsilon}\|_\infty) \leq 3$, we obtain from Lemma A.4

$$\|\hat{\mathbf{w}}_G\|_1 \leq 3\|\Delta\hat{\mathbf{w}} - \hat{\mathbf{w}}_G\|_1.$$

Therefore $\|\Delta\hat{\mathbf{w}} - \Delta\hat{\mathbf{w}}_I\|_\infty \leq \|\Delta\hat{\mathbf{w}}_G\|_1/\ell \leq 3\|\Delta\hat{\mathbf{w}} - \hat{\mathbf{w}}_G\|_1/\ell$, which implies that

$$\begin{aligned} \|\Delta\hat{\mathbf{w}} - \Delta\hat{\mathbf{w}}_I\|_2 &\leq (\|\Delta\hat{\mathbf{w}} - \Delta\hat{\mathbf{w}}_I\|_1 \|\Delta\hat{\mathbf{w}} - \Delta\hat{\mathbf{w}}_I\|_\infty)^{1/2} \\ &\leq 3^{1/2} \|\Delta\hat{\mathbf{w}} - \hat{\mathbf{w}}_G\|_1 \ell^{-1/2} \leq (3k/\ell)^{1/2} \|\Delta\hat{\mathbf{w}}_I\|_2. \end{aligned}$$

By rearranging this inequality, we obtain the desired bound. \blacksquare

Lemma A.6 *Let the conditions of Lemma A.4 hold, and let $k = |G^c|$. If $t = 1 - \pi(k + \ell, \ell)k^{1/2}\ell^{-1} > 0$, and $(\lambda_0 + 2\|\hat{\epsilon}\|_\infty)/(\lambda_G - 2\|\hat{\epsilon}\|_\infty) \leq (4 - t)/(4 - 3t)$, then*

$$\|\Delta\hat{\mathbf{w}}\|_2 \leq \frac{1 + (3k/\ell)^{0.5}}{t\rho_-(k + \ell)} \left[2\|\hat{\epsilon}_{G^c}\|_2 + \left(\sum_{j \in F} (\lambda_j^{(\ell-1)})^2 \right)^{1/2} \right].$$

Proof Let J be the indices of the largest ℓ coefficients (in absolute value) of $\hat{\mathbf{w}}_G$, and $I = G^c \cup J$. The conditions of the lemma imply that

$$\begin{aligned} \max(0, \Delta\hat{\mathbf{w}}_I^\top A \Delta\hat{\mathbf{w}}) &\geq \rho_-(k + \ell) [\|\Delta\hat{\mathbf{w}}_I\|_2 - \pi(k + \ell, \ell) \|\hat{\mathbf{w}}_G\|_1/\ell] \|\Delta\hat{\mathbf{w}}_I\|_2 \\ &\geq \rho_-(k + \ell) [1 - (1 - t)(4 - t)(4 - 3t)^{-1}] \|\Delta\hat{\mathbf{w}}_I\|_2^2 \\ &\geq 0.5t\rho_-(k + \ell) \|\Delta\hat{\mathbf{w}}_I\|_2^2. \end{aligned}$$

In the above derivation, the first inequality is due to Lemma A.2, and the second inequality is due to Lemma A.4. The last inequality follows from $1 - (1 - t)(4 - t)(4 - 3t)^{-1} \geq 0.5t$.

If $\Delta\hat{\mathbf{w}}_I^\top A \Delta\hat{\mathbf{w}} \leq 0$, then the above inequality, together with Lemma A.5, imply the lemma. Therefore in the following, we can assume that

$$\Delta\hat{\mathbf{w}}_I^\top A \Delta\hat{\mathbf{w}} \geq 0.5t\rho_-(k + \ell) \|\Delta\hat{\mathbf{w}}_I\|_2^2.$$

Moreover, let $\lambda_j = \lambda_j^{(\ell-1)}$. We obtain from (12) with $\mathbf{v} = \Delta\hat{\mathbf{w}}_I$ the following:

$$\begin{aligned} 2\Delta\hat{\mathbf{w}}_I^\top \hat{A} \Delta\hat{\mathbf{w}} &\leq -2\Delta\hat{\mathbf{w}}_I^\top \hat{\epsilon} - \sum_{j \in I} \lambda_j \Delta\hat{\mathbf{w}}_j \text{sgn}(\hat{\mathbf{w}}_j) \\ &\leq 2\|\Delta\hat{\mathbf{w}}_I\|_2 \|\hat{\epsilon}_{G^c}\|_2 + 2\|\hat{\epsilon}_G\|_\infty \sum_{j \in G} |\Delta\hat{\mathbf{w}}_j| + \sum_{j \in F} \lambda_j |\Delta\hat{\mathbf{w}}_j| - \sum_{j \in G} \lambda_j |\Delta\hat{\mathbf{w}}_j| \\ &\leq 2\|\Delta\hat{\mathbf{w}}_I\|_2 \|\hat{\epsilon}_{G^c}\|_2 + \left(\sum_{j \in F} \lambda_j^2 \right)^{1/2} \|\Delta\hat{\mathbf{w}}_I\|_2. \end{aligned}$$

Now by combining the above two estimates, we obtain

$$\|\Delta\hat{\mathbf{w}}_I\|_2 \leq \frac{1}{t\rho_-(k + \ell)} \left[2\|\hat{\epsilon}_{G^c}\|_2 + \left(\sum_{j \in F} \lambda_j^2 \right)^{1/2} \right].$$

The desired bound now follows from Lemma A.5. \blacksquare

Lemma A.7 Consider $g(\cdot)$ that satisfies the conditions of Theorem 5.1. Let $\lambda_j = \lambda g'(|\tilde{\mathbf{w}}_j|)$ for some $\tilde{\mathbf{w}} \in R^d$, then

$$\left(\sum_{j \in F} \lambda_j^2 \right)^{1/2} \leq \lambda \left(\sum_{j \in F} g'(|\bar{\mathbf{w}}_j| - \theta)^2 \right)^{1/2} + \lambda \theta^{-1} \left(\sum_{j \in F} |\bar{\mathbf{w}}_j - \tilde{\mathbf{w}}_j|^2 \right)^{1/2}.$$

Proof By assumption, if $|\bar{\mathbf{w}}_j - \tilde{\mathbf{w}}_j| \geq \theta$, then

$$g'(|\tilde{\mathbf{w}}_j|) \leq 1 \leq \theta^{-1} |\bar{\mathbf{w}}_j - \tilde{\mathbf{w}}_j|;$$

otherwise, $g'(|\tilde{\mathbf{w}}_j|) \leq g'(|\bar{\mathbf{w}}_j| - \theta)$. It follows that the following inequality always holds:

$$g'(|\tilde{\mathbf{w}}_j|) \leq g'(|\bar{\mathbf{w}}_j| - \theta) + \theta^{-1} |\bar{\mathbf{w}}_j - \tilde{\mathbf{w}}_j|.$$

The desired bound is a direct consequence of the above result and the 2-norm triangle inequality $(\sum_j (x_j + \Delta x_j)^2)^{1/2} \leq (\sum_j x_j^2)^{1/2} + (\sum_j \Delta x_j^2)^{1/2}$. \blacksquare

Lemma A.8 Under the conditions of Theorem 5.1, we have for all $\ell \geq 2\bar{k}$:

$$\|\hat{\mathbf{w}}^{(\ell)} - \bar{\mathbf{w}}\|_2 \leq \frac{7}{\rho_-(2\bar{k} + \ell)} \sqrt{|F|} \lambda.$$

Proof Let $t = 0.5$, then using Lemma A.3, the condition of the theorem implies that

$$\frac{\lambda + 2\|\hat{\epsilon}\|_\infty}{\lambda g'(\theta) - 2\|\hat{\epsilon}\|_\infty} \leq \frac{4 - t}{4 - 3t}.$$

Moreover, Lemma A.1 implies that the condition

$$t = 0.5 \leq 1 - \pi(2\bar{k} + \ell, \ell)(2\bar{k})^{0.5}/\ell$$

is also satisfied.

For each ℓ , we prove by induction for $\ell = 1, 2, \dots$. Let $G = \{j \notin F : \lambda_j^{(\ell-1)} \geq \lambda g(\theta)\}$. Then when $\ell = 1$, $G = F^c$. Otherwise, by induction hypothesis, we have

$$\sqrt{|G^c - F|} \leq \frac{\|\hat{\mathbf{w}}^{(\ell-1)} - \bar{\mathbf{w}}\|_2}{\theta} \leq \frac{7\lambda}{\rho_-(2\bar{k} + \ell)\theta} \sqrt{|F|} \leq \sqrt{|F|}.$$

Therefore, we have $|G^c| \leq 2\bar{k}$. We obtain from Lemma A.6 that

$$\|\hat{\mathbf{w}}^{(\ell)} - \bar{\mathbf{w}}\|_2 \leq \frac{1 + \sqrt{3}}{t\rho_-(2\bar{k} + \ell)} \left[2\sqrt{|G^c|}\|\hat{\epsilon}\|_\infty + \sqrt{|F|}\lambda \right] \leq \frac{3.2}{t\rho_-(2\bar{k} + \ell)} \sqrt{|F|}\lambda,$$

where we have used the fact that $\lambda \geq 20\|\hat{\epsilon}\|_\infty$ in the derivation of the second inequality. \blacksquare

Proof of Theorem 5.1

As in the proof of Lemma A.8, if we let $t = 0.5$, then using Lemma A.3, the condition of the theorem implies that

$$\frac{\lambda + 2\|\hat{\epsilon}\|_\infty}{\lambda g'(\theta) - 2\|\hat{\epsilon}\|_\infty} \leq \frac{4-t}{4-3t}.$$

Moreover, Lemma A.1 implies that the condition

$$t = 0.5 \leq 1 - \pi(2\bar{k} + \ell)(2\bar{k})^{0.5}/\ell$$

is also satisfied.

We prove by induction: for $\ell = 1$, the result follows from Lemma A.8. For $\ell > 1$, we let $G^c = F \cup \{j : |\hat{\mathbf{w}}_j^{(\ell-1)}| \geq \theta\}$. By Lemma A.8, we know that

$$k = |G^c| \leq \bar{k} + \frac{\|\hat{\mathbf{w}}^{(\ell-1)} - \bar{\mathbf{w}}\|_2}{\theta} \leq 2\bar{k}.$$

Let $u = \sqrt{\rho_+(\bar{k})}\sigma[\sqrt{7.4\bar{k}/n} + \sqrt{2.7\ln(2/\eta)/n}]$. We know from Lemma A.3, and $\lambda \geq 20\|\hat{\epsilon}\|_\infty$ that with probability $1 - 2\eta$,

$$\begin{aligned} \|\hat{\epsilon}_{G^c}\|_2 &\leq \|\hat{\epsilon}_F\|_2 + \sqrt{|G^c - F|}\|\hat{\epsilon}\|_\infty \\ &\leq u + \sqrt{|G^c - F|}\lambda/20 \\ &\leq u + \lambda(20\theta)^{-1}\|\hat{\mathbf{w}}^{(\ell-1)} - \bar{\mathbf{w}}\|_2. \end{aligned}$$

Now, using Lemma A.6 and Lemma A.7, we obtain

$$\begin{aligned} \|\Delta\hat{\mathbf{w}}^{(\ell)}\|_2 &\leq \frac{1 + \sqrt{3}}{t\rho_-(k + \ell)} \left[2\|\hat{\epsilon}_{G^c}\|_2 + \left(\sum_{j \in F} (\lambda_j^{(\ell-1)})^2 \right)^{1/2} \right] \\ &\leq \frac{1 + \sqrt{3}}{t\rho_-(k + \ell)} \left[2\|\hat{\epsilon}_{G^c}\|_2 + \lambda \left(\sum_{j \in F} g'(|\bar{\mathbf{w}}_j| - \theta)^2 \right)^{1/2} + \lambda\theta^{-1} \left(\sum_{j \in F} |\bar{\mathbf{w}}_j - \hat{\mathbf{w}}_j^{(\ell-1)}|^2 \right)^{1/2} \right] \\ &\leq \frac{1 + \sqrt{3}}{t\rho_-(k + \ell)} \left[2u + \lambda \left(\sum_{j \in F} g'(|\bar{\mathbf{w}}_j| - \theta)^2 \right)^{1/2} + 1.1\lambda\theta^{-1}\|\bar{\mathbf{w}} - \hat{\mathbf{w}}^{(\ell-1)}\|_2 \right] \\ &\leq \frac{1 + \sqrt{3}}{t\rho_-(k + \ell)} \left[2u + \lambda \left(\sum_{j \in F} g'(|\bar{\mathbf{w}}_j| - \theta)^2 \right)^{1/2} \right] + 0.67\|\bar{\mathbf{w}} - \hat{\mathbf{w}}^{(\ell-1)}\|_2. \end{aligned}$$

We solve this recursion to obtain the desired bound.