

Implicit Regularization via Hadamard Product Over-Parametrization in High-Dimensional Linear Regression

Peng Zhao¹, Yun Yang^{*2}, and Qiao-Chu He³

¹Department of Statistics, Florida State University

²Department of Statistics, University of Illinois Urbana-Champaign

³Southern University of Science and Technology

Abstract

We consider Hadamard product parametrization as a change-of-variable (over-parametrization) technique for solving least square problems in the context of linear regression. Despite the non-convexity and exponentially many saddle points induced by the change-of-variable, we show that under certain conditions, this over-parametrization leads to implicit regularization: if we directly apply gradient descent to the residual sum of squares with sufficiently small initial values, then under proper early stopping rule, the iterates converge to a nearly sparse rate-optimal solution with relatively better accuracy than explicit regularized approaches. In particular, the resulting estimator does not suffer from extra bias due to explicit penalties, and can achieve the parametric root- n rate (independent of the dimension) under proper conditions on the signal-to-noise ratio. We perform simulations to compare our methods with high dimensional linear regression with explicit regularizations. Our results illustrate advantages of using implicit regularization via gradient descent after over-parametrization in sparse vector estimation.

Keywords: High-dimensional regression; Implicit regularization; Over-parametrization; Early Stopping; Gradient Descent; Variable selection.

^{*}Y. Yang's research is supported in part by NSF DMS-1810831.

1 Introduction

In high dimensional linear regression

$$y = X\beta^* + w, \quad \text{with noise } w \sim \mathcal{N}(0, \sigma^2 I),$$

the goal is to parsimoniously predict the response $y \in \mathbb{R}^n$ as a linear combination of a large number of covariates $X = (X_1, X_2, \dots, X_p) \in \mathbb{R}^{n \times p}$, and conduct statistical inference on the linear combination coefficients $\beta^* = (\beta_1^*, \dots, \beta_p^*)^T \in \mathbb{R}^p$ (Tibshirani, 1996; Donoho, 2006). By leveraging on certain lower dimensional structure in the regression coefficient vector $\beta^* \in \mathbb{R}^p$ such as a sparsity constraint $s = \|\beta^*\|_0 \ll n$, where $\|\beta^*\|_0$ counts the number of nonzeros in β^* , the number p of covariates is allowed to be substantially larger than the sample size n . Due to the intrinsic computational hardness in dealing with the ℓ_0 metric reflecting sparsity, people instead use different metrics as surrogates, and cast the estimation problem into various convex or nonconvex optimization problems. Many approaches have been proposed for high dimensional regression by solving certain penalized optimization problem, including basis pursuit (Chen et al., 2001), the Lasso (Tibshirani, 1996), the Dantzig selector (Candes and Tao, 2007), SCAD (Fan and Li, 2001), MCP (Zhang, 2010) and so on. In this work, we focus on the recovery of $\beta^* \in \mathbb{R}^p$ without explicitly specifying a penalty.

Recent work (Hoff, 2017) shows that through a change-of-variable (over-parametrization) via Hadamard product parametrization, the Lagrangian (dual) form of the non-smooth convex optimization problem for the Lasso (1):

$$\min_{\beta} \frac{1}{2n} \|X\beta - y\|^2 + \lambda \|\beta\|_1, \quad \text{with } \|\beta\|_1 := \sum_{j=1}^p |\beta_j|, \quad (1)$$

can be reformulated as a smoothed optimization problem at a cost of introducing non-convexity. Due to the smoothness feature, simple and low-cost first-order optimization

methods such as gradient descent and coordinate descent can be applied to recover β^* . Despite the non-convexity and exponentially many stationary points induced by the change-of-variable, these first-order algorithms exhibit encouraging empirical performance (Hoff, 2017).

In this work, we consider the same Hadamard product over-parametrization $\beta = g \circ l$ as in Hoff (2017), where $g, l \in \mathbb{R}^p$ and \circ denotes the Hadamard product (element-wise product). Instead of solving the penalized optimization problem (1), we consider directly applying the gradient descent to the quadratic loss function

$$f(g, l) = \frac{1}{2n} \|X(g \circ l) - y\|^2. \quad (2)$$

In the noiseless case where $\sigma = 0$, minimizing $f(g, l)$ jointly over (g, l) is a highly non-convex optimization problem with exponentially many saddle points. To see this, notice that each non-zero pattern of β corresponds to at least one saddle point by choosing $g_j = l_j = 0$ for each j such that $\beta_j = 0$. In addition, due to the column rank deficiency of the design matrix X (for example, when $p > n$), there are infinitely many global minimizers of (2) as potential convergent points of the gradient descent. Interestingly, we show that despite these seemingly hopeless difficulties, in the noiseless case if we initialize the gradient descent arbitrarily close to $g = l = 0$, then under the prominent Restricted Isometry Property (RIP) condition (Candes, 2008) on the design matrix X , a properly tuned gradient descent converges to least ℓ_1 -norm solution within error ε in $\mathcal{O}(\log \frac{C}{\varepsilon})$ iterations, where constant C depends on the RIP constant, step size of the gradient descent, and some other characteristics of the problem. Our proofs borrow ideas from Li et al. (2018), where they prove the algorithmic convergence of matrix factorized gradient descent in the context of noiseless matrix sensing under the RIP.

In high dimensional regression, the usual regularized least square is known to suffer from the so-called saturation phenomenon (Vito et al., 2005; Yao et al., 2007), where the overall estimation error is dominated by a bias term due to the penalty. In particular, since regu-

larization is artificially introduced for restricting the “effective size” of the parameter space, the resulting estimator may be deteriorated and the estimation error cannot fall below the penalty level to adapt to a possibly faster convergence rate. For example, the estimator by solving the Lasso achieves the minimax rate of $\sqrt{s}\lambda \asymp \sqrt{s \log p/n}$. However, this worse-case rate only happens when there exist weak signals, meaning that some nonzero β_j^* ’s have a borderline magnitude of order $\sqrt{s \log p/n}$. In fact, if all signals are sufficiently strong, or significantly larger this borderline magnitude, then the faster dimension-independent parametric rate $\sqrt{s/n}$ is attainable. For regularized approaches such the Lasso, one possible way to remove the penalty-induced bias term (whose order is λ) is to refit the model with the selected variables. However, this two stage procedure requires stringent assumptions on the minimal signal strength to guarantee variable selection consistency for the first stage, and will suffer from weak signals. Interestingly, we show that by combining the Hadamard product over-parametrization with early stopping, a widely used regularization technique in boosting (Zhang and Yu, 2005) and nonparametric regression (Raskutti et al., 2014), our method can overcome the saturation issue and lead to more accurate estimation. More precisely, in the presence of random noise w in the linear model, the solution path of minimizing the quadratic loss function (2) as we increase the gradient descent iteration still tends to converge to the least ℓ_1 -norm solution that will overfit the data. Fortunately, by terminating the gradient descent updating procedure earlier within a proper number of iterations, we may find a solution that optimally balances between the model complexity (reflected by the increasing ℓ_1 -norm of the iterate) and goodness fit of the model, akin the bias-variance trade-off. In particular, we show that the estimator can adapt to an optimal convergence rate of $\sqrt{s/n}$ when all signals are relatively strong. Generally, when both strong signals and weak signals exist, our estimator attains the rate $\sqrt{s_1/n} + \sqrt{s_2 \log p/n}$ (with s_1, s_2 denoting the number of strong signals and weak signals, respectively).

Our result also complements the recent surge of literature on over-parametrization for implicit regularization of the first-order iterative method for non-convex optimization in

machine learning. [Gunasekar et al. \(2017\)](#) introduce the phenomenon of implicit regularization in matrix factorization, where they empirically observe the convergence of gradient methods in matrix factorization problem to the minimal nuclear norm solution as the initial value tends to zero. However, they only provide some heuristic illustration under some hard-to-check assumptions such as the continuity of the solution relative to the change in the initialization. Later, [Li et al. \(2018\)](#) rigorously prove the implicit regularization in matrix sensing problem under a matrix RIP condition. Some other very recent works such [Gunasekar et al. \(2018\)](#) and [Soudry et al. \(2018\)](#) study implicit regularizations in mirror descent and in classification problems. Note that all above implicit regularization literatures only consider data that are either noiseless (regression) or perfectly separable (classification). To our best knowledge, we are the first to rigorously study and utilize implicit regularization in high dimensional linear regression where responses are noisy.

In a nutshell, we show that through a simple change-of-variable, the non-smooth ℓ_1 -penalized optimization problem (1) can be transformed to an unconstrained smooth quadratic loss minimization one; moreover, a simple gradient descent initialized near zero efficiently solves this non-convex optimization problem with provable guarantees. Furthermore, our method enjoy several advantages over existing procedures for high dimensional linear regression under sparsity constraints. First, our method is computationally efficient — its time complexity is $O(np)$, which is linear in both n and p . Second, despite the non-convexity nature, our method has a natural initialization that provably leads the optimal solution. In comparison, penalized M -estimators based on non-convex penalties such as SCAD and MCP require stringent conditions on their initializations: to obtain good estimators, they require good initial values that are sufficiently closed to the truth (theoretically) or satisfy some restricted strong convexity conditions ([Zhao et al., 2018](#)), otherwise their optimization algorithms will suffer from bad local minima with bad generalization errors. In contrast, our algorithm only requires the initialization to be closed to zero. Moreover, unlike penalized approaches such as SCAD and MCP, where both parameters for the noise

level and the concavity of the penalty need to be tuned, our method only need to tune the number of iterations.

To conclude, our main contributions with respect to the relative literatures are as follows:

1. We propose an estimator by combining early stopping with implicit regularization to overcome the saturation issues in high dimensional regression with explicit regularizations;
2. Unlike recent implicit regularization literatures that exclusively focus on noiseless data, we are the first to rigorously study the effect of implicit regularization for noisy data;
3. From computational perspective, we transform the non-smooth optimization problem to an unconstrained smooth quadratic loss minimization problem for which standard optimization tools can be applied.

2 Background and Our Method

To begin with, we formally introduce the setup and notations used throughout the paper. After that, we introduce the intuition for our new implicit regularized algorithm for high dimensional linear regression via Hadamard product parameterization.

2.1 Setup and notations

Recall that β^* is the unknown s -sparse signal in \mathbb{R}^p to be recovered. Let $S \subset \{1, \dots, p\}$ denote the index set that corresponds to the nonzero components of β^* , and the size $|S|$ of S is then s . For two vectors $g, l \in \mathbb{R}^p$, we call $\beta = g \circ l \in \mathbb{R}^p$ as their Hadamard product, whose components are $\beta_j = g_j l_j$ for $j = 1, \dots, p$. For two vectors $a, b \in \mathbb{R}^p$, we use the notation $a \geq b$ to indicate element-wise “great than or equal to”. When there

is no ambiguity, we use $\beta^2 = \beta \circ \beta$ to denote the self-Hadamard product of β . For a function $f : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$, $(g, l) \mapsto f(g, l)$, we use $\nabla_g f$ and $\nabla_l f$ to denote its partial derivative relative to g and l , respectively. For any index set $J \subset \{1, \dots, p\}$ and vector $a \in \mathbb{R}^p$, we use $a_J = (a_j : j \in J)$ to denote the sub-vector of a formed by concatenating the components indexed by J . Let $\mathbf{1} \in \mathbb{R}^p$ denote the vector with all entries as 1, and I as the identity matrix in \mathbb{R}^p . Let I_J be the diagonal matrix with one on the j th diagonal for $j \in J$ and 0 elsewhere. For a vector $a \in \mathbb{R}^p$, we use $\|a\|$ to denote its vector- ℓ_2 -norm, and $\|a\|_\infty = \max_j |a_j|$ its ℓ_∞ -norm. Let $\text{Unif}(a, b)$ to denote the uniform distribution over interval (a, b) . For a symmetric matrix A , let $\lambda_{\min}(A)$ denote its smallest eigenvalue. For two sequences $\{a_n\}$ and $\{b_n\}$, we use the notation $a_n \lesssim b_n$ or $a_n \gtrsim b_n$ to mean there exist some constant c and C independent of n such that $a_n \leq Cb_n$ or $a_n \geq cb_n$ for all $n < 0$, respectively, and $a_n \asymp b_n$ to mean $a_n \lesssim b_n$ and $b_n \lesssim a_n$.

2.2 Gradient descent with Hadamard product parametrization

As we mentioned in the introduction, we consider augmenting the p -dimensional vector β into two p -dimensional vectors g, l through $\beta = g \circ l$. Instead of solving the Lasso problem (1) with β replaced with $g \circ l$, we consider directly applying gradient descent to the quadratic loss function $f(g, l) = (2n)^{-1} \|X(g \circ l) - y\|^2$. In particular, we apply the updating formula $g_{t+1} = g_t - \eta \nabla_g f(g_t, l_t)$, $l_{t+1} = l_t - \eta \nabla_l f(g_t, l_t)$, with random initial values g_0 and l_0 chosen close enough to 0 (notice that $(0, 0)$ is a saddle point of the objective function, so we need to apply a small perturbation α on the initial values). This leads to the following algorithm:

Algorithm 1: Gradient Descent for linear regression

Data: Design matrix $X \in \mathbb{R}^{n \times p}$, measurement vector $y \in \mathbb{R}^n$, initialization magnitude α , step size η , and stopping threshold ϵ ;

Initialize variables $[g_0]_j \stackrel{iid}{\sim} \text{Unif}(-\alpha, \alpha)$, $[l_0]_j \stackrel{iid}{\sim} \text{Unif}(-\alpha, \alpha)$ for $j = 1, \dots, p$, and iteration number $t = 0$;

while $\|X(g_t \circ l_t) - y\|/\sqrt{n} > \epsilon$ **do**

$g_{t+1} = g_t - \eta l_t \circ [n^{-1} X^T (X(g_t \circ l_t) - y)]$;
 $l_{t+1} = l_t - \eta g_t \circ [n^{-1} X^T (X(g_t \circ l_t) - y)]$;
 $t = t + 1$;

end

Result: Output the final estimate $\hat{\beta} = g_t \circ l_t$;

Algorithm 1 is the standard gradient descent algorithm, and the iterates (g_{t+1}, l_{t+1}) tend to converge to a stationary point (g_∞, l_∞) of $f(g, l)$ that satisfies the first order optimality condition $\nabla f_g(g_\infty, l_\infty) = 0$ and $\nabla f_l(g_\infty, l_\infty) = 0$. However, stationary points of $f(g, l)$ can be local minimum, local maximum, or saddle points (when the Hessian matrix $\nabla_{g,l}^2 f(g, l)$ contains both positive and negative eigenvalues). The following result provides the optimization landscape of $f(g, l)$, showing that $f(g, l)$ does not have local maximum, all its local minimums are global minimum, and all saddle points are strict. The strict saddle points are saddle points with negative smallest eigenvalues for Hessian matrix.

Lemma 2.1. $f(g, l) = (2n)^{-1} \|X(g \circ l) - y\|^2$ does not have local maximum, and all its local minimums are global minimum. In particular, (\bar{g}, \bar{l}) is a global minimum of $f(g, l)$ if and only if

$$X^T (X(\bar{g} \circ \bar{l}) - y) = 0.$$

In addition, any saddle point (g^\dagger, l^\dagger) of $f(g, l)$ is a strict saddle, that is, $\lambda_{\min}(\nabla_{g,l}^2 f(g^\dagger, l^\dagger)) < 0$.

According to the first order condition associated with $f(g, l)$

$$g \circ \left[X^T (X(g \circ l) - y) \right] = l \circ \left[X^T (X(g \circ l) - y) \right] = 0,$$

there could be exponentially many (at least $2^p - 1$) saddle points as a solution to this equation, for example, for those (g, l) satisfying

$$g_A = l_A = 0 \in \mathbb{R}^{|A|}, \quad \text{and} \quad \left[X^T (X(g \circ l) - y) \right]_{A^c} = 0 \in \mathbb{R}^{p-|A|},$$

for any non-empty subset A of $\{1, \dots, p\}$. Consequently, the gradient descent algorithm may converge to any of these bad saddle points. To see this, if we initialize (g, l) in a way such that the components in the index set A are zero, or $[g_0]_A = [l_0]_A = 0$, then these components will remain zero forever in the gradient iterations, or $[g_t]_A = [l_t]_A = 0$ for all $t > 0$. Fortunately, the following result implies that as long as we use a random initialization for (g, l) with continuous pdf over \mathbb{R}^{2p} as in Algorithm 1, then the gradient descent almost surely converges to a global minimum.

Lemma 2.2. *Suppose the step size η is sufficiently small. Then with probability one, Algorithm 1 converges to a global minimum of $f(g, l)$.*

In the low-dimensional regime where the design matrix X has full column rank, the solution $\bar{\beta}$ to the normal equation $X^T(X\beta - y) = 0$ is unique, which is also the least squares estimator. Under this scenario, Lemma 2.1 and Lemma 2.2 together certify that Algorithm 1 will converge to this optimal least squares estimator. However, in the high-dimensional regime which is the main focus in the paper, the normal equation $X^T(X\beta - y) = 0$ have infinitely many solutions, and it is not clear which solution the algorithm tends to converge to. For example, if we consider instead applying the gradient descent to the original parameter β in the objective function $(2n)^{-1} \|X\beta - y\|^2$ with initialization $\beta_0 = 0$, then the iterates will converge to the minimal ℓ_2 -norm solution of the normal equation (see below for details). Interestingly, as we will illustrate in the following, under the Hadamard

parametrization the gradient descent Algorithm 1 now tends to converge to the minimal ℓ_1 -norm solution under certain conditions for initialization, thereby inducing sparsity and naturally facilitating variable selection.

2.3 Gradient descent converges to sparse solution

In this subsection, we provide two different perspectives for understanding the following informal statement on the behavior of simple gradient descent for the loss function $f(g, l)$ defined in (2) under the Hadamard product parameterization $\beta = g \circ l$. For simplicity, we assume that the response y in the linear model is perfect, that is, the noise variance $\sigma^2 = 0$, throughout this subsection. Then in the next subsection, we turn to general noisy observations, and propose methods that lead to optimal estimation when the true regression coefficient β^* is sparse.

Informal Statement: *If we initialize the algorithm to be arbitrarily close to $g = l = 0$, then under suitable conditions on design X , a simple gradient descent converges to a solution of basis pursuit problem:*

$$\min_{\beta \in \mathbb{R}^p} \|\beta\|_1 \quad \text{subject to} \quad X\beta = y. \quad (3)$$

Our first perspective is based on the ℓ_2 -norm implicit regularization in linear system, and the second is based on analyzing the gradient dynamical system as the limit of the gradient descent algorithm as the step size $\eta \rightarrow 0_+$. However, both perspectives in this section are heuristic, and formal statements and proofs (based on a different strategy) will be provided in Section 3.

ℓ_2 -norm implicit regularization perspective: Consider the under-determined system $X\beta = y$, where $X \in \mathbb{R}^{n \times p}$ has full row rank. Our first intuition comes from the fact that

a zero-initialized gradient descent algorithm over $\beta \in \mathbb{R}^p$ for solving

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{2n} \|X\beta - y\|^2 := g(\beta)$$

finds a minimal ℓ_2 -norm solution to $X\beta = y$.

In fact, we know that any solution to the linear system $X\beta = y$ takes the form of

$$\beta = X^+y + [I - X^+X]w, \quad \text{over all } w \in \mathbb{R}^p,$$

where $X^+ = \lim_{\lambda \rightarrow 0^+} (X^T X + \lambda I)^{-1} X^T$ is the Moore-Penrose inverse of X . Since $X(I - X^+X) = 0$, we have

$$\|\beta\|^2 = \|X^+y\|^2 + \|[I - X^+X]w\|^2 \geq \|X^+y\|^2,$$

implying that X^+y is the unique solution of $X\beta = y$ in the column space of X^T , which is also the minimal ℓ_2 -norm solution. Now since the gradient updating formula for β , $\beta_{t+1} = \beta_t - \eta X^T(X\beta_t - y)/n$, implies that the difference $\beta_t - \beta_0$ always lies in the column span of X^T . Let $\beta_\infty := \lim_{t \rightarrow \infty} \beta_t$ be the limiting point of the gradient algorithm. Then β_∞ must be a solution to $X\beta = y$. On the other hand, when β_0 is initialized at zero, β_∞ should also belong to the column span of X^T . These two properties combined imply that β_∞ must be the minimal ℓ_2 -norm solution X^+y .

In high dimensional linear regression with perfect observations, a popular class of penalization methods attempts find the minimal ℓ_1 -norm solution to $X\beta = y$. Under the Hadamard product parameterization $\beta = g \circ l$, the fact that gradient descent tends to find the minimal ℓ_2 -norm solution suggests (this is not rigorous) that the gradient descent algorithm for jointly minimizing $f(g, l)$ over (g, l) tends to converge to a solution to $X(g \circ l) = y$ with a minimal ℓ_2 -norm $\sqrt{\|g\|^2 + \|l\|^2}$. However, a minimal ℓ_2 -norm solution to $X(g \circ l) = y$ must satisfy $|g_j| = |l_j|$ for each $j = 1, \dots, p$ (otherwise we can always construct another so-

lution with strictly smaller ℓ_2 -norm), which implies $\sqrt{\|g\|^2 + \|l\|^2} = \sqrt{2} \|g \circ l\|_1 = \sqrt{2} \|\beta\|_1$. As a consequence, $\beta_\infty = g_\infty \circ l_\infty$ should be the minimal ℓ_1 -norm solution to $X\beta = y$.

Another way to understand the difference in the evolutions of gradient descents for $f(g, l)$ and $h(\beta)$ is by noticing that the gradient $\nabla_{g_j} f(g, l) = l_j \cdot \nabla_{\beta_j} h(\beta) \Big|_{\beta=g \circ l}$ in the new parametrization, for each $j = 1, \dots, p$, has an extra multiplicative factor of l_j than the gradient $\nabla_{\beta_j} h(\beta)$ in the usual least squares of minimizing $g(\beta)$. It is precisely this extra multiplicative factor l_j that helps select important signals (nonzero regression coefficients) and prevent unimportant signals (zero regression coefficients) to grow too fast at the early stage of the evolution when both g and l are close to zero. Precisely, as we will show in our theory part (Section 3), under suitable conditions on the model, all unimportant signals remain to stay in a $\mathcal{O}(p^{-1})$ neighbourhood of zero, while important ones tend to grow exponentially fast.

Gradient dynamical system perspective: Our second perspective comes from considering the limiting gradient dynamical system of the problem (i.e. gradient descent with an infinitesimally small step size), which is motivated by the interpretation for matrix factorization problems in [Gunasekar et al. \(2017\)](#) and [Gunasekar et al. \(2018\)](#). In particular, the behavior of this limiting dynamical system is captured by the ordinary differential equations

$$\begin{cases} \dot{g}(t) = -[X^T r(t)] \circ l(t), \\ \dot{l}(t) = -[X^T r(t)] \circ g(t), \end{cases} \quad \text{with initialization} \quad \begin{cases} g(0) = \alpha \mathbf{1}, \\ l(0) = 0, \end{cases} \quad (4)$$

where $r(t) = n^{-1}[X(g(t) \circ l(t)) - y] \in \mathbb{R}^p$, and for simplicity we fixed the initialization. To emphasize the dependence of the solution on α , we instead write $g(t)$, $l(t)$, $r(t)$ as $g(t, \alpha)$, $l(t, \alpha)$, $r(t, \alpha)$. For illustration purposes, we assume that the limiting point of this system is continuous and bounded as the initialization value $\alpha \rightarrow 0_+$, that is, both limits $g_\infty = \lim_{t \rightarrow \infty, \alpha \rightarrow 0_+} g(t, \alpha)$ and $l_\infty = \lim_{t \rightarrow \infty, \alpha \rightarrow 0_+} l(t, \alpha)$ exist in \mathbb{R}^p and are finite.

Let $s(t, \alpha) = \int_0^t r(\tau, \alpha) d\tau \in \mathbb{R}^p$, then simple calculation leads to the relation

$$\begin{bmatrix} g_j(t, \alpha) + l_j(t, \alpha) \\ g_j(t, \alpha) - l_j(t, \alpha) \end{bmatrix} = \alpha \begin{bmatrix} \exp(-X_j^T s(t, \alpha)) \\ \exp(X_j^T s(t, \alpha)) \end{bmatrix}, \quad \text{for each } j = 1, \dots, p.$$

Under the aforementioned assumption on the existence of limits as $t \rightarrow \infty$ and $\alpha \rightarrow 0_+$, the preceding display implies one of the following three situations for each j :

Case 1: $g_{j,\infty} = l_{j,\infty} \neq 0$, and $\lim_{t \rightarrow \infty, \alpha \rightarrow 0_+} X_j^T s(t, \alpha) / \log(1/\alpha) = 1$.

Case 2: $g_{j,\infty} = -l_{j,\infty} \neq 0$, and $\lim_{t \rightarrow \infty, \alpha \rightarrow 0_+} X_j^T s(t, \alpha) / \log(1/\alpha) = -1$.

Case 3: $g_{j,\infty} = l_{j,\infty} = 0$, and $\lim_{t \rightarrow \infty, \alpha \rightarrow 0_+} X_j^T s(t, \alpha) / \log(1/\alpha) = \gamma_j \in [-1, 1]$.

Denote s_∞ as the limit $\lim_{t \rightarrow \infty, \alpha \rightarrow 0_+} s(t, \alpha) / \log(1/\alpha)$. Recall $\beta_\infty = g_\infty \circ l_\infty$, and the previous three cases can be unified into

$$X_j^T s_\infty = \begin{cases} \text{sign}(\beta_{j,\infty}), & \text{if } \beta_{j,\infty} \neq 0, \\ \gamma_j \in [-1, 1], & \text{if } \beta_{j,\infty} = 0, \end{cases} \quad \text{for each } j = 1, \dots, p.$$

This identity together with the limiting point condition $X\beta_\infty = y$ coincides with the KKT condition for the basis pursuit problem (3).

Again, this argument is based on the hard-to-check solution continuity assumption. In the next section, we provide a formal proof of the result without making this assumption, albeit under a somewhat strong RIP condition on X . We conjecture this gradient descent implicit regularization property to be generally true for a much wider class of design matrices (see our simulation section for numerical evidences).

2.4 Gradient descent with early stopping

In this subsection, we consider the general case where the response y contains noise, or $\sigma^2 \neq 0$. In particular, we propose the use of early stopping, a widely used implicit regu-

larization technique (Zhang and Yu, 2005; Raskutti et al., 2014), to the gradient descent Algorithm 1. As the name suggests, we will use certain criterion to decide whether to terminate the gradient descent updating to prevent overfitting of the data. Algorithm 2 below summarizes a particular stopping criterion widely-used in the machine learning community via validation. In principal, we can also treat the number of iterations as a hyperparameter and repeat this procedure multiple times in same spirits as data splitting and cross validation.

Algorithm 2: Gradient Descent for Linear Regression with Validation Data

Data: Training design matrix $X \in \mathbb{R}^{n \times p}$, measurement vector $y \in \mathbb{R}^n$, validation data X' , y' , initialization magnitude α , step size η , and maximal number of iterations T_{max} ;

Initialize variables $[g_0]_j \stackrel{iid}{\sim} \text{Unif}(-\alpha, \alpha)$, $[l_0]_j \stackrel{iid}{\sim} \text{Unif}(-\alpha, \alpha)$ for $j = 1, \dots, p$, and iteration number $t = 0$;

while $t < T_{max}$ **do**

$$\left| \begin{array}{l} g_{t+1} = g_t - \eta \, l_t \circ \left[n^{-1} X^T (X(g_t \circ l_t) - y) \right]; \\ l_{t+1} = l_t - \eta \, g_t \circ \left[n^{-1} X^T (X(g_t \circ l_t) - y) \right]; \\ t = t + 1; \end{array} \right.$$

end

Result: Choose the first \tilde{t} such that $\|X'(g_{\tilde{t}} \circ l_{\tilde{t}}) - y'\| > \|X'(g_{\tilde{t}+1} \circ l_{\tilde{t}+1}) - y'\|$ or $\|X'(g_{\tilde{t}} \circ l_{\tilde{t}}) - y'\|$ is minimized over all iterations, then output the final estimate $\hat{\beta} = g_{\tilde{t}} \circ l_{\tilde{t}}$.

Recall that in the introduction, we discussed about the saturation issue suffered by explicit penalized methods such as the Lasso. Now we turn to our method and illustrate that it is unaffected, or at least less suffered from the saturation issue. In the next section, we will provide rigorous theory showing that our method can achieve a faster $\sqrt{s/n}$ rate of convergence then the typical $\sqrt{s \log p/n}$ rate when all nonzero signals are relatively large.

Due to the connection of our method with the basis pursuit problem (3), one may naturally think that our method in the noisy case should be equivalent to a basis pursuit

denoising problem:

$$\min \|\beta\|_1 \quad \text{subject to} \quad \|X\beta - y\|_2 \leq \epsilon, \quad (5)$$

with some error tolerance level ϵ depending on the stopping criterion, and therefore is equivalent to the Lasso. Surprisingly, a simulation example below shows that the iterate path of the gradient descent Algorithm 1 contains estimates with much smaller error than the Lasso. Precisely, we adopt the simulation setting S2 in section 4.2. As comparisons, we also report the Lasso solution path (as a function of the regularization parameter λ) solved by ISTA and FISTA (Beck and Teboulle, 2009). For our gradient descent algorithm, we set $\alpha = 10^{-5}$ in the random initialization. From figure 1, when the iteration number is around 1000, even though the prediction error in panel (c) of our algorithm and the Lasso (with an optimally tuned λ , see panel (b) for the entire Lasso solution path), the estimation error in panel (a) of our method is significantly lower than that of the Lasso, illustrating the occurrence of the saturation phenomenon of the Lasso. Moreover, the stabilized region (iterations 200–1000) of our method GD in panel (a) is relatively wide, and therefore the performance tends to be robust to the stopping criterion.

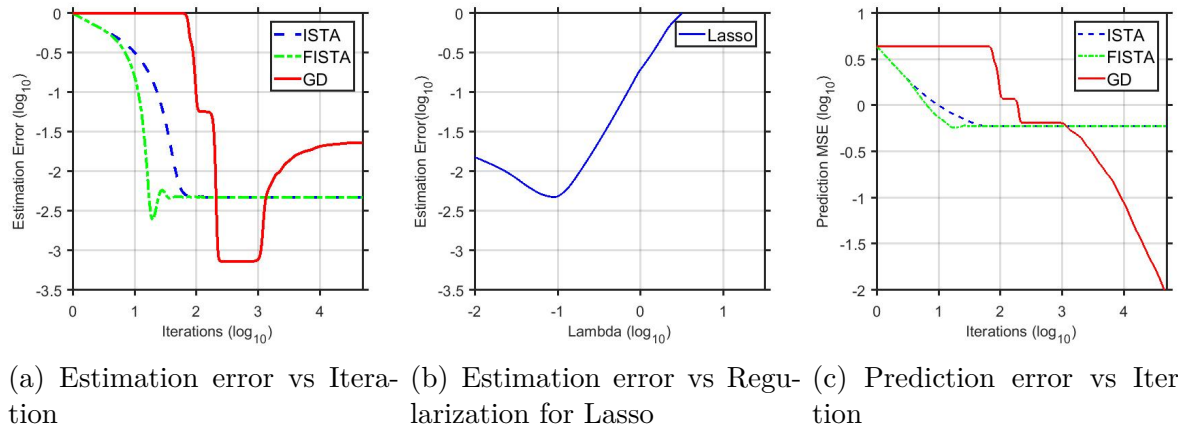


Figure 1: Panel (a) is a log-log plot of standardized estimation error $\|\hat{\beta} - \beta^*\|_2^2 / \|\beta^*\|_2^2$ versus iteration number t . Panel (b) is a log-log plot of standardized estimation error versus regularization parameter λ for Lasso. Panel (c) is a log-log plot of mean prediction error $\sqrt{\|\hat{y} - y\|_2^2 / n}$ versus iteration number t .

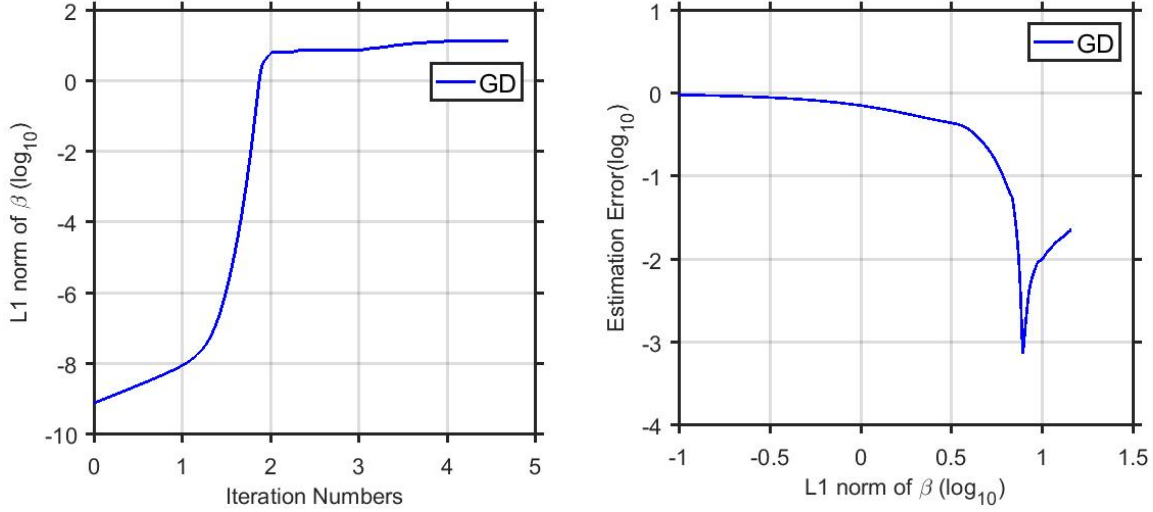
Next, let us briefly illustrate why implicit regularization with early stopping works,

while explicit regularized methods may fail. We know that early stopping, serving as algorithmic regularization, is based on the intuition that as the iteration number grows, the bias keeps decreasing while the variance increasing. Consequently, the iteration number T , acting as an implicit regularization parameter, aims to optimally balance between the bias and the variance, akin to the bias-variance trade-off. In our parametrization, the iteration number T controls the ℓ_1 norm of the solution, reflecting the model space complexity. To see this, we plot the ℓ_1 norm versus the iteration number, and also the estimation errors versus the ℓ_1 norm, all in logarithmic scales, in figure 2. As we expected, as the number of iterations increases, the ℓ_1 norm of the iterate also increases. When the logarithm of the iteration number is within $(2.3, 3)$, the ℓ_1 norm of the estimated coefficients tends to be stabilized at the ℓ_1 norm of the true β^* as 0.9, corresponding to the most accurate estimation region in panel (a) of figure 1. In contrast, as we can see from panel (b) of figure 2, the estimation error is very sensitive in the regularization parameter domain — the region corresponds to smallest estimation accuracy is very narrow, and a small change in the ℓ_1 norm in the solution leads to a drastic deterioration in the estimation accuracy. This numerical example provides evidences of why explicit regularized approaches may suffer from large bias and low accuracy.

Finally, we discuss several commonly used early stopping rules by treating the iteration number as a tuning parameter.

Hold-out or cross validation: The simplest method is to use hold-out data as validation: for example, we can split the training data into half-half, and then run gradient descent on the first half D_1 of the data while calculate the prediction error $R(t) = \sum_{i \in D_2} (X^{(i)T}(g_t \circ l_t) - y_i)^2$ for all $t \leq T_{max}$ on the second half D_2 , then the final iteration number is decided by (cf. Algorithm 2):

$$\tilde{t} := \arg \min \{t \leq T_{max} \mid R(t+1) > R(t)\} \quad \text{or} \quad (6)$$



(a) ℓ_1 norm vs Iteration

(b) Prediction error vs Iteration Comparison

Figure 2: Panel (a) is a log-log plot of ℓ_1 norm of the estimated coefficients versus iteration number t . Panel (b) is a log-log plot of standardized estimation error versus ℓ_1 norm of the estimated coefficients.

$$\tilde{t} := \arg \min \{t \leq T_{max} \mid R(t) = \min_{\tau \leq T_{max}} R(\tau)\}. \quad (7)$$

To make use of the whole dataset, we can perform cross validation: first split data into K folds, then apply gradient descent on all possible combinations of $K - 1$ folds without replacement and evaluate at the corresponding rest 1 folds. The final risk $R(t)$ can be the sum of all the evaluations on each fold, and the criterion (6) or (7) can be used to obtain the iteration number. Finally we can apply the same iteration number obtained from cross validation to approximate the optimal one for the entire training data.

Stein's unbiased risk estimate (SURE): Stein (1981) suggested the use of degrees of freedom as surrogate to the true prediction error given the standard derivation σ . Under our settings, ignoring the second order term of step size η , the updating of the prediction error (up to rescaling) $r_t = n^{-1} [X(g_t \circ l_t) - y] \in \mathbb{R}^n$ through gradient descent can be approximated by (by ignoring second order terms of order η^2):

$$r_{t+1} \approx [I - 2\eta n^{-1} X \text{diag}(|g_t \circ l_t|) X^T] r_t, \quad (8)$$

where for a vector u , $\text{diag}(u)$ denotes the diagonal matrix with components of u in the its diagonals. Define $S_t = \Pi_{i=1}^{t-1}(I - 2\eta n^{-1}X\text{diag}(|g_t \circ l_t|)X^T)$, then the estimated degrees of freedom at time t can be approximated by $n - \text{trace}(S_t)$. Consequently, the risk based on the C_p -type statistic (Efron, 2004) is

$$R(t) = \frac{\|r_t\|^2}{n} + \left(2 - \frac{2\text{trace}(S_t)}{n}\right)\sigma^2. \quad (9)$$

The total iteration number as a tunign parameter can then be selected by minimizing $R(t)$ in equation (6) or (7) . In practice, we can use the plug-in estimator $\hat{\sigma}$ to replace the unknown σ in $R(t)$. According to our simulation studies (for example, see figure 3), early stopping based on SURE generally works not as good as the hold-out or cross validation method.

Measure of model complexity: (Raskutti et al., 2014) proposed an early stopping rule based on local empirical Rademacher complexity of the Reproducing kernel Hilbert space. However, their method can not be directly applied in our case: their stopping rule is based on the eigenvalues of empirical kernel matrix, which is $n^{-1}X\text{diag}(|g_t \circ l_t|)X^T$ in our settings. Since our empirical kernel matrix keeps updated throughout the iterations, their method is not directly applicable.

In the end of this subsection, we adopt the simulation framework S1-S4 (only change the standard derivation to $\sigma = 0.1$) in section 4.2 to compare different early stopping criteria. We record the mean estimation errors averaging over 50 trials and report the errors in figure 3. From the figure, we can see that cross-validation tends to be more robust than SURE. Therefore, we recommend using hold-out or cross validation to select the iteration number, and will stick to this method in the rest of the paper.

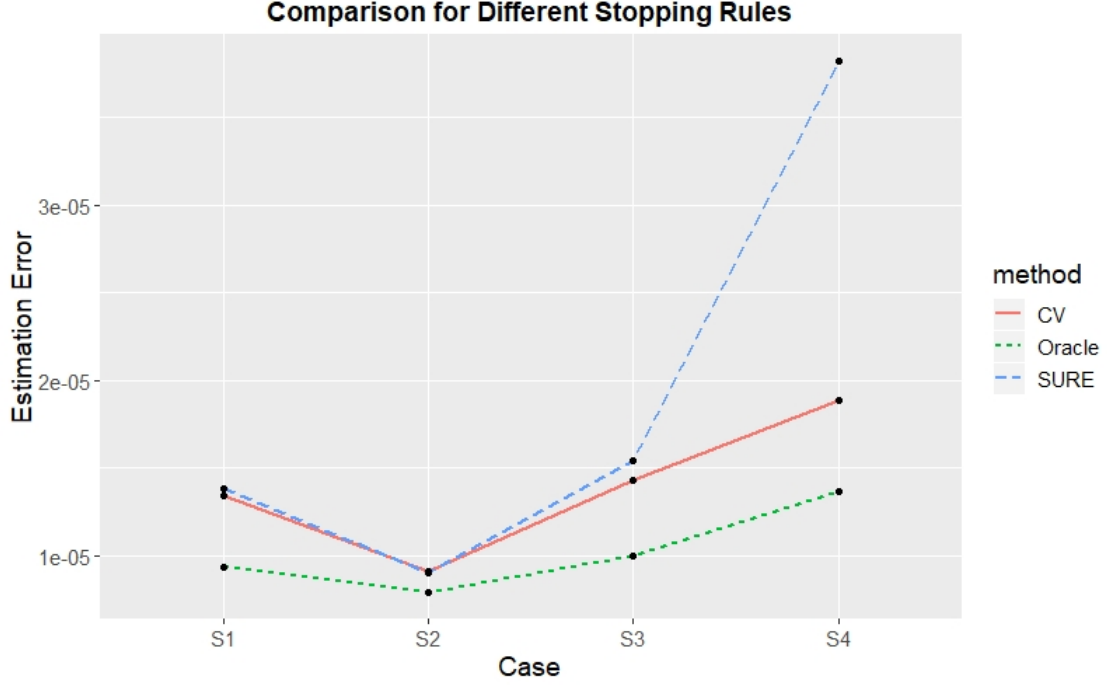


Figure 3: Comparison of the estimation errors for different early stopping rules, ‘Oracle’ stands for the optimal early stopping rule with knowledge on the truth. ‘CV’ stand for early stopping through 5 fold cross validation.

2.5 Adaptive step size and variable selection

A nature extension of gradient descent algorithm 1 is to assign different weights (step sizes) to different coordinates of β , which is related to the adaptive Lasso (Zou, 2006). It can be seen from the differential equation interpretation: by inserting a constant weighting matrix $D(\Omega) = \text{diag}(\omega_1, \dots, \omega_p)$ into the equation (4), we obtain the limiting dynamical system as

$$\begin{cases} \dot{g}(t) = -[D(\Omega)X^T r(t)] \circ l(t), \\ \dot{l}(t) = -[D(\Omega)X^T r(t)] \circ g(t). \end{cases}$$

Based on similar heuristic analysis as in Section 2.3 for the noiseless case, the limiting point of the dynamic system satisfies:

$$X_j^T s_\infty = \begin{cases} \text{sign}(\beta_{j,\infty})/\omega_j, & \text{if } \beta_{j,\infty} \neq 0, \\ \gamma_j \in [-\frac{1}{\omega_j}, \frac{1}{\omega_j}], & \text{if } \beta_{j,\infty} = 0, \end{cases} \quad \text{for each } j = 1, \dots, p.$$

which is the KKT condition for the dual form of the adaptive Lasso

$$\min_{\beta \in \mathbb{R}^p} \sum_{j=1}^p \frac{|\beta_j|}{w_j} \quad \text{subject to } X\beta = y.$$

In the limiting case when the step size ω_j of a particular component β_j tends to 0, we are equivalently adding an $+\infty$ when $\beta_j \neq 0$. In contrast, if we apply a larger step size ω_j to β_j , then $\beta_j = g_j \circ l_j$ tend to move faster and more freely in the parameter space, which is equivalent to a smaller penalty on β_j . The original paper in Zou (2006) constructed the weights based on the ordinary least square solution, which requires $n \geq p$. In practice when $p > n$, we can construct weights through a preprocessing step. For example, variable screening can be applied to introduce sparse weights.

To enable variable selection in our method, we can perform a component-wise hard thresholding operation to the final estimator $\hat{\beta} = g_{\hat{t}} \circ l_{\hat{t}}$. Based on our theoretical analysis, since our method tries to shrink both weak signals and errors into very small order p^{-2} , it is more robust to false detections than other explicit regularizations when the same tuning parameter for noise level is applied. Let us consider a simple example to illustrate the basic idea: we set $n = 10$, $p = 20$, $X_{ij} \stackrel{iid}{\sim} \mathbb{N}(0, 1)$ for $i = 1, 2, \dots, n$ and $j = 1, 2, \dots, p$, $\beta_1^* = 0.5\sigma\sqrt{\log p/n}$, $\beta_2^* = 5\sigma\sqrt{\log p/n}$, and all other components are zeros in the data generating model $y = X\beta^* + w$ with $w \sim \mathbb{N}(0, I)$. Since the strength of the first components of truth is weak, it is hard to be detected by all methods we have tried. However, the effect of the weak signals on y still pertains. In particular, when applying the cross validation, traditional penalized methods tends to over-select the predictors, leading a lot of false

discoveries. In comparison, due to the implicit regularization our method tends to be more robust to the weak signals—our estimate is typically non-sparse, the effect of the non-detected weak signals can be distributed to all components of the estimated vector, and no component is particularly spiked. As a consequence, our method tends to be more robust to false discoveries after applying the hard thresholding. The variable selection results are shown in figure 4. As we can see, the Lasso can not detect the weak signal, and two components, indexed by 6, 19, appear to be false detected through cross validation (note that in Lasso, a soft thresholding has already been applied). In contrast, in our method most unimportant signals remain small. Performing a hard thresholding with the same regularization parameter selected by the Lasso can erase all false detections, leading to the selection of strong signal only.

2.6 Related literatures

Li et al. (2018) studies the theory for implicit regularization in matrix sensing, which requires the data to be perfectly measured and has different geometric structures as linear regression. Hoff (2017) considers the Hadamard product parametrization to optimize the parameters in high-dimensional linear regression. However, his method is computational-wise in order to reformulate the non-smooth Lasso optimization problem into a smooth one. In particular, their objective function involves an ℓ_2 penalty on (g, l) (which is equivalent to the ℓ_1 penalty on β), and the solution is precisely the Lasso solution.

3 Theoretical Analysis

In this section, we provide formal statements for characterizing the behavior of the gradient descent algorithm for minimizing the Hadamard product parametrized quadratic loss $f(g, l)$ as defined in (2). We start with a description of our assumptions. Then we turn to the case of non-negative parameters, where a simpler parametrization $\beta = u \circ u$ can be applied,

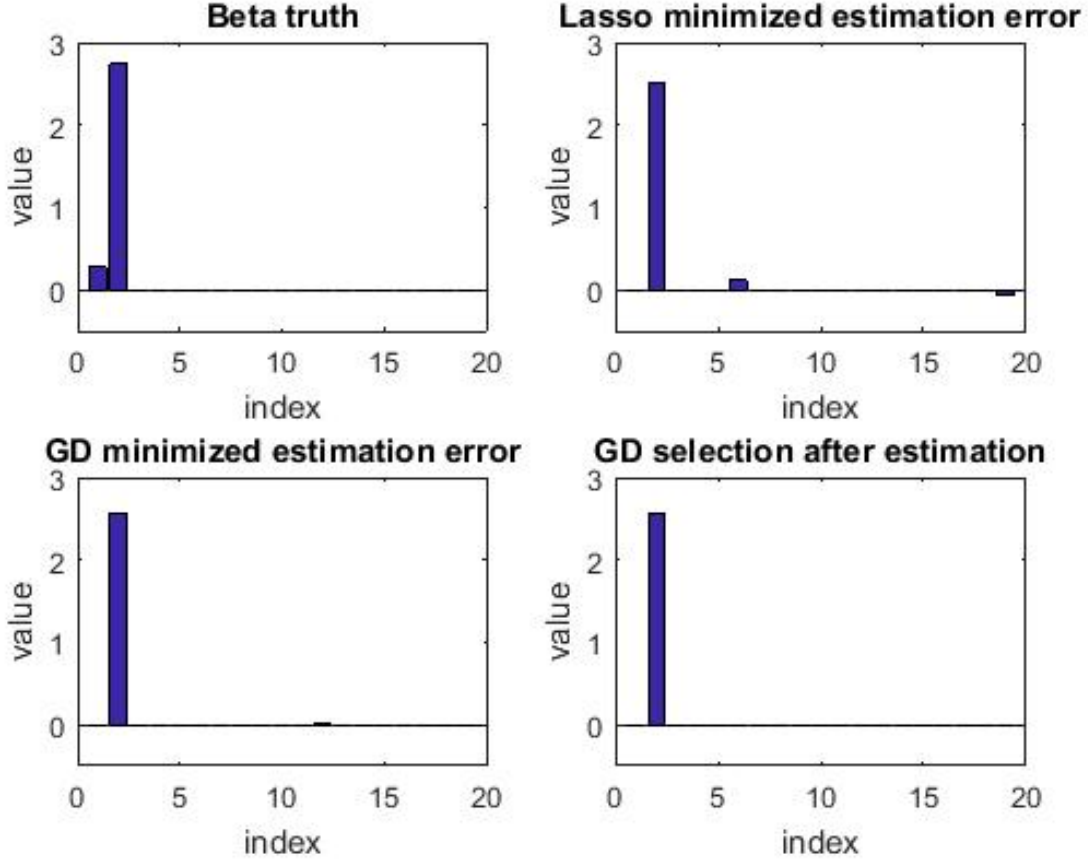


Figure 4: The values versus index for truth β^* , β estimator through lasso by minimizing the cross validation error in prediction, β estimator through gradient descent by minimizing the cross validation error in prediction and β estimator through ‘post estimation’ selection for gradient descent.

as a warm-up to convey the main proof ideas. Finally, we consider the general signal case and illustrate when the fast parametric root- n rate independent of the dimension can be achieved. All proofs are deferred to the appendix in the supplementary material of the paper.

3.1 Assumptions

Recall that the underlying true data generating model is $y = X\beta^* + w$ with $w \sim \mathcal{N}(0, \sigma^2 I)$, where the true parameter β^* is s -sparse. Within the s nonzero signal components of β^* , we define the index set of strong signals as $S_1 = \{i \in S : |\beta_i^*| \geq 2\sigma \log p \sqrt{\frac{\log p}{n}}\}$ and weak

signals as $S_2 = \{i \in S : |\beta_i^*| \leq 2\sigma\sqrt{\frac{\log p}{n}}\}$, where $|S_1| = s_1$, $|S_2| = s_2$. According to the information-theoretic limits from [Wainwright \(2009\)](#), weak signals of order $\sigma\sqrt{\log p/n}$ in sparse linear regression are generally impossible to be jointly recovered or selected (but can be detected in terms of the type I/II error control in hypothesis testings, e.g. see [Jin and Ke \(2016\)](#)). Therefore, our primary focus would be the estimation and selection of strong signals. We use the notation $\theta_{s_1}(\beta)$ to denote the s_1 -th largest absolute component value of β , and let $m = \theta_{s_1}(\beta^*)$, which reflects the minimal strength for strong signals. We also use κ to denote the strong signal-condition number as the ratio between the largest absolute signal value to the smallest strong signal. We will also make use of the notation of Restricted Isometry Property (RIP, [Candes \(2008\)](#)), which is a commonly used assumption (e.g. [Candes and Tao \(2007\)](#)) in the high dimensional linear regression literatures.

Definition 3.1 (Restricted Isometry Property). *A matrix $X \in \mathbb{R}^{n \times p}$ is said to satisfy the (s, δ) -Restricted Isometry Property (RIP) if for any s -sparse vector u in \mathbb{R}^p , we have:*

$$(1 - \delta)\|u\|^2 \leq \frac{1}{n} \|Xu\|^2 \leq (1 + \delta)\|u\|^2$$

As an easy consequence, if matrix X satisfies $(2s, \delta)$ -RIP, then Euclidean inner-product is also approximately preserved, that is, $\left|n^{-1} \langle Xu, Xv \rangle - \langle u, v \rangle\right| \leq \delta \|u\| \cdot \|v\|$ holds for any two s -sparse vectors $u, v \in \mathbb{R}^p$.

With these preparations, we make the following assumptions on the true parameter β^* , design matrix X , initialization parameter α and step size η in [Algorithm 1](#).

Assumption (A): The true parameter β^* is s -sparse, and $s = s_1 + s_2$, that is, each nonzero signal in β^* is either weak or strong. In addition, $\kappa m \lesssim 1$.

Assumption (B): The design matrix X satisfies $(s + 1, \delta)$ -RIP condition with $\delta \lesssim 1/(\kappa\sqrt{s} \log \frac{p}{\alpha})$.

Assumption (C): The initialization parameter α satisfies $0 < \alpha \lesssim p^{-1}$, and the step size η satisfies $0 < \eta \lesssim (\kappa \log \frac{p}{\alpha})^{-1}$.

Some remarks are in order. First, our current proof heavily relies on the RIP condition as in Assumption (B), which is satisfied if the predictors are iid and $s \log p \ll n$. However, the extensive simulation studies in the next section provide a strong numerical evidence suggesting that our conclusions remain valid even when the RIP condition is violated. We leave the formal theoretical investigation as an open problem for our future studies. Second, Assumption (A) is made mainly for illustrating the possibility of achieving the fast parametric root n rate of convergence when $s_1 = 0$. In fact, our current proof can still lead to the typical high-dimensional rate of $\sqrt{s \log p / n}$ without introducing the notion of strong and weak signals. And due to space constraint we omit the details.

3.2 Non-negative Signals

To start with, we demonstrate the key ideas of our proofs and give an analysis of the non-negative case as a warm-up. More specifically, we consider the case when all components of true signal β^* are non-negative. To exploit this non-negativeness, we may instead use the self-Hadamard product parametrization $\beta = u^2 = u \circ u$ for $u \in \mathbb{R}^p$, and write $\beta^* = (u^*)^2 = u^* \circ u^*$. Now, we have the optimization problem:

$$\min_{u \in \mathbb{R}^p} f(u) = \frac{1}{2n} \|Xu^2 - y\|^2,$$

with gradient descent updating formula $u_{t+1} = u_t - 2\eta u_t \circ [n^{-1} X^T (Xu_t^2 - y)]$. For technical simplicity, we instead focus on the deterministic initialization $u_0 = \alpha \mathbf{1} \in \mathbb{R}^p$. This case is simpler for the analysis than the general case since components of u_t will not change sign during the iterations, and will always stay away from saddle points. We summarize our result in the following main theorem. Since the non-negative signal constraint resembles the positive semi-definiteness constraint in matrix sensing, our proof utilizes the proof

strategy in [Li et al. \(2018\)](#) for analyzing matrix factorized gradient descent for noiseless matrix sensing by dividing the convergence into two stages (more details are provided after the theorem).

Theorem 3.1. *Under the above assumptions (A), (B) and (C). Let $\epsilon = \max\{\alpha^2, \sigma^2 \frac{Ms_1}{n}, \sigma^2 \frac{s_2 \log p}{n}\}$, $\tau = \max\{\delta\alpha, \sigma\sqrt{\frac{\log p}{n}}\}$ and any $M \geq 1$. Then there exist positive constants $(c_1, c_2, c_3, c_4, c_5)$ such that for every time t satisfying $c_1 \log(\frac{p}{\alpha})/(\eta m) \leq t \leq c_2/(\eta\tau)$, with probability at least $1 - p^{-c_4} - e^{-c_5 Ms}$, the time t -iterate u_t satisfies*

$$\|u_t^2 - \beta^*\|^2 \leq c_5 \epsilon,$$

This theorem tells us in high dimension linear regression, combining early stopping with implicit regularization can significantly improve the estimation accuracy. In particular, when all signals are strong ($s_1 = s$ and $s_2 = 0$), the estimate $\hat{\beta} = u_t^2$ attains a parametric rate $\sigma\sqrt{s_1/n}$ of convergence that is independent of the dimension p . In general when weak signals exist, then the overall rate $\sigma\sqrt{\frac{s_1}{n}} + \sigma\sqrt{\frac{s_2 \log p}{n}}$ depends on the number of weak (strong) signals, which is still minimax-optimal ([Zhao et al., 2018](#)). The same remark also applies to our theory in the general case.

Our proof strategy is to divide the convergence into two stages. Recall that $S = \{j : \beta_j^* \neq 0\}$ is the support set of true signal β^* , and $S_1 \subset S$ corresponds to the subset of all strong signals. In the first “burn-in” stage, we show that each component of the strong signal part u_{t,S_1} increases at an exponential rate in t until hitting $\sqrt{m}/2$, while the weak signal and error part u_{t,S_1^c} remains bounded by $\mathcal{O}(p^{-1})$. In the second stage, iterate u_t enters a geometric convergence region where u_t converges towards u^* at a linear rate up to some high-order error term, and then stay in a $O(\epsilon)$ neighborhood of u^* up to the time $\Theta(1/\tau)$. Therefore, the time interval $c_1 \log(\frac{p}{\alpha})/(\eta m) \leq t \leq c_2/(\eta\delta\tau)$ would be the theoretical “best solution region” corresponding to the stabilized region in [figure 1](#).

More specifically, in the proof we consider the decomposition of u_t into three parts:

strong signal part z_t , weak signal part d_t and error part e_t :

$$u_t = z_t + d_t + e_t, \quad \text{with} \quad z_t := I_{S_1} u_t \in \mathbb{R}^p, \quad d_t := I_{S_2} u_t \in \mathbb{R}^p \quad \text{and} \quad e_t := I_{S^c} u_t \in \mathbb{R}^p,$$

where recall that I_E is the diagonal matrix with ones on the positions indexed by subset $E \subset \{1, \dots, p\}$ and zero elsewhere. We use induction to prove the following results characterizing the gradient dynamics in the first stage. Recall that $\theta_{s_1}(b)$ denote the s_1 -th largest absolute component value of vector $b \in \mathbb{R}^p$ and m is the minimal strength of the strong signals.

Proposition 3.1 (Stage one dynamics). *Under assumptions of theorem 3.1, there are constants (c_7, c_8) and (c'_7, c'_8) , such that for each $t < T_1 = c_7 \log(\frac{p}{\alpha})/(\eta m)$, we have:*

$$\begin{aligned} \text{strong signal dynamics:} \quad & \theta_s(z_{t+1}) \geq \min \left\{ \frac{\sqrt{m}}{2}, \left(1 + \frac{\eta m}{4}\right) \theta_s(z_t) \right\}, \\ & \|z_t\|_\infty \leq c'_7, \quad \text{and} \quad \|z_t^2 - I_S \beta^*\| \leq c'_7 \sqrt{s}; \\ \text{weak signal dynamics:} \quad & \|d_{t+1}\|_\infty \leq \left(1 + c_8 \eta m / \log(p\alpha)\right) \|d_t\|_\infty \leq c'_8/p. \\ \text{error dynamics:} \quad & \|e_{t+1}\|_\infty \leq \left(1 + c_8 \eta m / \log(p\alpha)\right) \|e_t\|_\infty \leq c'_8/p. \end{aligned}$$

Define $\Theta_{LR} = \{(z, d, e) \in \mathbb{R}^{3p} : z_{S_1^c} = 0, d_{S_2^c} = 0, \theta_s(z) \geq \sqrt{m}/2, \|z\|_\infty \leq c'_7, \|z^2 - I_S \beta^*\| \leq c'_7 \sqrt{s}, \|d\|_\infty \leq c'_8/p, \|e\|_\infty \leq c'_8/p\}$ as a good parameter region for (z_t, d_t, e_t) . Proposition 3.1 tells that for all $t < T_1$, iterate (z_t, d_t, e_t) satisfies all constraints of Θ_{LR} except $\theta_s(z) \geq \sqrt{m}/2$, but will enter this good region in at most $\mathcal{O}(\log(\frac{m}{\alpha^2})/(\eta m))$ iterations. The second stage starts when (z_t, d_t, e_t) first enters Θ_{LR} . The next result summarizes the behavior of the gradient dynamics in the second stage—once it enters Θ_{LR} , it will stay in Θ_{LR} up to time $t = \Theta(1/\tau)$, where u_t converges towards u^* at a linear rate up to the statistical error ϵ .

Proposition 3.2 (Stage two dynamics). *Under assumptions of theorem 3.1, there exists*

some constant c_9 such that if $(z_t, d_t, e_t) \in \Theta_{LR}$, then $(z_{t+1}, d_{t+1}, e_{t+1}) \in \Theta_{LR}$, and

$$\|u_{t+1}^2 - \beta^*\|^2 \leq (1 - \eta m) \|u_t^2 - \beta^*\|^2 + c_9 \epsilon.$$

A combination of these two propositions, whose proofs are deferred to the supplementary material, leads to a proof of Theorem 3.1.

3.3 General Signals

Now let us consider the general case where signs of nonzero components in signal β^* are unknown. In this case, we need the full Hadamard product parameterization $\beta = g \circ l$ and want to solve

$$\min_{g, l \in \mathbb{R}^p} f(g, l) = \frac{1}{2n} \|X(g \circ l) - y\|^2,$$

where the associated gradient descent updating rule becomes:

$$\begin{cases} g_{t+1} = g_t - \eta l_t \circ [n^{-1} X^T (X g_t \circ l_t - y)] \\ l_{t+1} = l_t - \eta g_t \circ [n^{-1} X^T (X g_t \circ l_t - y)] \end{cases} \quad \text{with initial condition} \quad \begin{cases} g_0 = \alpha \mathbf{1} \\ l_0 = \mathbf{0} \end{cases},$$

Here again for technical simplicity, we focus on a deterministic initialization as above. Our main result is summarized in the following theorem.

Theorem 3.2. *Under the above assumptions (A), (B) and (C). Let $\epsilon = \max\{\alpha^2, \sigma^2 \frac{Ms_1}{n}, \sigma^2 \frac{s_2 \log p}{n}\}$, $\tau = \max\{\delta\alpha, \sigma\sqrt{\frac{\log p}{n}}\}$ and any $M \geq 1$. Then there exist positive constants $(\tilde{c}_1, \tilde{c}_2, \tilde{c}_3, \tilde{c}_4, \tilde{c}_5)$ such that for every time t satisfying $\tilde{c}_1 \log(\frac{p}{\alpha})/(\eta m) \leq t \leq \tilde{c}_2/(\eta\tau)$, with probability at least $1 - p^{-\tilde{c}_4} - e^{-\tilde{c}_5 Ms}$, the time t -iterate $\beta_t = g_t \circ l_t$ satisfies*

$$\|\beta_t - \beta^*\|^2 \leq \tilde{c}_5 \epsilon.$$

The proof strategy of this result is similar to the previous nonnegative signal case

by considering two stages. However, due to the lack of knowledge on signs of nonzero components in signal β^* , the proof of Theorem 3.2 has the extra component in showing that β_t tends to shoot towards the right direction (positive versus negative) during the first “burn-in” stage before entering the second geometric convergence region as in the proof for nonnegative signals in the previous subsection. To deal with this extra complication, we introduce another reparametrization in the proof:

$$a_t = (g_t + l_t)/2 \in \mathbb{R}^p, \quad \text{and} \quad b_t = (g_t - l_t)/2 \in \mathbb{R}^p \quad \text{for} \quad t = 0, 1, \dots$$

Notice that $\beta_t = a_t^2 - b_t^2$, implying that the sign of each component $\beta_{t,j}$ is determined by the relative magnitude of $|a_{t,j}|$ and $|b_{t,j}|$. These new variables (a_t, b_t) are especially convenient for understanding the dynamics of the sign of β_t due to a more tractable updating formula:

$$a_{t+1} = a_t - \eta a_t \circ \left[n^{-1} X^T (X \beta_t - y) \right], \quad \text{and} \quad b_{t+1} = b_t + \eta b_t \circ \left[n^{-1} X^T (X \beta_t - y) \right].$$

More precisely, we have the following results characterizing the gradient dynamics in the first stage—showing that depending on the sign of β_j^* for $j \in S_1$, either $a_{j,t}^2$ or $b_{j,t}^2$ will increase at an exponential rate while the other remaining small. Recall that S_1 denotes the support of strong signals of β^* .

Proposition 3.3 (Stage one dynamics). *Under assumptions of theorem 3.2, there are constants $(\tilde{c}_7, \tilde{c}_8, \tilde{c}_9)$ and $(\tilde{c}'_7, \tilde{c}'_8, \tilde{c}'_9)$, such that for each $t < T_1 = \tilde{c}_7 \log(\frac{p}{\alpha})/(\eta m)$, we have:*

$$\text{signal dynamics:} \quad \text{sign}(\beta_{S_1}^*) \circ \beta_{t,S_1} \geq \min \left\{ \frac{m}{2}, \left(1 + \eta m/4 \right)^t \alpha^2 - \tilde{c}'_9 \alpha^2 \right\} \mathbf{1} \in \mathbb{R}^{s_1},$$

$$\max \left\{ \|a_{t,S_1}\|_\infty, \|b_{t,S_1}\|_\infty \right\} \leq \tilde{c}'_7, \quad \text{and} \quad \|\beta_{t,S_1} - \beta_{S_1}^*\| \leq \tilde{c}'_7 \sqrt{s_1};$$

$$\text{weak signals and error dynamics:} \quad \|a_{t+1,S_1^c}\|_\infty \leq \left(1 + \tilde{c}_8 \eta m / \log(p/\alpha) \right) \|a_{t,S_1^c}\|_\infty \leq \tilde{c}'_8/p,$$

$$\text{and} \quad \|b_{t+1,S_1^c}\|_\infty \leq \left(1 + \tilde{c}_8 \eta m / \log(p/\alpha) \right) \|b_{t,S_1^c}\|_\infty \leq \tilde{c}'_8/p.$$

Similarly, let $\Theta_{LR}^G = \{(a, b) \in \mathbb{R}^{2p} : \text{sign}(\beta_{S_1}^*) \circ (a_{S_1}^2 - b_{S_1}^2) \geq (m/2) \mathbf{1}, \max \{\|a_{S_1}\|_\infty, \|b_{S_1}\|_\infty\} \leq \tilde{c}'_7, \|a_{S_1}^2 - b_{S_1}^2 - \beta_{S_1}^*\| \leq \tilde{c}'_7 \sqrt{s_1}, \max \{\|a_{S_1^c}\|_\infty, \|b_{S_1^c}\|_\infty\} \leq \tilde{c}'_8/p\}$ as a good parameter region for (a_t, b_t) . Proposition 3.3 tells that for all $t < T_1$, iterate (a_t, b_t) satisfies all constraints of Θ_{LR}^G except $\text{sign}(\beta_{S_1}^*) \circ (a_{S_1}^2 - b_{S_1}^2) \geq m/2$, but will enter this good region in at most $\mathcal{O}(\log(\frac{m}{\alpha^2})/(\eta m))$ iterations. The second stage starts when (a_t, b_t) first enters Θ_{LR}^G . The next result summarizes the behavior of the gradient dynamics in the second stage—once it enters Θ_{LR}^G , it will stay in Θ_{LR}^G for a long time $\Theta(1/\tau)$, where $\beta_t = a_t^2 - b_t^2$ converges toward β^* at a linear rate up to the statistical error ϵ .

Proposition 3.4 (Stage two dynamics). *Under assumptions of theorem 3.2, there exists some constant \tilde{c}_9 such that if $(a_t, b_t) \in \Theta_{LR}^G$, then $(a_{t+1}, b_{t+1}) \in \Theta_{LR}^G$, and $\beta_t = a_t^2 - b_t^2$ satisfies*

$$\|\beta_{t+1} - \beta^*\|^2 \leq (1 - \eta m) \|\beta_t - \beta^*\|^2 + \tilde{c}_9 \epsilon.$$

A combination of these two propositions leads to a proof of Theorem 3.2.

Finally, we provide a theorem regarding the strong signal selection consistency.

Theorem 3.3 (Variable selection consistency). *Under assumptions of theorem 3.2, for every time t satisfying $\tilde{c}_1 \log(\frac{p}{\alpha})/(\eta m) \leq t \leq \tilde{c}_2/(\eta \tau)$, with probability at least $1 - p^{-\tilde{c}_3} - n^{-\tilde{c}_4}$, the following hard thresholded estimator β'_t ,*

$$\beta'_{t,j} = \begin{cases} g_{t,j} l_{t,j}, & \text{if } |g_{t,j} l_{t,j}| \geq \lambda \\ 0, & \text{if } |g_{t,j} l_{t,j}| < \lambda, \end{cases},$$

for any $\lambda \in \left[\frac{\tilde{c}_{10}}{p}, \tilde{c}_{11} \sigma \sqrt{\frac{\log p}{n}} \right]$ where \tilde{c}_{10} and \tilde{c}_{11} are positive constants, satisfies variable selection consistency, that is, $\{j : \beta'_{t,j} \neq 0, j = 1, \dots, p\} = S_1$.

This theorem implies the variable selection consistency of our method is not sensitive to the choice of tuning parameter of noise level λ . As shown in the toy example in section 2.5, the Lasso suffers from false positives due to non-discoveries of weak signals when the cross

validation error in prediction is applied as the criterion. However, even when we apply the same regularization to our method, since the tuning parameter is still in the region $\left[\frac{\tilde{c}_{10}}{p}, \tilde{c}_{11}\sigma\sqrt{\frac{\log p}{n}}\right]$, the variable selection consistency of our method can still achieve. As we will illustrate in the numerical part below, our method simultaneously achieves the best false discovery control and attains the most accurate prediction.

4 Simulations and Real Data Analysis

We start with examples with noiseless data to show that under this setting, our gradient descent Algorithm 1 truly converges to the least ℓ_1 -norm global minimizer. These examples provide strong numerical evidences to our statements in Section 3 even when the RIP condition is violated. After that, we consider the general noisy settings, and compare our methods with state-of-the-art approaches. This section ends up with a real data analysis.

4.1 Simulations for Noiseless Case

In this section, we conduct some numerical experiments to illustrate the performance of our method. In our first example, we study the performance as we change the initialization under settings where the RIP may or may not hold. In the second example, we study the behavior of the algorithm when the null space property (Cohen et al., 2009) is violated, under which the ℓ_1 -norm minimizer in the basis pursuit problem differs from the sparsest solution (the ℓ_0 -norm minimizer).

Impact of initialization: We consider two classes of random designs $X \in \mathbb{R}^{n \times p}$ with i.i.d. rows $X^{(i)} \stackrel{iid}{\sim} \mathcal{N}(0, \Sigma) \in \mathbb{R}^p$, where population p.s.d covariance matrix $\Sigma \in \mathbb{R}^{p \times p}$ is

independent design: $\Sigma_{jk} = \mathbb{I}(j = k)$, or correlated design: $\Sigma_{jk} = 0.5 + 0.5\mathbb{I}(j = k)$,

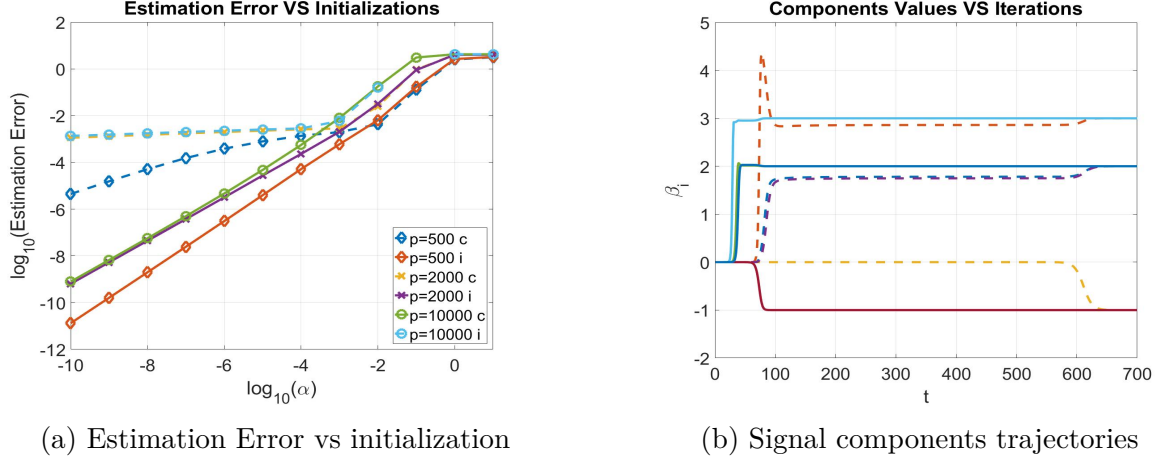


Figure 5: Panel (a) is a log-log plot of estimation error $\|\hat{\beta} - \beta^*\|_2$ versus initialization level α . Panel (b) shows trajectories of signal components (β_1 to β_4) at $(n, p) = (200, 2000)$. In both plots, solid curves correspond to the independent (i) design, and dashed curves the correlated (c) design.

for $j, k = 1, \dots, p$, where $\mathbb{I}(A)$ denotes the indicator function of event A . We choose sample size $n = 200$, sparsity level $s = 4$, and signal dimension $p \in \{500, 2000, 10000\}$. For the independent (correlated) design, the RIP is satisfied (fails) with high probability for some small $\delta > 0$. In both scenarios, we choose true signal $\beta^* = (-1, 2, 2, 3)^T \in \mathbb{R}^p$, and set $y = X\beta^*$. When implementing gradient descent, we choose step size $\eta = 0.2$ (0.1) for the independent (correlated) design, $\alpha \in \{10^{-10}, 10^{-9}, \dots, 10^1\}$, and stopping threshold $\epsilon = 0.01\alpha$. Figure 5 shows the estimation error $\|\hat{\beta} - \beta^*\|_2$ versus α in log-log plots. As we can see, they tend to have a linear trend under the log-scale, which is consistent with our theoretical error bound estimate in Section 4. In addition, in the correlated design scenario where the RIP does not hold, the algorithm is still able to recover β^* as $\alpha \rightarrow 0$, albeit under a slower convergence (due to a smaller allowable step size and a larger condition number of X). This observation provides evidence to the correctness of our informal statement made at the beginning of Section 3 even without RIP condition. We leave the proof of this conjecture open.

ℓ_0 -norm minimizer differs from ℓ_1 -norm minimizer: In this example, we study the empirical performance of the algorithm when the least ℓ_1 -norm in the basis pursuit

problem (1) is not the sparsest solution of $X\beta = y$ (the null space property is violated). In particular, we choose

$$X = \begin{bmatrix} 0.2 & 1 & 0 \\ 0.2 & 0 & -1 \end{bmatrix}, \quad \beta^* = \begin{bmatrix} 5 \\ 0 \\ 0 \end{bmatrix}, \quad \text{and} \quad y = \begin{bmatrix} 1 \\ 1 \end{bmatrix},$$

so that $X\beta^* = y$. It is easy to verify that for this example, the sparsest solution of $X\beta = y$ is β^* , while the least ℓ_1 -norm solution is $\beta^\dagger = [0, 1, -1]^T$. We use the same setting as before for implementing the algorithm with $\alpha \in \{10^{-10}, 10^{-5}, 10^{-3}, 10^{-1}, 10^0, 10^1\}$. Table 1 reports final outputs $\beta = (\beta_1, \beta_2, \beta_3)^T$ of the algorithm. Due to our intuition in Section 3, as expected, the algorithm still converges to the least ℓ_1 -norm solution β^\dagger instead of the least ℓ_0 -norm solution β^* . Again, the estimation error decreases as the initialization level α decreases. We conjecture this phenomenon of convergence to the least ℓ_1 -norm solution to be generally true, and leave a formal proof as a future direction.

Table 1: Convergent point without null space property

α	10^{-10}	10^{-5}	10^{-3}	0.1	1	10
β_1	$7.433e-13$	$5.703e-7$	$1.289e-4$	$2.884e-2$	$2.987e-1$	$8.823e-1$
$1 - \beta_2$	$1.492e-13$	$1.141e-7$	$2.577e-5$	$5.769e-3$	$5.974e-2$	$1.765e-1$
$1 + \beta_3$	$1.492e-13$	$1.141e-7$	$2.577e-5$	$5.769e-3$	$5.974e-2$	$1.765e-1$

4.2 Simulations for Noisy Case

Comparison with other high dimensional estimators: We further demonstrate the advantages of our algorithm by considering the following 8 simulation settings, the sparsity level $s = 4$, signals $-1, 2, 2, 3$ and noise level σ with $\sigma = 0.15 * \|\beta^*\|$. We generate $3n$ observations independently and split into 3 even parts, then use the first part for training, the second part for validation and the final part for testing. The evaluation metric is standardized estimation error $\|\hat{\beta} - \beta^*\|_2^2 / \|\beta^*\|_2^2$ and mean prediction error $\sqrt{\|y - \hat{y}\|^2 / n}$

for the test data set. We compare the median of the standardized estimation errors and prediction errors with the Lasso, SCAD and MCP by repeating 50 times. We implement the Lasso using *glmnet* R package (Friedman et al., 2010) while for SCAD and MCP, we use the R package *ncvreg* (Breheny and Huang, 2011). The standard error of medians are calculated by bootstrapping the calculated errors 1000 times. For our algorithm, we use the initialization $\alpha = 10^{-5} \times \mathbf{1}$. The simulation results in Table 2 and 3 indicate that our methods consistently have the best performance over all explicit penalization-based competitors across all settings.

1. **S1:** $n = 200, p = 500, \Sigma_{jk} = 1$ for $j = k$ while $\Sigma_{jk} = 0$ for $j \neq k$;
2. **S2:** $n = 200, p = 500, \Sigma_{jk} = 0.1^{|j-k|}$;
3. **S3:** $n = 200, p = 500, \Sigma_{jk} = 0.2^{|j-k|}$;
4. **S4:** $n = 200, p = 500, \Sigma_{jk} = 0.5^{|j-k|}$;
5. **S5:** $n = 200, p = 2000, \Sigma_{jk} = 1$ for $j = k$ while $\Sigma_{jk} = 0$ for $j \neq k$;
6. **S6:** $n = 200, p = 2000, \Sigma_{jk} = 0.1^{|j-k|}$;
7. **S7:** $n = 200, p = 2000, \Sigma_{jk} = 0.2^{|j-k|}$;
8. **S8:** $n = 200, p = 2000, \Sigma_{jk} = 0.5^{|j-k|}$.

	S1	S2	S3	S4	S5	S6	S7	S8
GD	0.520 (0.0428)	0.448 (0.0530)	0.510 (0.0607)	0.568 (0.0850)	0.385 (0.0533)	0.290 (0.0465)	0.465 (0.0858)	0.460 (0.0863)
Lasso	3.11 (0.219)	3.10 (0.173)	3.42 (0.242)	4.51 (0.274)	4.62 (0.279)	4.04 (0.205)	4.40 (0.306)	6.98 (0.452)
SCAD	0.613 (0.0464)	0.533 (0.0679)	0.650 (0.0702)	0.691 (0.103)	0.519 (0.0527)	0.401 (0.0574)	0.595 (0.0837)	0.646 (0.0776)
MCP	0.628 (0.0392)	0.552 (0.0779)	0.594 (0.0809)	0.733 (0.0902)	0.484 (0.0706)	0.405 (0.0597)	0.595 (0.0741)	0.708 (0.0680)

Table 2: Simulation result for median of standardized estimation error of each method, with standard derivation in the parenthesis under the median. There are 10^{-3} factors for all medians and standard derivations.

	S1	S2	S3	S4	S5	S6	S7	S8
GD	0.638 (0.0597)	0.636 (0.0753)	0.640 (0.0718)	0.634 (0.0709)	0.646 (0.0491)	0.651 (0.0510)	0.641 (0.0406)	0.642 (0.0498)
Lasso	0.676 (0.0899)	0.672 (0.0981)	0.671 (0.0932)	0.685 (0.0510)	0.693 (0.0568)	0.699 (0.0918)	0.696 (0.0615)	0.708 (0.0488)
SCAD	0.638 (0.0530)	0.638 (0.0724)	0.637 (0.0716)	0.637 (0.0713)	0.650 (0.0470)	0.654 (0.0451)	0.643 (0.0402)	0.647 (0.0434)
MCP	0.637 (0.0516)	0.639 (0.0756)	0.638 (0.0684)	0.637 (0.0753)	0.650 (0.0475)	0.652 (0.0441)	0.644 (0.0435)	0.647 (0.0408)

Table 3: Simulation result for median of mean prediction error of each method, with standard derivation in the parenthesis under the median. There are 10^{-1} factors for all standard derivations.

Comparison in Variable Selection: Now we consider variable selection when there exists some weak signals. Suppose the simulation settings are similar with S3 above, with only the true signals changed. Let $s = 20$, and the strength of first 4 signals is $0.5\sigma\sqrt{\log p/n}$, while the other 16 are $5\sigma\sqrt{\log p/n}$, where $\sigma = 1$. Clearly the first 4 signals are too weak to be selected by all methods. However, since all methods are based minimizing the prediction error, the effect of these weak signals pertains, and may increase the false discovery rate. Under the above settings, we perform a model selection based on minimized prediction errors through 5-fold cross validation. For our method, we use the same regularization parameter as the lasso to perform hard threshold after estimation. We repeat the process 50 times and compare variable selection errors. From the figure 6, we can see our method is more robust to the enhancement of false detection due to failure on detecting weak signals: although the true negative errors of our method is 4, which means all weak signals can not be detected, the false detections of our methods are closed to zero. For other methods, although sometimes weak signals can be detected, the risk of false detections is high. Overall, our methods perform consistent variable selection for strong signals, and achieve better estimation than the competitors.

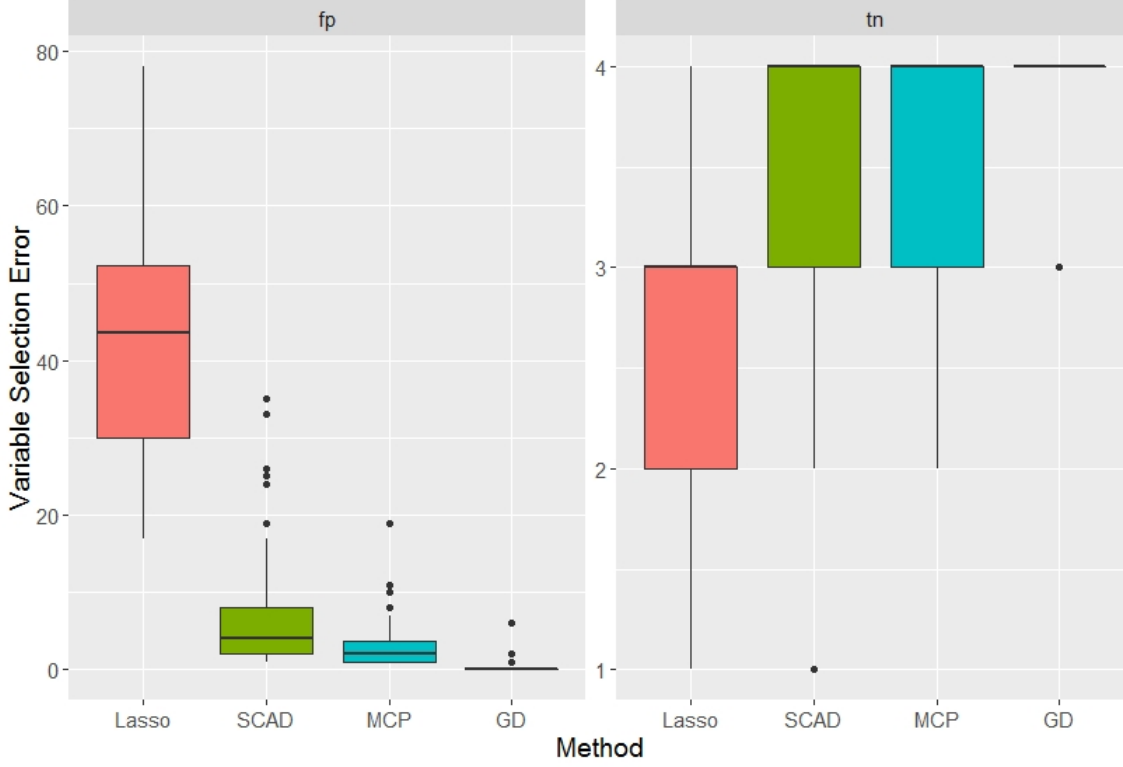


Figure 6: Variable selection errors for selected model based on minimized prediction cross validation errors. ‘fp’ stands for false positive when the truth is zero but detected as signal; ‘tn’ stands for true negative when the truth is nonzero but not detected.

4.3 Real Data Analysis

We compare our method with others to analyze the Riboflavin data set (Bühlmann et al., 2014), which is available in *hdi* R package. The dataset contains 71 observations of log-transformed riboflavin production rate verses the logarithm of expression level of 4088 genes. Before estimation, we first perform independence screening (Fan and Lv, 2008) based on the rank of the correlation strength for each predictor verses response to decrease the dimension of feature space into 500. Then we normalize and add the intercept column into the design matrix. For evaluation, we split the observations into 50 training samples and 21 testing samples, with performing 10-fold cross validation to select iteration steps and regularization parameters in the training data. Still, for our algorithm, we use the initial value $\alpha = 10^{-5} \times \mathbf{1}$ for all training processes. We record the prediction errors for

testing data set and repeat 50 times. From the figure 7 below, our method also obtain the least prediction errors, which implies the estimation of this high dimensional linear regression problems can also have the least errors.

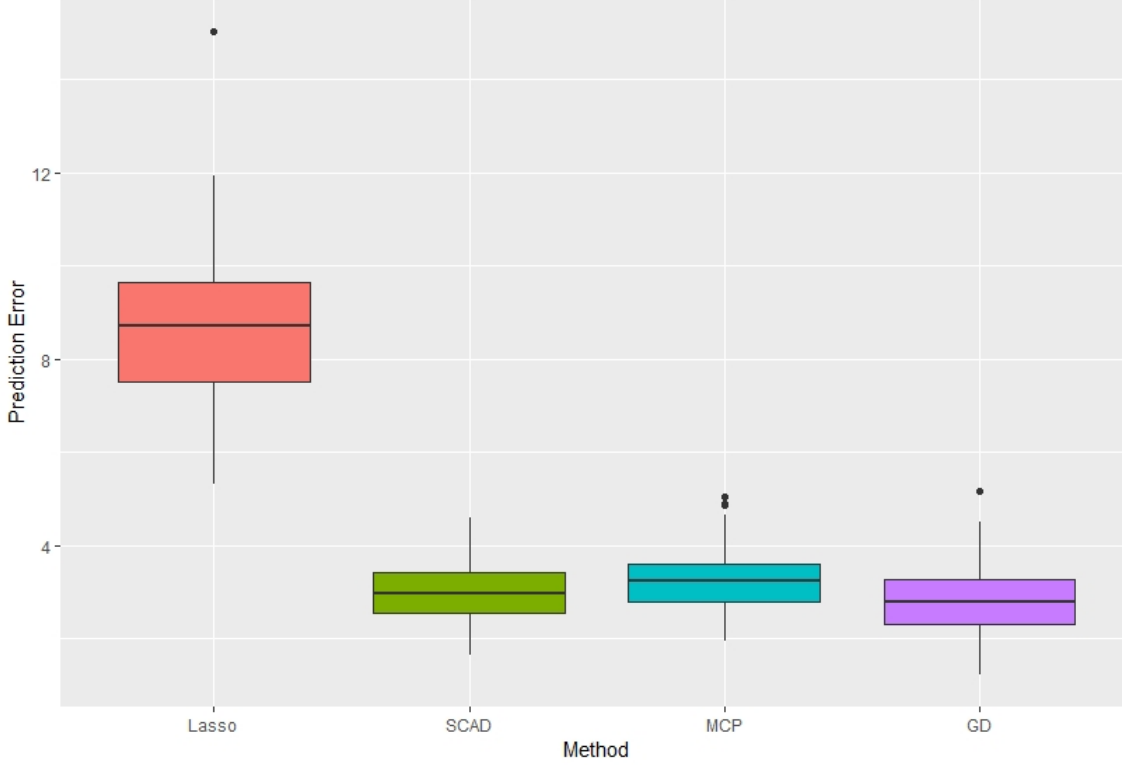


Figure 7: Prediction errors on the test data of Riboflavin data set for each method. x -axis stands for the methods used for estimation, and y -axis stands for the testing prediction error $\|y - \hat{y}\|$.

We also perform variable selection on the whole Riboflavin data set with the same tuning parameter obtained through minimized estimation errors based on cross validation. For our algorithm, when we get the number of iterations from cross validation, we run gradient descent on the whole data set with the same initial values and step size until the corresponding number of iterations. We use the same regularization parameter obtained by the lasso as the the hard thresholding value on the absolute value of obtained the ‘post-estimation’ selection. The comparison between different variable selection methods is given in Table 4. From the table, we can see except one variable, all other variables detected by

our method are also selected by other methods, illustrating that our methods tend to have lower false positive rate for variable selection without sacrificing estimation and prediction accuracy.

	Lasso	SCAD	MCP	GD
Lasso	33			
SCAD	11	14		
MCP	2	3	5	
GD	8	9	3	11
Independent	20	2	2	1

Table 4: Variable selection result for Riboflavin data set, each cell stands for the number of the same detected variables between row labels and column labels. ‘Independent’ means the number of detected variables for the corresponding column method that are not detected by others.

5 Discussion

In this paper, we illustrated the phenomenon of implicit regularization through Hadamard product change of variables in high dimensional linear regression, and demonstrated that a combination of implicit regularization with early stopping yields better estimation than state-of-the-art penalized approaches with explicit regularization. However, we still face several important open problems on implicit regularizations as our future directions. First, our theory heavily relies on the RIP condition, which is relatively strong comparing to the restricted eigenvalue condition as the minimal possible assumption on the design in the literature. It would be interesting to investigate whether our results remain valid without the RIP condition. Second, it is interesting to study whether any computationally efficient early stopping rule (rather than cross validation) based on certain data-driven model complexity measure can be applied and provably works.

6 Proof of the results in the paper

6.1 Overview

In this supplementary material, we provide proofs of the main theorems presented in the paper.

6.2 Notation

Recall that $\|v\| = \sqrt{\sum_{j=1}^p v_j^2}$ and $\|v\|_\infty = \max_j |v_j|$ denote the vector ℓ_2 -norm and ℓ_∞ -norm, respectively. Moreover, I is the identity matrix in \mathbb{R}^p , and for any subset S of $\{1, \dots, p\}$, I_S is the diagonal matrix with 1 on the j th diagonals for $j \in S$ and 0 elsewhere. We use bold letter $\mathbf{1} \in \mathbb{R}^p$ to denote an all-one vector. $\theta_s(\beta)$ denote the s -largest component of vector $\beta \in \mathbb{R}^p$ in absolute value. We use notation \lesssim and \gtrsim to denote \leq and \geq up to some positive multiplicative constant, respectively. For two vectors u and v of the same dimension, we use $a \geq b$ and $a \leq b$ to denote element-wise \geq and \leq . Denote $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$ be the maximal and minimal eigenvalues of matrix A . Through this document, letters c , c' and c'' denote some constants whose meaning may change from line to line.

6.3 Some Useful Results

In our proof, we will constantly deal with the Hadamard product $u \circ v$ and the operation $(n^{-1} X^T X u) \circ v$ for two vectors $u, v \in \mathbb{R}^p$. Therefore, we collect some useful properties in this section, some of which are consequences of the RIP condition.

The first property regarding the Hadamard product is a direct consequence of the Hölder inequality.

Lemma 6.1. *For any two vectors u and v in \mathbb{R}^p , we have:*

$$\|u \circ v\| \leq \|u\| \|v\|_\infty. \quad (10)$$

Proof. This follows since $\|u \circ v\|^2 = \sum_j u_j^2 v_j^2 \leq \|v\|_\infty^2 \sum_j u_j^2 = \|u\|^2 \|v\|_\infty^2$. \square

The second lemma shows that under the RIP, the product $(n^{-1} X^T X u) \circ v$ can be well-approximated by $u \circ v$ for all sparse vectors $u \in \mathbb{R}^p$ and any vector $v \in \mathbb{R}^p$.

Lemma 6.2. *Let X be a matrix in $\mathbb{R}^{n \times p}$ that satisfies $(s+1, \delta)$ -restricted isometry property (see Definition 2.1 in the paper). Then for any s -sparse vectors u and any v in \mathbb{R}^p , we have:*

$$\|(n^{-1} X^T X u) \circ v - u \circ v\|_\infty \leq \delta \|u\|_2 \|v\|_\infty. \quad (11)$$

Proof. Let $D(v)$ be the diagonal matrix in $\mathbb{R}^{n \times p}$ with diagonal elements the same as components of v correspondingly. Then $\|(n^{-1} X^T X u) \circ v - u \circ v\|_\infty$ can be represented as:

$$\max_{i=1,2,\dots,p} |e_i^T D(v) n^{-1} X^T X u - e_i^T D(v) u|,$$

where e_i is a p -dimensional vector whose i -th component is 1 and 0 elsewhere. Using the fact that X satisfies $(s+1, \delta)$ -RIP and $e_i^T D(v)$ is 1-sparse, we have (see the remark right after Definition 2.1 in the paper):

$$\|(n^{-1} X^T X u) \circ v - u \circ v\|_\infty \leq \max_{i=1,2,\dots,p} \delta \|e_i^T D(v)\| \|u\| = \delta \|u\|_2 \|v\|_\infty.$$

\square

Our third lemma considers the case when u and v are both arbitrary.

Lemma 6.3. *Let X be a matrix in $\mathbb{R}^{n \times p}$ that satisfies $(2, \delta)$ -restricted isometry property. Then for any vectors $u, v \in \mathbb{R}^p$, we have:*

$$\|(n^{-1} X^T X u) \circ v - u \circ v\|_\infty \leq \delta \|u\|_1 \|v\|_\infty. \quad (12)$$

Proof. Since we can decompose $u = \sum_j I_j u$, we have

$$\begin{aligned}
\|(n^{-1} X^T X u) \circ v - u \circ v\|_\infty &= \max_{i=1,2,\dots,p} |e_i^T D(v) n^{-1} X^T X u - e_i^T D(v) u| \\
&\leq \sum_{j=1}^p \max_{i=1,2,\dots,p} |e_i^T D(v) n^{-1} X^T X I_j u - e_i^T D(v) I_j u| \\
&\leq \sum_{j=1}^p \max_{i=1,2,\dots,p} \delta \|e_i^T D(v)\| \|u_j\| \\
&\leq \delta \|u\|_1 \|v\|_\infty.
\end{aligned}$$

□

Our fourth and fifth lemma consider the concentration behavior about the noise terms.

Lemma 6.4. $w \sim \mathcal{N}(0, \sigma^2 I_{n \times n})$, all ℓ_2 norm of column vectors of $X_{n \times s}$ are normalized to \sqrt{n} , $s < n$, then with probability $1 - e^{-Ms}$ for any $M \geq 1$, such that:

$$\frac{1}{n} \|X^\top w\| \leq 3\sigma \sqrt{\frac{Ms \lambda_{\max}(\frac{X^\top X}{n})}{n}}. \quad (13)$$

Proof. By directly applying Hanson-Wright inequality:

$$\mathbb{P}(w^\top X X^\top w \geq 9\sigma^2 n M s \lambda_{\max}(\frac{X^\top X}{n})) \leq \exp\left(-\frac{9\sigma^2 n M s \lambda_{\max}(\frac{X^\top X}{n})}{8\sqrt{2}\sigma^2 \lambda_{\max}(X^\top X)}\right) \leq e^{-Ms}. \quad (14)$$

□

Lemma 6.5. $w \sim \mathcal{N}(0, \sigma^2 I_{n \times n})$, all ℓ_2 norm of column vectors of $X_{n \times p}$ are normalized to \sqrt{n} , then with probability $1 - 2p^{-1}$ such that:

$$\frac{1}{n} \|X^\top w\|_\infty \leq \sigma \sqrt{\frac{4 \log p}{n}}. \quad (15)$$

Proof. $\frac{1}{n} \|X^\top w\|_\infty$ is the maximum over p normal random variables distributed as $\mathcal{N}(0, \frac{\sigma^2}{n})$,

so considering the union bound:

$$\mathbb{P}(\frac{1}{n}\|X^\top w\|_\infty \geq t) \leq 2e^{\frac{-nt^2}{2\sigma^2} + \log p},$$

the result is obtained by taking $t = \sigma\sqrt{\frac{4\log p}{n}}$. □

6.4 Proofs of Lemmas, Propositions and Theorems in the Paper

We prove the results in the paper one-by-one. For simplicity, we assume the largest signal strength is at most constant level, which implies $\kappa m \lesssim 1$.

6.5 Proof of Lemma 2.1

Let $\beta = g \circ l$. Since the least square problem $\min_{\beta \in \mathbb{R}^p} (2n)^{-1}\|X\beta - y\|^2$ is a convex optimization, it does not have local maximum and any local minimum is a global minimum.

Now suppose $f(g, l)$ has a local maximum (\tilde{g}, \tilde{l}) , meaning that there exists an open set $\mathcal{B} \in \mathbb{R}^{2p}$ centering at (\tilde{g}, \tilde{l}) such for any $(g, l) \in \mathcal{B}$, $f(g, l) \leq f(\tilde{g}, \tilde{l})$. It is easy to verify that the set $\mathcal{B}_\beta = \{g \circ l : (g, l) \in \mathcal{B}\}$ is also an open set in \mathbb{R}^p centering at $\tilde{\beta} = \tilde{g} \circ \tilde{l}$. Consequently, $\tilde{\beta}$ is also a local maximum of $(2n)^{-1}\|X\beta - y\|^2$, which is a contradiction.

Due to the same argument, any local minimum (\bar{g}, \bar{l}) of $f(g, l)$ corresponds to a local minimum $\bar{\beta} = \bar{g} \circ \bar{l}$ of $(2n)^{-1}\|X\beta - y\|^2$. Since all local minimums of $(2n)^{-1}\|X\beta - y\|^2$ are global and satisfies $X^T(X\bar{\beta} - y) = 0$, (\bar{g}, \bar{l}) must be a global minimum of $f(g, l)$ and satisfies $X^T(X(\bar{g} \circ \bar{l}) - y) = 0$.

Now let us prove the last part concerning saddle points. According to the previous paragraph and the discussion after lemma 2.1, a point (g^\dagger, l^\dagger) is a saddle point (that is, neither local minimum nor local maximum) if and only if there exists some $j_0 \in \{1, \dots, p\}$ such that

$$g_{j_0}^\dagger = l_{j_0}^\dagger = 0 \quad \text{and} \quad R_{j_0}^\dagger \neq 0,$$

where $R^\dagger = X^T(X(g^\dagger \circ l^\dagger) - y) \in \mathbb{R}^p$ (since otherwise this saddle point would satisfy $X^T(X(\bar{g} \circ \bar{l}) - y) = 0$, which is the sufficient and necessary condition for a global minimum). For any vector $u \in \mathbb{R}^p$, we use $\text{Diag}(u)$ to denote the p -by- p diagonal matrix whose diagonals are components of u . By direct calculations, we can express the $2p$ -by- $2p$ Hessian matrix $\nabla_{g,l}^2 f(g^\dagger, l^\dagger)$ of $f(g, l)$ at (g^\dagger, l^\dagger) as

$$\begin{pmatrix} \text{Diag}(l^\dagger) & 0 \\ 0 & \text{Diag}(g^\dagger) \end{pmatrix} \begin{pmatrix} n^{-1}X^TX & n^{-1}X^TX \\ n^{-1}X^TX & n^{-1}X^TX \end{pmatrix} \begin{pmatrix} \text{Diag}(l^\dagger) & 0 \\ 0 & \text{Diag}(g^\dagger) \end{pmatrix} + \begin{pmatrix} 0 & \text{Diag}(R^\dagger) \\ \text{Diag}(R^\dagger) & 0 \end{pmatrix}.$$

Consequently, it is easy to verify that

$$\begin{pmatrix} e_{j_0}^T & -\text{sgn}(R_{j_0}^\dagger) e_{j_0}^T \end{pmatrix} \nabla_{g,l}^2 f(g^\dagger, l^\dagger) \begin{pmatrix} e_{j_0} \\ -\text{sgn}(R_{j_0}^\dagger) e_{j_0} \end{pmatrix} = -2(R_{j_0}^\dagger)^2 < 0,$$

where e_j denotes the vector whose j -th component is one and zero elsewhere for $j = 1, \dots, p$, and $\text{sgn}(\cdot)$ is the sign function. Consequently, there exists some nonzero vector u such that the evaluation of the quadratic form $u^T \nabla_{g,l}^2 f(g^\dagger, l^\dagger) u < 0$, implying $\lambda_{\min}(\nabla_{g,l}^2 f(g^\dagger, l^\dagger)) < 0$.

6.6 Proof of Lemma 2.2

According to lemma 2.1, all saddle points of $f(g, l)$ are strict, and all local minimums are global. Therefore, we can apply Theorem 4 in Lee et al. (2016) to finish the proof, which states that for any twice continuously differentiable function, if it has only strict saddle points, then gradient descent with a random initialization and sufficiently small constant step size almost surely converges to a local minimizer.

6.7 Proof of Proposition 3.1

Recall the decomposition of u_t as

$$\begin{aligned} u_t &= z_t + d_t + e_t, \quad \text{with } z_t := I_{S_1} u_t \in \mathbb{R}^p \\ \text{and } d_t &:= I_{S_2} u_t \in \mathbb{R}^p, \quad e_t := (I - I_S) u_t \in \mathbb{R}^p. \end{aligned}$$

Moreover, we have the updating rule $u_{t+1} = u_t - 2\eta n^{-1} X^T (X u_t^2 - y) \circ u_t$. Therefore, the corresponding updating rules for strong signal component z_t , weak signal component d_t and error component e_t are

$$\begin{aligned} z_{t+1} &= z_t - 2\eta n^{-1} X^T (X u_t^2 - y) \circ z_t, \\ d_{t+1} &= d_t - 2\eta n^{-1} X^T (X u_t^2 - y) \circ d_t, \\ e_{t+1} &= e_t - 2\eta n^{-1} X^T (X u_t^2 - y) \circ e_t. \end{aligned}$$

Also recall the following conditions on the step size η , and RIP constant δ : $\eta \lesssim 1/(\kappa \log(p/\alpha))$, and $\delta \lesssim 1/(\kappa s^{1/2} \log(p/\alpha))$. All the following analysis is based on the conditional events:

$$\frac{1}{n} \|X_{S_1}^\top w\| \leq c\sigma \sqrt{\frac{M s_1}{n}} \quad \text{and} \quad \frac{1}{n} \|X^\top w\|_\infty \leq \sigma \sqrt{\frac{4 \log p}{n}}, \quad (16)$$

where the first concentration is based on lemma 6.4 and the RIP of $X^\top X/n$. We will use induction to prove that for all time $t < T_1 = \Theta(\frac{\log(p/\alpha)}{\eta m})$ in the first stage,

$$\|e_t\|_\infty \lesssim 1/p, \quad \|d_t\|_\infty \lesssim 1/p, \quad (17)$$

$$\|z_t\|_\infty \lesssim 1, \quad (18)$$

$$\|z_t^2 - I_{S_1} \beta^*\| \lesssim \sqrt{s}, \quad (19)$$

where recall that m is the smallest absolute value of components of $I_{S_1} \beta^*$, s_1 is the sparsity of $I_{S_1} \beta^*$, and κ is the ratio of largest and smallest of $I_{S_1} \beta^*$, $\kappa m \lesssim 1$ for simplicity. When

$t = 0$, we have $u_0 = \alpha \mathbf{1}$, $z_0 = \alpha I_S \mathbf{1}$, and $e_0 = \alpha I_{S^c} \mathbf{1}$. Therefore, by choosing $\alpha \leq c/p$ for some sufficiently small constant $c > 0$, we can ensure $\|e_0\|_\infty \lesssim 1/p$, and $\|z_0\|_\infty \lesssim 1$. Now suppose for time $t < T_1 = \Theta(\frac{\log(p/\alpha)}{\eta m})$, we have (17)-(19). We will show that these inequalities also hold for time $t + 1$.

Our strategy is to show both signal part $\|z_t\|_\infty$, $\|d_t\|_\infty$ and error part $\|e_t\|_\infty$ increases in an exponential rate when $\theta_{s_1}(z_t) \leq \frac{\sqrt{m}}{2}$. In addition, if $t + 1 < T_1$, then we show that induction hypothesis equations (17)–(19) still hold. In particular, we divide the proof into 3 steps:

Step 1: Analyzing Error Dynamics and Weak Signals

The proposition below gives us the dynamics of the error terms, we can use this dynamics to see the infinity norm of the error terms remains in the order of $O(1/p)$ in the first stage when $t < T_1$.

Proposition 6.1. *Under assumptions in theorem 3.1 and the induction hypothesis (17)–(19) at t , when $t < T_1$ we have:*

$$\begin{aligned}\|e_{t+1}\|_\infty &\leq (1 + c\eta m / \log(p/\alpha))\|e_t\|_\infty, \\ \|d_{t+1}\|_\infty &\leq (1 + c\eta m / \log(p/\alpha))\|d_t\|_\infty,\end{aligned}$$

when $\|z_t^2 - \beta^*\| \leq \sqrt{\epsilon}$, then let $\tau = \max\{\delta\alpha, \sigma\sqrt{\frac{\log p}{n}}\}$:

$$\begin{aligned}\|e_{t+1}\|_\infty &\leq (1 + c\eta\tau)\|e_t\|_\infty. \\ \|d_{t+1}\|_\infty &\leq (1 + c\eta\tau)\|d_t\|_\infty.\end{aligned}$$

Proof. For $\|e_{t+1}\|_\infty$, we bound the term $n^{-1}X^T(Xu_t^2 - y) \circ e_t$ in the updating formula:

$$\begin{aligned}
& \|n^{-1}X^T(Xu_t^2 - y) \circ e_t\|_\infty = \|e_t \circ (n^{-1}X^T(X(u_t^2 - \beta^*) - w))\|_\infty \\
& \stackrel{(i)}{=} \|e_t \circ ((n^{-1}X^T X(z_t^2 + d_t^2 - \beta^*)) - (z_t^2 + d_t^2 - \beta^*) + n^{-1}X^T X(e_t^2) - n^{-1}X^T w)\|_\infty \\
& \stackrel{(ii)}{\leq} \|e_t\|_\infty \left(\delta \|z_t^2 + d_t^2 - \beta^*\| + \|(e_t^2)\|_\infty + \delta \|(e_t^2)\|_1 + c\sigma \sqrt{\frac{\log p}{n}} \right) \\
& \stackrel{(iii)}{\leq} \|e_t\|_\infty (c \frac{m}{\log(p/\alpha)}),
\end{aligned}$$

where in step (i) we use the fact that $e_t \circ (z_t^2 + d_t^2 - \beta^*) = 0$; step (ii) follows by lemma 6.2:

$$\|e_t \circ ((n^{-1}X^T X(z_t^2 + d_t^2 - \beta^*)) - (z_t^2 + d_t^2 - \beta^*))\|_\infty \leq \|e_t\|_\infty (\delta \|z_t^2 + d_t^2 - \beta^*\|),$$

lemma 6.3:

$$\|e_t \circ (n^{-1}X^T X(e_t^2) - e_t^2 + e_t^2)\|_\infty \leq \|e_t\|_\infty (\delta \|e_t^2\|_1 + \|(e_t^2)\|_\infty),$$

and lemma 6.5:

$$\|e_t \circ (n^{-1}X^T w)\|_\infty \leq \|e_t\|_\infty c\sigma \sqrt{\frac{\log p}{n}}.$$

In step (iii), we applied the induction hypothesis (17). Finally, note that $\delta\sqrt{s}\kappa m \lesssim \frac{m}{\log(p/\alpha)}$ and $\sigma\sqrt{\frac{\log p}{n}} \lesssim m/\log(p/\alpha)$ by the minimal strength of m , then an application of the triangle inequality leads to the claimed bound on $\|e_{t+1}\|_\infty$.

Similarly, for $\|d_{t+1}\|_\infty$,

$$\begin{aligned}
& \|n^{-1}X^T(Xu_t^2 - y) \circ d_t\|_\infty = \|d_t \circ (n^{-1}X^T(X(u_t^2 - \beta^*) - w))\|_\infty \\
& \stackrel{(i)}{=} \|d_t \circ ((n^{-1}X^T X(z_t^2 + d_t^2 - \beta^*)) - (z_t^2 + d_t^2 - \beta^*) + (d_t^2 - I_{S_2}\beta^*) + n^{-1}X^T X(e_t^2) - n^{-1}X^T w)\|_\infty \\
& \stackrel{(ii)}{\leq} \|d_t\|_\infty \left(\delta\|z_t^2 + d_t^2 - \beta^*\| + \|d_t^2 - I_{S_2}\beta^*\|_\infty + \delta\|(e_t^2)\|_1 + c\sigma\sqrt{\frac{\log p}{n}} \right) \\
& \leq \|d_t\|_\infty (cm/\log(p/\alpha)),
\end{aligned}$$

where step (i) follows by $d_t \circ (z_t^2 - I_{S_1}\beta^*) = 0$; step (ii) follows from the orthogonality between d_t and e_t . The last inequality follows from the induction hypothesis.

Then, when $\|z_t^2 + d_t^2 - \beta^*\| \leq \sqrt{\epsilon}$:

$$\begin{aligned}
\|n^{-1}X^T(Xu_t^2 - y) \circ e_t\|_\infty & \leq \|e_t\|_\infty \left(\delta\|z_t^2 + d_t^2 - \beta^*\| + \|(e_t^2)\|_\infty + \delta\|(e_t^2)\|_1 + c\sigma\sqrt{\frac{\log p}{n}} \right) \\
& \leq c\tau\|e_t\|_\infty,
\end{aligned}$$

where the last inequality follows from $\delta \cdot \max\{\alpha, \sigma\sqrt{\frac{s \log p}{n}}\} \leq \max\{\delta\alpha, \sigma\sqrt{\frac{\log p}{n}}\}$ and $\delta p\alpha^2 \lesssim \delta\alpha$.

$$\begin{aligned}
\|n^{-1}X^T(Xu_t^2 - y) \circ d_t\|_\infty & \leq \|d_t\|_\infty \left(\delta\|z_t^2 + d_t^2 - \beta^*\| + \|d_t^2 - I_{S_2}\beta^*\|_\infty + \delta\|(e_t^2)\|_1 + \sigma\sqrt{\frac{\log p}{n}} \right) \\
& \leq c\tau\|d_t\|_\infty,
\end{aligned}$$

which is also based on $\|d_t^2 - I_{S_2}\beta^*\|_\infty \leq c \max\{\alpha^2, \sigma\sqrt{\frac{\log p}{n}}\}$, since $\|d_t\| \leq c\alpha$ and the signals in set S_2 are weak. \square

Step 2: Analyzing Strong Signal Dynamics

The proposition below tells us the first-stage convergence of the strong signal terms, which are much faster than the growing of the weak signal and error terms, under the condition

on δ in Theorem 3.1.

Proposition 6.2. *Under assumptions in theorem 3.1 and the induction hypothesis (17)-(19) at t , when $\theta_{s_1}(z_t) \leq \frac{\sqrt{m}}{2}$ and $t < T_1$, we have*

$$\theta_{s_1}(z_{t+1}) \geq (1 + \frac{\eta m}{4})\theta_{s_1}(z_t)$$

Proof. First, applying lemma 6.2 with $u = z_t^2 - I_{S_1}\beta^*$ (which is s_1 -sparse) and $v = \mathbf{1} \in \mathbb{R}^p$, we have

$$\begin{aligned} & \|n^{-1}X^T(Xu_t^2 - y) - (z_t^2 - I_{S_1}\beta^*)\|_\infty \\ &= \|n^{-1}X^T X(z_t^2 - I_{S_1}\beta^*) - (z_t^2 - I_{S_1}\beta^*) + n^{-1}X^T X(d_t^2 + e_t^2 - I_{S_2}\beta^*) - n^{-1}X^T w\|_\infty \\ &\leq \delta\|z_t^2 - I_{S_1}\beta^*\| + \|n^{-1}X^T X(d_t^2 - I_{S_2}\beta^*)\|_\infty + \|n^{-1}X^T X e_t^2\|_\infty + \|n^{-1}X^T w\|_\infty \\ &\stackrel{(i)}{\lesssim} m/\log(p/\alpha), \end{aligned} \tag{20}$$

where step (i) follows from induction hypothesis (17) and (19), and our condition on δ and m . Let z_{t+1}^+ be the element-wise pseudo-inverse (meaning $1/0 = 0$) of z_{t+1} .

First, we show $\|z_{t+1}^+ \circ z_t\|_\infty \lesssim 1$. By rewriting the updating rule of z_t as

$$z_{t+1} = z_t \circ (\mathbf{1} - 2\eta n^{-1}X^T(X(u_t^2 - \beta^*) - w)), \tag{21}$$

we have the element-wise lower bound

$$z_{t+1} \geq z_t \circ (\mathbf{1} - 2\eta(z_t^2 - I_{S_1}\beta^*)) - c\eta \frac{m}{\log(p/\alpha)}|z_t|,$$

where we applied (20). Since $\eta \leq 1/4$ and $m \leq 1$ under the assumptions of the proposition, the preceding display implies that each component of z_{t+1} will larger than $1/4$ times the corresponding component of z_t in absolute value, which implies $\|z_t \circ z_{t+1}^+\|_\infty \leq 4 \lesssim 1$.

Using the updating rule (21) for z_t , equation (20), and the induction hypothesis (18),

we have:

$$\|z_{t+1} - z_t \circ (\mathbf{1} - 2\eta(z_t^2 - I_{S_1}\beta^*))\|_\infty \lesssim \eta \frac{m}{\log(p/\alpha)} |z_t|. \quad (22)$$

Since $\|z_t \circ z_{t+1}^+\|_\infty \leq 4$, by conducting Hadamard product with z_{t+1}^+ on both sides, we obtain:

$$\begin{aligned} \|\mathbf{1} - (\mathbf{1} - 2\eta(z_t^2 - I_{S_1}\beta^*)) \circ z_{t+1}^+ \circ z_t\|_\infty &\lesssim \eta \frac{m}{\log(p/\alpha)}. \\ \|\mathbf{1} - (\mathbf{1} - 2\eta z_t^2)(\mathbf{1} + 2\eta I_{S_1}\beta^*) \circ z_{t+1}^+ \circ z_t\|_\infty &\lesssim \eta \frac{m}{\log(p/\alpha)} + \eta^2 \kappa^2 m^2. \end{aligned}$$

We write

$$z_{t+1} - z_t \circ (\mathbf{1} - 2\eta z_t^2)(\mathbf{1} + 2\eta I_{S_1}\beta^*) = \xi_t z_{t+1},$$

where $\xi_t = \mathbf{1} - (\mathbf{1} - 2\eta z_t^2)(\mathbf{1} + 2\eta I_{S_1}\beta^*) \circ z_{t+1}^+ \circ z_t$, where we drop the η^2 term since $\eta\kappa \lesssim \frac{1}{\log(p/\alpha)}$ and $\kappa m \lesssim 1$, using the display after equation (22), we have the following bound on ξ_t ,

$$\|\xi_t\|_\infty \lesssim \eta \frac{m}{\log(p/\alpha)}.$$

Consider the function $f(x) = x(1 - 2\eta x^2)$, which is nondecreasing for $x \in [-(6\eta)^{-1}, (6\eta)^{-1}]$. Therefore, when $\|z_t\|_\infty \lesssim 1$ and $\eta \leq 1/6$ (by choosing the constant c_3 in Theorem 3.1 in the paper sufficiently small), we must have

$$\theta_{s_1}(z_t \circ (\mathbf{1} - 2\eta(z_t^2))) = \theta_{s_1}(z_t)(1 - 2\eta\theta_{s_1}(z_t)^2). \quad (23)$$

The last three displays together imply that when $\theta_{s_1}(z_t) \leq \frac{\sqrt{m}}{2}$,

$$\begin{aligned} \theta_{s_1}(z_{t+1}) &\stackrel{(i)}{\geq} \frac{(1 + 2\eta m)\theta_{s_1}(z_t)(1 - 2\eta\theta_{s_1}(z_t)^2)}{1 + c\eta \frac{m}{\log(p/\alpha)}} \\ &\geq (1 + 2\eta m)(1 - 2\eta\theta_{s_1}(z_t)^2)\theta_{s_1}(z_t)(1 - c\eta \frac{m}{\log(p/\alpha)}) \\ &\stackrel{(ii)}{\geq} (1 + \frac{\eta m}{4})\theta_{s_1}(z_t), \end{aligned}$$

where in step (i) we used the definition of m as $\theta_{s_1}(\beta^*)$, and step (ii) follows from the assumption $\theta_{s_1}(z_t) \leq \frac{\sqrt{m}}{2}$ and our conditions on (η, m) , so that

$$\begin{aligned} (1 + 2\eta m)(1 - \frac{\eta m}{2}) &= (1 + \frac{3\eta m}{2} - \eta^2 m^2) \stackrel{(iii)}{\geq} 1 + \eta m, \quad \text{and} \\ (1 + \eta m)(1 - c\eta \frac{m}{\log(p/\alpha)}) &= (1 + \eta m - c\eta \frac{m}{\log(p/\alpha)} - c\eta^2 \frac{m^2}{\log(p/\alpha)}) \\ &\stackrel{(iv)}{\geq} (1 + \frac{\eta m}{2} - \frac{\eta^2 m^2}{2}) \stackrel{(v)}{\geq} (1 + \frac{\eta m}{4}), \end{aligned}$$

where step (iii) follows from $\eta m \leq 1/2$, step (iv) by $\frac{c}{\log(p/\alpha)} \leq 1/2$, and in step (v) we used $\eta m \leq 1/2$ and $\frac{c}{\log(p/\alpha)} \leq 1/2$ again. □

Step 3: Prove Induction Hypothesis

Our last piece shows that the induction hypothesis is kept for $t + 1$. This proposition combined with propositions 6.1-6.2 leads to a proof for proposition 3.1 in the paper.

Proposition 6.3. *Under assumptions in theorem 3.1 and the induction hypothesis (17)-(19) at $t < T_1$, we have:*

$$\|e_{t+1}\|_\infty \lesssim 1/p, \quad \|d_{t+1}\|_\infty \lesssim 1/p, \quad \|z_{t+1}\|_\infty \lesssim 1, \quad \text{and} \quad \|z_{t+1}^2 - I_{S_1}\beta^*\| \lesssim \sqrt{s}.$$

When $T_1 < t < T_2$, we also have:

$$\|e_{t+1}\|_\infty \lesssim 1/p, \quad \|d_{t+1}\|_\infty \lesssim 1/p.$$

Proof. For error component e_{t+1} , we have:

$$\begin{aligned} \|e_{t+1}\|_\infty &\leq (1 + c\eta m / \log(p/\alpha)) \|e_t\|_\infty \stackrel{(i)}{\leq} (1 + c/T_1)^{T_1} \|e_0\|_\infty \\ &\leq e^c \|e_0\|_\infty \lesssim \alpha \lesssim 1/p, \end{aligned}$$

where in step (i) we apply $\frac{m}{\log(p/\alpha)} \lesssim 1/(\eta T_1)$. When $\|z_t^2 - \beta^*\| \leq \sqrt{\epsilon}$, since $T_2 = c/(\eta\tau)$, when $T_1 < t < T_2$:

$$\begin{aligned}\|e_{t+1}\|_\infty &\leq (1 + c\eta\tau)\|e_t\|_\infty \stackrel{(i)}{\leq} (1 + c/T_2)^{T_2} \cdot \alpha \\ &\leq e^c \cdot \alpha \lesssim \alpha \lesssim 1/p.\end{aligned}$$

Similarly we can show the induction hypothesis for weak signals $\|d_{t+1}\|$. For strong signal component z_{t+1} , we apply equation (21) to get element-wise upper bound:

$$z_{t+1,i} \leq z_{t,i}(1 - 2\eta(z_{t,i}^2 - \beta_i^*)) + c\eta\delta\sqrt{s}z_{t,i},$$

Since $\delta\sqrt{s} \lesssim 1$, and the function $f(x) = 2\eta x^3 - c\eta x$ is nonnegative when $x \gtrsim 1$. Therefore, as long as $\|z_t\|_\infty \gtrsim 1$, we always have $\|z_{t+1}\|_\infty \leq \|z_t\|_\infty$, thus

$$\|z_{t+1}^2 - \beta^*\| \leq \|z_{t+1}^2\| + \|\beta^*\| \lesssim \sqrt{s}.$$

□

6.8 Proof of Proposition 3.2

Proof. We decompose $\|u_{t+1}^2 - \beta^*\|^2 = \|(u_t - \eta\nabla f(u_t)) \circ (u_t - \eta\nabla f(u_t)) - \beta^*\|^2$ into 5 parts, according to the order of η :

$$\begin{aligned}\|u_{t+1}^2 - \beta^*\|^2 &= \|u_t^2 - \beta^*\|^2 - 4\eta\langle \nabla f(u_t) \circ u_t, u_t^2 - \beta^* \rangle \\ &\quad + 4\eta^2\langle \nabla f(u_t)^2, u_t^2 - \beta^* \rangle + 2\eta^2\|f(u_t) \circ u_t\|^2 \\ &\quad - 4\eta^3\langle \nabla f(u_t) \circ u_t, \nabla f(u_t)^2 \rangle + \eta^4\|\nabla f(u_t)^2\|^2.\end{aligned}$$

In addition, we have the decomposition,

$$\begin{aligned}\|u_{t+1}^2 - \beta^*\|^2 &= \|z_{t+1}^2 - I_{S_1}\beta^*\|^2 + \|d_{t+1}^2 - I_{S_2}\beta^*\|^2 + \|e_t^2\|^2 \\ &\leq \|z_{t+1}^2 - I_{S_1}\beta^*\|^2 + cs_2 \frac{\log p}{n} + c'p\alpha^4,\end{aligned}\quad (24)$$

where the rate are based on the dynamics for weak signals and error terms, so we only need to consider the estimation errors on the index set S_1 .

Before we bound each term, we need the following bound for the gradient:

$$\begin{aligned}\|\nabla f(u_t)\|_\infty &= \|2u_t \circ n^{-1}X^T(Xu_t^2 - y)\|_\infty \\ &\leq c(\|n^{-1}X^TX(u_t^2 - \beta^*)\|_\infty + \|X^Tw\|_\infty) \\ &= c(\|n^{-1}X^TX(z_t^2 + d_t^2 - \beta^* + e_t^2)\|_\infty + \|X^Tw\|_\infty) \\ &\stackrel{(i)}{\leq} c(\|z_t^2 + d_t^2 - \beta^*\|_\infty + \delta\|z_t^2 + d_t^2 - \beta^*\| + \|n^{-1}X^TX(e_t)^2\|_\infty + \|X^Tw\|_\infty) \\ &\stackrel{(ii)}{\lesssim} 1,\end{aligned}\quad (25)$$

where (i) follows by lemma 6.2 and lemma 6.3 with $v = \mathbf{1} \in \mathbb{R}^p$, and (ii) follows by the assumption that $(z_t, d_t, e_t) \in \Theta_{LR}$, and

$$\begin{aligned}\|n^{-1}X^TX(e_t)^2\|_\infty &\leq \|e_t^2\|_\infty + \delta\|e_t^2\|_1 \leq c\delta p\alpha^2 \\ \|n^{-1}(X^TX(d_t^2 - I_{S_2}\beta^*))\|_\infty &\leq \|d_t^2 - I_{S_2}\beta^*\|_\infty + \delta\|d_t^2 - I_{S_2}\beta^*\|_2 \leq c\alpha^2 + c'\sigma\sqrt{\frac{\log p}{n}} \\ \|n^{-1}(X^TX(d_t^2 - I_{S_2}\beta^* + e_t^2) - X^Tw)\|_\infty &\leq c\delta p\alpha^2 + c'\sigma\sqrt{\frac{\log p}{n}}.\end{aligned}\quad (26)$$

Note that z_t is $s_1 \leq s$ sparse, we have the core error term characterization:

$$\begin{aligned}\|z_t \circ n^{-1}(X^TX(d_t^2 - I_{S_2}\beta^* + e_t^2) - X^Tw)\|^2 &\leq c(Ms_1)\delta^2 p^2 \alpha^4 + c'\sigma^2 \frac{Ms_1}{n} \\ &\leq c(\alpha^2 + \sigma^2 \frac{Ms_1}{n}),\end{aligned}\quad (27)$$

where we used $\delta^2 s \leq 1$, $\|e_t\|_\infty \lesssim \alpha \lesssim 1/p$ and lemma 6.3 with $v = \mathbf{1} \in \mathbb{R}^p$ and lemma 6.4.

In the rest of the proof, we will repeatedly apply lemma 6.2, 6.3 to get rid of the $n^{-1}X^T X$ term; we will also repeatedly use lemma 6.1 and the assumption $(z_t, d_t, e_t) \in \Theta_{LR}$ to bound each term in the decomposition of $\|u_{t+1}^2 - \beta^*\|^2$. For quantifying the errors, we will use equation (26) and (27). Note that based on the decomposition (24), we only need to care about terms contributed by the index set S_1 .

First, we need a lower bound for the order one term,

$$\begin{aligned}
& \langle \nabla f(u_t) \circ u_t, I_{S_1}(u_t^2 - \beta^*) \rangle = \langle \nabla f(u_t) \circ z_t, z_t^2 - I_{S_1}\beta^* \rangle \\
& \geq \langle n^{-1}X^T X(z_t^2 - I_{S_1}\beta^*) \circ z_t^2, z_t^2 - I_{S_1}\beta^* \rangle \\
& \quad - \|z_t \circ n^{-1}(X^T X(d_t^2 - I_{S_2}\beta^* + e_t^2) - X^T w)\| \|z_t \circ (z_t^2 - I_{S_1}\beta^*)\| \\
& \geq \|z_t \circ (z_t^2 - I_{S_1}\beta^*)\|^2 - \delta \|z_t^2 - I_{S_1}\beta^*\|^2 \\
& \quad - \|z_t \circ n^{-1}(X^T X(d_t^2 - I_{S_2}\beta^* + e_t^2) - X^T w)\|^2/2 - \|z_t \circ (z_t^2 - I_{S_1}\beta^*)\|^2/2, \\
& \geq \frac{1}{2} \|z_t \circ (z_t^2 - I_{S_1}\beta^*)\|^2 - \delta \|z_t^2 - I_{S_1}\beta^*\|^2 - c(\alpha^2 + \sigma^2 \frac{Ms_1}{n}),
\end{aligned}$$

where the second inequality is based on Cauchy-Schwarz inequality and the last step follows by equation (27). The rest of the terms are also based on the bound:

$$\begin{aligned}
\|I_{S_1} \nabla f(u_t)\|^2 &= \|z_t \circ n^{-1}X^T(Xu_t^2 - y)\|^2 \\
&\stackrel{(i)}{\leq} 2\|z_t \circ n^{-1}X^T X(z_t^2 - I_{S_1}\beta^*)\|^2 + 2\|z_t \circ n^{-1}(X^T X(d_t^2 - I_{S_2}\beta^* + e_t^2) - X^T w)\|^2 \\
&\stackrel{(ii)}{\leq} c\|z_t^2 - I_{S_1}\beta^*\|^2 + c'(\alpha^2 + \sigma^2 \frac{Ms_1}{n}),
\end{aligned}$$

where in step (i) we use induction hypothesis and Cauchy-Schwarz inequality, and step (ii) we use the bound for equation (27).

For the order two term, we apply a similar idea to bound $\|\langle I_{S_1} \nabla f(u_t)^2, u_t^2 - \beta^* \rangle\|$ and

$\|I_{S_1}f(u_t)^2\|^2$. More precisely, since $\|u_t^2 - \beta^*\|_\infty \lesssim 1$, we have

$$\begin{aligned} \langle \nabla f(u_t)^2, I_{S_1}(u_t^2 - \beta^*) \rangle &= \langle \nabla f(u_t) \circ (u_t^2 - \beta^*), \nabla I_{S_1}f(u_t) \rangle \\ &\leq c\|I_{S_1}\nabla f(u_t)\|^2 \leq c\|z_t^2 - I_{S_1}\beta^*\|^2 + c'(\alpha^2 + \sigma^2 \frac{Ms_1}{n}). \end{aligned}$$

Similarly $\|u_t\|_\infty \lesssim 1$, we bound $\|I_{S_1}\nabla f(u_t) \circ u_t\|^2$ through

$$\|I_{S_1}\nabla f(u_t) \circ u_t\|^2 \leq c\|I_{S_1}\nabla f(u_t)\|^2 \leq c\|z_t^2 - I_{S_1}\beta^*\|^2 + c'(\alpha^2 + \sigma^2 \frac{Ms_1}{n}).$$

For the order three term, using the fact that $\|u_t\|_\infty \lesssim 1$ and $\|\nabla f(u_t)\|_\infty \lesssim 1$, we have:

$$\begin{aligned} \langle I_{S_1}\nabla f(u_t) \circ u_t, \nabla f(u_t)^2 \rangle &\leq c\|\nabla f(u_t)\|^2 \leq c\|I_{S_1}\nabla f(u_t)\|^2 \\ &\leq c\|z_t^2 - I_{S_1}\beta^*\|^2 + c'(\alpha^2 + \sigma^2 \frac{Ms_1}{n}). \end{aligned}$$

For the order four term, we have

$$\|I_{S_1}\nabla f(u_t)^2\|^2 \lesssim \|I_{S_1}\nabla f(u_t)\|^2 \lesssim c\|z_t^2 - I_{S_1}\beta^*\|^2 + c'(\alpha^2 + \sigma^2 \frac{Ms_1}{n}).$$

Putting all pieces together, since $\eta \lesssim 1$ under our assumption on η , we get

$$\begin{aligned} \|z_{t+1}^2 - I_{S_1}\beta^*\|^2 &\leq \|z_t^2 - I_{S_1}\beta^*\|^2 - 2\eta\|(z_t^2 - I_{S_1}\beta^*) \circ z_t\|^2 \\ &\quad + c\eta(\alpha^2 + \sigma^2 \frac{Ms_1}{n}) + 4\delta\eta\|z_t^2 - I_{S_1}\beta^*\|^2, \end{aligned}$$

with $\|z_t^2 - I_{S_1}\beta^*\|^2 \cdot m/4 \leq \|(z_t^2 - I_{S_1}\beta^*) \circ z_t\|^2$. Since $\theta_s(z_t) \geq \sqrt{m}/2$ for $(z_t, d_t, e_t) \in \Theta_{LR}$, and $\delta\kappa \leq \frac{1}{4}$,

under our assumptions on (η, δ) , we reach

$$\|z_{t+1}^2 - I_{S_1}\beta^*\|^2 \leq (1 - \eta m)\|z_t^2 - I_{S_1}\beta^*\|^2 + c(\alpha^2 + \sigma^2 \frac{Ms_1}{n}).$$

□

6.9 Proof of Theorem 3.1

Proof. As mentioned after the theorem in the paper, we divide the convergence into two stages. In the first “burn-in” stage, we show that the smallest component of the strong signal part z_t increases at an exponential rate in t until hitting $\sqrt{m}/2$,

$$\theta_{s_1}(z_{t+1}) \geq (1 + \frac{\eta m}{4})\theta_s(z_t), \text{ when } \theta_{s_1}(z_t) \leq \sqrt{m}/2.$$

In the second stage, after iterate u_t enters Θ_{LR} , we have the geometric convergence up to some high-order error term

$$\|z_{t+1}^2 - I_{S_1}\beta^*\|^2 \leq (1 - \eta m)\|z_t^2 - I_{S_1}\beta^*\|^2 + c(\alpha^2 + \sigma^2 \frac{Ms_1}{n}), \text{ when } \theta_{s_1}(z_t) \geq \sqrt{m}/2,$$

When $t \leq \Theta(\log \frac{\sqrt{m}}{\alpha}/(\eta m))$, the convergence is in the first stage; and proposition 6.1 implies that the iterate u_t enters Θ_{LR} (corresponding to the second stage) in at most $\Theta(\log \frac{c}{p\alpha^2}/(\eta m))$ iterations. proposition 6.1 and proposition 3.2 together imply for any $t = \Theta(\log \frac{p}{\alpha}/(\eta m))$,

$$\|z_t^2 - I_{S_1}\beta^*\|^2 \leq c(\alpha^2 + \sigma^2 \frac{Ms_1}{n}). \quad (28)$$

For any $T_2 > t > T_1$, since $\|e_t\|, \|d_t\| \lesssim 1/p$ is still controlled, combined with bound (24), we have:

$$\|u_t^2 - \beta^*\|^2 \leq c(\alpha^2 + \sigma^2 \frac{Ms_1}{n} + \sigma^2 \frac{s_2 \log p}{n}).$$

□

6.10 Proof of Proposition 3.3

Similar to the proof for the nonnegative case, we use induction to show that for each $t \leq T_1$,

$$\|a_{t,S_1^c}\|_\infty \lesssim 1/p, \quad \|b_{t,S_1^c}\|_\infty \lesssim 1/p, \quad (29)$$

$$\|a_{t,S_1}\|_\infty \lesssim 1, \quad \|b_{t,S_1}\|_\infty \lesssim 1, \quad (30)$$

$$\|\beta_{t,S_1} - \beta_{S_1}^*\| \lesssim \sqrt{s}, \quad (31)$$

where the set S_1^c is the union of weak signals and errors. When $t = 0$, we have $g_0 = \alpha \mathbf{1}$, $l_0 = 0$. Therefore, under the assumption $\alpha \lesssim 1/p$, we have $\|a_{0,S_1^c}\|_\infty \lesssim 1/p$, $\|a_{0,S_1}\|_\infty \lesssim 1$, and similar bounds for b . Now suppose for time $t < T_1$, we have (29)-(31). We divide into 3 steps to show the same bounds hold for time $t + 1$. Note that the following analysis is still based on the conditional events:

$$\frac{1}{n} \|X_{S_1}^\top w\| \leq c\sigma \sqrt{\frac{Ms_1}{n}} \quad \text{and} \quad \frac{1}{n} \|X^\top w\|_\infty \leq \sigma \sqrt{\frac{4 \log p}{n}}, \quad (32)$$

still the first concentration is based on lemma 6.4 and the RIP of $X^\top X/n$.

Step 1: Analyzing Error Dynamics

Proposition 6.4. *Based on RIP and the constant assumption in theorem 3.2 and the induction hypothesis (29)-(31), when $t < T_1$ we have:*

$$\|a_{t+1,S_1^c}\|_\infty \leq (1 + c\eta m / \log(p/\alpha)) \|a_{t+1,S_1^c}\|_\infty,$$

$$\|b_{t+1,S_1^c}\|_\infty \leq (1 + c\eta m / \log(p/\alpha)) \|b_{t+1,S_1^c}\|_\infty.$$

When $\|\beta_{t,S_1} - \beta_{S_1}^*\| \leq c\sqrt{\epsilon}$, let $\tau = \max\{\sqrt{\frac{\log p}{n}}, \delta\alpha\}$, then

$$\|a_{t+1,S_1^c}\|_\infty \leq (1 + c\eta\tau)\|a_{t+1,S_1^c}\|_\infty,$$

$$\|b_{t+1,S_1^c}\|_\infty \leq (1 + c\eta\tau)\|b_{t+1,S_1^c}\|_\infty.$$

Proof. By the definition of a_t and b_t , we have the their updating rules:

$$a_{t+1,S_1^c} = a_{t,S_1^c} - \eta a_{t,S_1^c} \circ n^{-1} X^T (X\beta_t - y),$$

$$b_{t+1,S_1^c} = b_{t,S_1^c} + \eta b_{t,S_1^c} \circ n^{-1} X^T (X\beta_t - y).$$

The proofs of the two claimed bound follows exactly the same lines as in the nonnegative case, since we uses the triangle inequality $\|a_{t+1,S_1^c}\|_\infty \leq \|a_{t,S_1^c}\|_\infty + \|\eta a_{t,S_1^c} \circ n^{-1} X^T (X\beta_t - y)\|_\infty$ at the beginning, which does not change when the sign becomes positive as in the second line of b_{t+1,S_1^c} . Consequently, the dynamics of weak signals and error terms a_{t,S_1^c} and b_{t,S_1^c} will behave like e_t in the nonnegative case—the proof of proposition 6.1 still works when $\|a_{t,S_1^c}\|_\infty, \|b_{t,S_1^c}\|_\infty \lesssim 1/p$. \square

Step 2: Analyzing Signal Dynamics

The proof becomes more complicated for the signal dynamics. Our goal is to show:

- if the true signal β_i^* is positive and $\beta_{t,i} < \beta_i^*/2$, then $a_{t,i}^2$ increases and $b_{t,i}^2$ decreases, both in an exponential rate. Overall, the sign of i th component $\beta_{t,i} = a_{t,i}^2 - b_{t,i}^2$ of β_t tends to grow to positive in an exponential rate;
- if the true signal β_i^* is negative and $\beta_{t,i} > \beta_i^*/2$, then $a_{t,i}^2$ decreases and $b_{t,i}^2$ increases, both in an exponential rate. Overall, the sign of i th component $\beta_{t,i} = a_{t,i}^2 - b_{t,i}^2$ of β_t tends to fall to negative in an exponential rate.

Our analysis will be based on the resemblance between update rules of (a_t, b_t) and u_t in the nonnegative case. Recall that we assumed that the RIP constant $\delta \lesssim 1/(\kappa\sqrt{s}\log \frac{p}{\alpha})$,

step size $\eta \lesssim 1/(\kappa \log \frac{p}{\alpha})$, and $T_1 = \Theta(\frac{\log(p/\alpha)}{\eta m})$.

Proposition 6.5. *Under assumptions in theorem 3.2 and the induction hypothesis (29)-(31) at $t < T_1$, when $|\beta_{t,i} - \beta_{t,i}^*| \leq \frac{1}{2}|\beta_i^*|$, then for any $i \in S$:*

$$\begin{aligned}\beta_{t,i} &\geq (1 + c\eta\beta_i^*)^t \alpha^2 - c'\alpha^2, \quad \text{if } \beta_i^* > 0; \\ \beta_{t,i} &\leq (1 - c\eta\beta_i^*)^t (-\alpha^2) + c'\alpha^2, \quad \text{if } \beta_i^* < 0.\end{aligned}$$

Proof. Similar with the proof of proposition ??, first we approximate $(n^{-1}X^T Xu) \circ v$ by $u \circ v$ based on the RIP condition via lemmas 6.2 and 6.3,

$$\begin{aligned}&\|n^{-1}X^T(X\beta_t - y) - (I_{S_1}\beta_t - I_{S_1}\beta^*)\|_\infty \\ &\leq \|n^{-1}X^T X(I_{S_1^c}r_t)\|_\infty + \delta\|r_{t,S_1}\| + \|X^T w\|_\infty \lesssim \frac{m}{\log(p/\alpha)},\end{aligned}\tag{33}$$

implying that under the condition $m \lesssim 1$, we have

$$\begin{aligned}&\|n^{-1}X^T X(\beta_t - y)\|_\infty \\ &\leq \|\beta_{t,S_1} - \beta_{S_1}^*\|_\infty + \|n^{-1}X^T(X(I_{S_1^c}r_t) - w)\|_\infty + \delta\|\beta_{t,S_1} - \beta_{S_1}^*\| \lesssim 1,\end{aligned}$$

where the last inequality uses $\|\beta_{t,S_1^c}\|_\infty \lesssim 1/p$ and $\|\beta_{t,S_1} - \beta_{S_1}^*\| \lesssim \sqrt{s}$.

In order to analyze $\beta_{t,S_1} = a_{t,S_1}^2 - b_{t,S_1}^2$, let us focus on a_{t,S_1}^2 and b_{t,S_1}^2 , separately. According to the updating rule of a_{t,S_1} , we have

$$a_{t+1,S_1}^2 = a_{t,S_1}^2 - 2\eta a_{t,S_1}^2 \circ [n^{-1}X^T(X\beta_t - y)]_{S_1} + \eta^2 a_{t,S_1}^2 \circ [n^{-1}X^T(X\beta_t - y)]_{S_1}^2,$$

where recall that for a vector $a \in \mathbb{R}^p$, a_{S_1} denote the sub-vector of a with indices in S_1 .

Applying lemmas 6.2 with $v = \mathbf{1}$, we obtain

$$\|a_{t+1,S_1}^2 - a_{t,S_1}^2 - 2\eta a_{t,S_1}^2(\beta_{t,S_1} - \beta_{S_1}^*)\|_\infty \lesssim \eta \frac{m}{\log(p/\alpha)} + \eta^2 \kappa^2 m^2 \stackrel{(i)}{\lesssim} \eta \frac{m}{\log(p/\alpha)}.$$

where in step (i) we used $\eta\kappa m \lesssim \frac{m}{\log(p/\alpha)}$ and $\kappa m \lesssim 1$. Similar to the nonnegative case, since $\eta m \leq 1/2$, $\frac{m}{\log(p/\alpha)} \leq 1/2$, we have $a_{t,i}^2/a_{t+1,i}^2 \leq 4$ for $i \in S_1$. Therefore, we can obtain an element-wise bound for $\xi_t = (\xi_{t,i})_{i \in S_1}$,

$$\xi_{t,i} := 1 - a_{t,i}^2/a_{t+1,i}^2 \circ (2\eta(\beta_{t,i} - \beta_i^*)),$$

as $\|\xi_t\|_\infty \lesssim \eta \frac{m}{\log(p/\alpha)}$. Equivalently, we can write

$$a_{t+1,i}^2 = a_{t,i}^2(1 - 2\eta(\beta_{t,i} - \beta_i^*)) + \xi_{t,i}a_{t+1,i}^2. \quad (34)$$

Now let us divide into two cases depending on the sign of β_i^* , $i \in S_1$:

Case $\beta_i^ > 0$:* When $\beta_{t,i} - \beta_i^* \leq -\beta_i^*/2$, since $\beta_i^* \geq m$, we have by equation (34),

$$a_{t+1,i}^2 \geq \frac{a_{t,i}^2(1 + \eta\beta_i^*)}{1 + c\eta \frac{m}{\log(p/\alpha)}} \geq a_{t,i}^2(1 + \eta\beta_i^*)(1 - c\eta \frac{\beta_i^*}{\log(p/\alpha)}) \geq a_{t,i}^2(1 + \eta\beta_i^*/4),$$

where the last inequality follows since $1/\log(p/\alpha) \leq 1/2$ and $\eta\beta_i^* \leq 1/2$. Similarly, we can analyze $b_{t,S}^2$ to get

$$b_{t+1,i}^2 \leq \frac{b_{t,i}^2(1 - \eta\beta_i^*)}{1 - c\eta\delta\sqrt{s}} \leq b_{t,i}^2(1 - \eta\beta_i^*)(1 + c\eta\beta_i^*/\log(p/\alpha)) \leq b_{t,i}^2(1 - \eta\beta_i^*/4).$$

Therefore, $a_{t+1,i}^2$ increases at an exponential rate faster than the noise term a_{t+1,S_1^c} while $b_{t+1,i}^2$ decreases to zero at an exponential rate, and when $a_{t+1,i}$ increases to $\beta_i^*/2$, $b_{t+1,i}$ decreases to $O(\alpha^4)$ correspondingly. A combination of these two leads to the first claimed bound for $\beta_i^* > 0$.

Case $\beta_i^ < 0$:* The analysis for the case is similar: when $\beta_{t,i} - \beta_i^* \geq -\beta_i^*/2$, we have:

$$a_{t+1,i}^2 \leq a_{t,i}^2(1 - \eta\beta_i^*/4), \quad \text{and} \quad b_{t+1,i}^2 \geq b_{t,i}^2(1 + \eta\beta_i^*/4),$$

which leads to the second claimed bound for $\beta_i^* < 0$. □

As a consequence of the proof in this step, after at most $T \geq \Theta(\frac{\log(m/\alpha^2)}{\eta m})$ iterations, we are guaranteed to have $|\beta_{T,i}| \geq |\beta_i^*|/2$ with $\text{sign}(\beta_{T,i}) = \text{sign}(\beta_i^*)$ and $\min\{a_{T,i}^2, b_{T,i}^2\} \leq c\alpha^4$.

Step 3: Prove Induction Hypothesis

Proposition 6.6. *Under assumptions in theorem 3.2, the induction hypothesis (29)-(31) at $t < T_1$, we have:*

$$\begin{aligned} \|a_{t+1, S_1^c}\|_\infty &\lesssim 1/p, \quad \|b_{t+1, S_1^c}\|_\infty \lesssim 1/p, \\ \|a_{t+1, S_1}\|_\infty &\lesssim 1, \quad \|b_{t+1, S_1}\|_\infty \lesssim 1, \\ \text{and } \|\beta_{t+1, S_1} - \beta_{S_1}^*\| &\lesssim 1. \end{aligned}$$

Then for $T_1 < t < T_2$, we still have:

$$\|a_{t+1, S_1^c}\|_\infty \lesssim 1/p, \quad \|b_{t+1, S_1^c}\|_\infty \lesssim 1/p.$$

Proof. Our induction hypothesis

$$\begin{aligned} \|a_{t, S_1^c}\|_\infty &\lesssim 1/p, \quad \|b_{t, S_1^c}\|_\infty \lesssim 1/p, \\ \|a_{t, S_1}\|_\infty &\lesssim 1, \quad \|b_{t, S_1}\|_\infty \lesssim 1, \end{aligned}$$

implies $\|\beta_{t, S_1^c}\|_\infty \lesssim 1/p^2$ and $\|\beta_{t, S_1^c} - \beta_{t, S_1^c}^*\| \lesssim \sqrt{s}$. By updating rules:

$$\begin{aligned} a_{t+1} &= a_t - \eta a_t \circ n^{-1} X^T (X \beta_t - y), \\ b_{t+1} &= b_t + \eta b_t \circ n^{-1} X^T (X \beta_t - y), \end{aligned}$$

and our condition $\eta m / \log(p/\alpha) \lesssim 1/T_1$, we can prove along similar lines as in the proof of

?? that

$$\begin{aligned}\|a_{t+1, S_1^c}\|_\infty &\lesssim 1/p, & \|b_{t+1, S_1^c}\|_\infty &\lesssim 1/p, \\ \|a_{t+1, S_1}\|_\infty &\lesssim 1, & \|b_{t+1, S_1}\|_\infty &\lesssim 1.\end{aligned}$$

These inequalities also imply $\|\beta_{t+1, S_1^c}\|_\infty \lesssim 1/p^2$ and $\|\beta_{t+1, S_1} - \beta_{t+1, S_1}^*\| \lesssim \sqrt{s}$. Then similarly for $T_1 < t < T_2$, we can also show:

$$\|a_{t+1, S_1}\|_\infty \lesssim 1, \quad \|b_{t+1, S_1}\|_\infty \lesssim 1.$$

□

6.11 Proof of Proposition 3.4

Proof. The proof is similar to that of proposition 3.2.

In particular, we can decompose $\|\beta_{t+1} - \beta^*\|^2 = \|a_{t+1}^2 - b_{t+1}^2 - \beta^*\|^2$ into 5 parts according to the order of η :

$$\begin{aligned}\|\beta_{t+1} - \beta^*\|^2 &= \|(a_t - \eta a_t \circ n^{-1} X^T(Xr_t - w))^2 - (b_t + \eta b_t \circ n^{-1} X^T(Xr_t - w))^2 - \beta^*\|^2, \\ &= \|\beta_t - 2\eta(a_t^2 + b_t^2) \circ n^{-1} X^T(Xr_t - w) + \eta^2 \beta_t \circ (n^{-1} X^T(Xr_t - w))^2 - \beta^*\|^2, \\ &= \|\beta_t^2 - \beta^*\|^2 - 4\eta \langle (a_t^2 + b_t^2) \circ n^{-1} X^T(Xr_t - w), \beta_t - \beta^* \rangle, \\ &\quad + 4\eta^2 \|(a_t^2 + b_t^2) \circ n^{-1} X^T(Xr_t - w)\|^2 + 2\eta^2 \langle \beta_t \circ (n^{-1} X^T(Xr_t - w))^2, \beta_t - \beta^* \rangle, \\ &\quad - 4\eta^3 \langle \beta_t \circ (n^{-1} X^T(Xr_t - w))^2, (a_t^2 + b_t^2) \circ n^{-1} X^T(Xr_t - w) \rangle, \\ &\quad + \eta^4 \|\beta_t \circ (n^{-1} X^T(Xr_t - w))^2\|^2,\end{aligned}$$

where for notational simplicity, we denoted $\beta_t - \beta^*$ by r_t . First look at the term $a_t^2 + b_t^2$. We have the component-wise bound $a_{t,i}^2 + b_{t,i}^2 = g_{t,i}^2 + l_{t,i}^2 \geq |g_{t,i} \circ l_{t,i}| = |\beta_{t,i}|$, implying

$$a_{t,i}^2 + b_{t,i}^2 - |\beta_{t,i}| = 2 \min\{a_{t,i}^2, b_{t,i}^2\} \stackrel{(i)}{\lesssim} \alpha^4, \quad (35)$$

where (i) follows by proof of proposition ??, since the smaller one of $\{a_{t,i}^2, b_{t,i}^2\}$ will be its initial value decreasing at an exponential rate until the rate α^4 . Similarly, we have the decomposition:

$$\begin{aligned} \|\beta_{t+1} - \beta^*\|^2 &= \|\beta_{t+1,S_1} - \beta_{S_1}^*\|^2 + \|\beta_{t+1,S_2} - \beta_{S_2}^*\|^2 + \|\beta_{t+1,S^c} - \beta_{S^c}^*\|^2 \\ &\leq \|\beta_{t+1,S_1} - \beta_{S_1}^*\|^2 + cs_2 \frac{\log p}{n} + c'p\alpha^4, \end{aligned}$$

which is based on the orthogonality between different sub vectors and the induction hypothesis, so we only need to calculate the errors induced by the strong signals.

For the gradient, we have:

$$\begin{aligned} \|n^{-1}X^T(Xr_t - w)\|_\infty &\leq \|\beta_{t,S} - \beta^*\|_\infty + \delta\|\beta_{t,S} - \beta^*\| \\ &\quad + \|\beta_{t,S^c}\|_\infty + \delta p\|\beta_{t,S^c}\|_\infty + \|X^Tw\|_\infty \leq c; \end{aligned} \quad (36)$$

$$\|n^{-1}X^T(X(I_{S_1^c}r_t) - w)\|_\infty \leq c\delta p\alpha^2 + c\sigma\sqrt{\frac{\log p}{n}}. \quad (37)$$

Note that $I_{S_1}\beta_t$ is s_1 sparse, so we have the following core error term characterization:

$$\|I_{S_1}\beta_t \circ n^{-1}X^T(Xr_{t,S_1^c} - w)\| \leq cs\delta^2p^2\alpha^4 + c'\sigma^2Ms_1\frac{\log p}{n} \leq c(\alpha^2 + \sigma^2\frac{Ms_1}{n}). \quad (38)$$

Similarly, our analysis for errors are based on equation (37) and (38).

For the order one term, recall the bound

$$\begin{aligned}
& \langle (a_t^2 + b_t^2) \circ n^{-1} X^T (Xr_t - w), I_{S_1} r_t \rangle \\
& \stackrel{(i)}{\geq} \langle |I_{S_1} \beta_t| \circ n^{-1} X^T (Xr_t - w), I_{S_1} r_t \rangle - s\alpha^4 \\
& = \langle |I_{S_1} \beta_t| \circ n^{-1} X^T X(I_{S_1} r_t), I_{S_1} r_t \rangle + \langle |I_{S_1} \beta_t| \circ n^{-1} X^T (X(I_{S_1^c} r_t) - w), I_{S_1} r_t \rangle \\
& \geq \left\| r_{t,S_1} \circ \sqrt{|\beta_{t,S_1}|} \right\|^2 - \delta \|r_{t,S_1}\|^2 \\
& \quad - \left\| \sqrt{|\beta_{t,S_1}|} \circ n^{-1} X^T (X(I_{S_1^c} r_t) - w) \right\|^2 / 2 - \left\| \sqrt{|\beta_{t,S_1}|} \circ r_{t,S_1} \right\|^2 / 2 \\
& \geq \frac{1}{2} \left\| r_{t,S_1} \circ \sqrt{|\beta_{t,S_1}|} \right\|^2 - \delta \|r_{t,S_1}\|^2 - c(\alpha^2 + \sigma^2 \frac{Ms_1}{n}).
\end{aligned}$$

where (i) is based on Cauchy-Schwarz inequality and we use the fact $(a_t, b_t) \in \Theta_{LR}^G$, lemma 6.1, equation (38).

For the rest of terms, we can see they are all based on the bound:

$$\begin{aligned}
& \|I_{S_1} \sqrt{\beta_t} \circ n^{-1} X^T (Xr_t - w)\|^2 \\
& \leq \|I_{S_1} \sqrt{\beta_t} \circ n^{-1} X^T X(I_{S_1} r_t)\|^2 + \|I_{S_1} \sqrt{\beta_t} \circ n^{-1} X^T (I_{S_1^c} r_t - w)\|^2 \\
& \lesssim c \|r_{t,S_1}\|^2 + c'(\alpha^2 + \sigma^2 \frac{Ms_1}{n}).
\end{aligned}$$

First consider the order two terms $\langle \beta_t \circ (n^{-1} X^T Xr_t)^2, I_{S_1} (\beta_t - \beta^*) \rangle$ and $\|I_{S_1} (a_t^2 + b_t^2) \circ n^{-1} X^T Xr_t\|^2$. For the former one, since $\|r_t\|_\infty \lesssim 1$, we have

$$\begin{aligned}
\langle \beta_t \circ (n^{-1} X^T Xr_t)^2, I_{S_1} r_t \rangle & = \langle I_{S_1} \beta_t \circ n^{-1} X^T (Xr_t - w), n^{-1} X^T (Xr_t - w) \circ r_t \rangle \\
& \leq c \|I_{S_1} \sqrt{\beta_t} \circ n^{-1} X^T (Xr_t - w)\|^2 \leq c \|r_{t,S_1}\|^2 + c'(\alpha^2 + \sigma^2 \frac{Ms_1}{n}).
\end{aligned}$$

where we use the condition $(a_t, b_t) \in \Theta_{LR}^G$. The analysis of the latter one $\|I_{S_1} (a_t^2 + b_t^2) \circ$

$n^{-1}X^T X r_t\|^2$ is similar based on $\|a_t^2 + b_t^2 - \beta_t\|_\infty \lesssim \alpha^2$ and $\|\sqrt{\beta_t}\| \lesssim 1$:

$$\begin{aligned} & \|I_{S_1}(a_t^2 + b_t^2) \circ n^{-1}X^T(Xr_t - w)\|^2 \\ & \leq \|(I_{S_1}|\beta_t| \circ n^{-1}X^T(Xr_t - w))^2 + s\alpha^4 \leq c\|r_{t,S_1}\|^2 + c'(\alpha^2 + \sigma^2 \frac{Ms_1}{n}). \end{aligned}$$

For the order three term, using inequality $\|I_{S_1}(a_t^2 + b_t^2) \circ (n^{-1}X^T X r_t)\|_\infty \lesssim 1$, we obtain:

$$\begin{aligned} & \langle \beta_t \circ (n^{-1}X^T X r_t)^2, I_{S_1}(a_t^2 + b_t^2) \circ (n^{-1}X^T X r_t) \rangle \\ & \leq c\|(I_{S_1}\sqrt{\beta_t}) \circ n^{-1}X^T(Xr_t - w)\|^2 \leq c\|r_{t,S_1}\|^2 + c'(\alpha^2 + \sigma^2 \frac{Ms_1}{n}). \end{aligned}$$

For the order four term, we have

$$\|I_{S_1}\beta_t \circ (n^{-1}X^T X r_t)^2\|^2 \leq c\|I_{S_1}\sqrt{\beta_t}(n^{-1}X^T X r_t)\|^2 \leq c\|r_{t,S_1}\|^2 + c'(\alpha^2 + \sigma^2 \frac{Ms_1}{n}).$$

Putting all pieces together, and using our condition $\eta \lesssim 1$, we can obtain

$$\begin{aligned} \|I_{S_1}r_{t+1}\|^2 & \leq \|I_{S_1}r_t\|^2 - 2\eta \left\| r_{t,S_1} \circ \sqrt{|\beta_{t,S_1}|} \right\|^2 \\ & \quad + c(\alpha^2 + \sigma^2 \frac{Ms_1}{n}) + 4\eta\delta\|r_{t,S_1}\|^2. \end{aligned}$$

Since for $(a_t, b_t) \in \Theta_{LR}^G$, we have $\|r_{t,S_1} \circ \sqrt{|\beta_{t,S_1}|}\|^2 \geq \|r_{t,S_1}\|^2 \cdot m/2$. Therefore, under our conditions on (δ, η) with $\delta\kappa \leq 1/4$, we have:

$$\|\beta_{t+1,S_1} - \beta_{S_1}^*\|^2 \leq (1 - \eta m)\|\beta_{t,S_1} - \beta_{S_1}^*\|^2 + c(\alpha^2 + \sigma^2 \frac{Ms_1}{n}). \quad (39)$$

□

6.12 Proof of Theorem 3.2

The proof is the same as that of theorem 3.1, but now using propositions 3.3 and 3.4.

6.13 Proof of Theorem 3.3

According to the proof of Theorem 3.2, we have $\|\beta_{t,S_1} - \beta_{S_1}^*\|_\infty \leq \|\beta_{t,S_1} - \beta_{S_1}^*\| \leq \sigma\sqrt{s/n}$, therefore $\beta_{t,j} > \beta_j^* - \sigma\sqrt{s/n} \geq \lambda$ based on the definition of the strong signals for $j \in S_1$, while for errors and weak signals $j \in S_1^c$, we have $\|\beta_{t,S_1^c}\|_\infty \leq c\alpha^2 < \lambda$. Consequently, after the component-wise hard thresholding operation at level λ , all strong signals remains nonzero while all weak signals and errors become zero.

References

- Beck, A. and Teboulle, M. (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202.
- Breheny, P. and Huang, J. (2011). Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection. *The annals of applied statistics*, 5(1):232.
- Bühlmann, P., Kalisch, M., and Meier, L. (2014). High-dimensional statistics with a view toward applications in biology. *Annual Review of Statistics and Its Application*, 1:255–278.
- Candes, E. and Tao, T. (2007). The dantzig selector: Statistical estimation when p is much larger than n. *The Annals of Statistics*, 35(6):2313–2351.
- Candes, E. J. (2008). The restricted isometry property and its implications for compressed sensing. *Comptes Rendus Mathématique*, 346(9-10):589–592.
- Chen, S. S., Donoho, D. L., and Saunders, M. A. (2001). Atomic decomposition by basis pursuit. *SIAM Review*, 43(1):129–159.
- Cohen, A., Dahmen, W., and DeVore, R. (2009). Compressed sensing and best k-term approximation. *Journal of the American Mathematical Society*, 22:211–231.

- Donoho, D. L. (2006). Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306.
- Efron, B. (2004). The estimation of prediction error: covariance penalties and cross-validation. *Journal of the American Statistical Association*, 99(467):619–632.
- Fan, J. and Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360.
- Fan, J. and Lv, J. (2008). Sure independence screening for ultrahigh dimensional feature space. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(5):849–911.
- Friedman, J., Hastie, T., and Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1):1.
- Gunasekar, S., Lee, J., Soudry, D., and Srebro, N. (2018). Characterizing implicit bias in terms of optimization geometry. In Dy, J. and Krause, A., editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1832–1841, Stockholmsmässan, Stockholm Sweden. PMLR.
- Gunasekar, S., Woodworth, B. E., Bhojanapalli, S., Neyshabur, B., and Srebro, N. (2017). Implicit regularization in matrix factorization. In *Advances in Neural Information Processing Systems*, pages 6152–6160.
- Hoff, P. D. (2017). Lasso, fractional norm and structured sparse estimation using a hadamard product parametrization. *Computational Statistics & Data Analysis*, 115:186–198.
- Jin, J. and Ke, Z. T. (2016). Rare and weak effects in large-scale inference: methods and phase diagrams. *Statistica Sinica*, pages 1–34.

- Lee, J. D., Simchowitz, M., Jordan, M. I., and Recht, B. (2016). Gradient descent only converges to minimizers. In *Conference on Learning Theory*, pages 1246–1257.
- Li, Y., Ma, T., and Zhang, H. (2018). Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pages 2–47.
- Raskutti, G., Wainwright, M. J., and Yu, B. (2014). Early stopping and non-parametric regression: an optimal data-dependent stopping rule. *The Journal of Machine Learning Research*, 15(1):335–366.
- Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. (2018). The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878.
- Stein, C. M. (1981). Estimation of the mean of a multivariate normal distribution. *The annals of Statistics*, pages 1135–1151.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288.
- Vito, E. D., Rosasco, L., Caponnetto, A., Giovannini, U. D., and Odone, F. (2005). Learning from examples as an inverse problem. *Journal of Machine Learning Research*, 6(May):883–904.
- Wainwright, M. J. (2009). Information-theoretic limits on sparsity recovery in the high-dimensional and noisy setting. *IEEE Transactions on Information Theory*, 55(12):5728–5741.
- Yao, Y., Rosasco, L., and Caponnetto, A. (2007). On early stopping in gradient descent learning. *Constructive Approximation*, 26(2):289–315.

- Zhang, C.-H. (2010). Nearly unbiased variable selection under minimax concave penalty. *The Annals of statistics*, 38(2):894–942.
- Zhang, T. and Yu, B. (2005). Boosting with early stopping: Convergence and consistency. *The Annals of Statistics*, 33(4):1538–1579.
- Zhao, T., Liu, H., and Zhang, T. (2018). Pathwise coordinate optimization for sparse learning: Algorithm and theory. *The Annals of Statistics*, 46(1):180–218.
- Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American statistical association*, 101(476):1418–1429.