

# Exponential Lasso: robust sparse penalization under heavy-tailed noise and outliers with exponential-type loss

The Tien Mai 

Norwegian Institute of Public Health, Oslo, 0456, Norway  
email: the.tien.mai@fhi.no

## Abstract

In high-dimensional statistics, the Lasso is a cornerstone method for simultaneous variable selection and parameter estimation. However, its reliance on the squared loss function renders it highly sensitive to outliers and heavy-tailed noise, potentially leading to unreliable model selection and biased estimates. To address this limitation, we introduce the Exponential Lasso, a novel robust method that integrates an exponential-type loss function within the Lasso framework. This loss function is designed to achieve a smooth trade-off between statistical efficiency under Gaussian noise and robustness against data contamination. Unlike other methods that cap the influence of large residuals, the exponential loss smoothly redescends, effectively down-weighting the impact of extreme outliers while preserving near-quadratic behavior for small errors. We establish theoretical guarantees showing that the Exponential Lasso achieves strong statistical convergence rates, matching the classical Lasso under ideal conditions while maintaining its robustness in the presence of heavy-tailed contamination. Computationally, the estimator is optimized efficiently via a Majorization-Minimization (MM) algorithm that iteratively solves a series of weighted Lasso subproblems. Numerical experiments demonstrate that the proposed method is highly competitive, outperforming the classical Lasso in contaminated settings and maintaining strong performance even under Gaussian noise.

Our method is implemented in the R package `heavylasso` available on Github: <https://github.com/tienmt/heavylasso>.

Keywords: heavy-tailed noise; Lasso; robust regression; sparsity; soft-thresholding; non-asymptotic bounds, outliers.

## 1 Introduction

In modern data analysis, it is common to encounter datasets where the number of features ( $p$ ) greatly exceeds the number of samples ( $n$ ), a setting that breaks down classical statistical methods. When faced with this high-dimensionality, the traditional least squares estimator becomes unreliable and ill-posed. This has spurred the development of regularization methods, which are designed to impose structure, prevent overfitting, and improve the model's interpretability by favoring simpler, sparser solutions [11, 3, 9].

A pioneering and widely adopted technique in this domain is the Lasso (least absolute shrinkage and selection operator) [30]. It provides a powerful framework for performing both variable selection and parameter estimation simultaneously. The Lasso finds the coefficient vector  $\beta$  by solving the

following optimization problem:

$$\hat{\boldsymbol{\beta}}^{\text{lasso}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\}$$

In this formulation, the first term is the standard least squares loss function, while the second is an  $L_1$  penalty weighted by a tuning parameter  $\lambda > 0$ . By penalizing the sum of the absolute values of the coefficients, the Lasso effectively shrinks many of them to exactly zero, thus performing automated feature selection. The foundational principles of the Lasso have ignited a vast body of research, leading to numerous advancements in sparse estimation and high-dimensional inference [38, 4, 36, 20, 2].

However, the squared loss underlying the classical Lasso implicitly assumes light-tailed, approximately Gaussian errors. In many real-world datasets—such as those arising in genomics, finance, and environmental monitoring—this assumption is frequently violated. Outliers or heavy-tailed noise can exert a disproportionate influence on the squared loss, leading to biased estimates and poor variable selection performance [17, 19]. To mitigate such sensitivity, numerous robust variants of the Lasso have been proposed, often by replacing the squared loss with more robust alternatives such as the Huber loss [26, 17, 18], Tukey’s biweight loss [5, 28], Student’s loss [22], or rank-based and median-of-means losses [24, 15, 32]. These methods reduce the impact of extreme residuals but may introduce additional tuning complexity or require nontrivial optimization schemes when the loss becomes nonconvex.

In this work, we consider a novel robust loss function—the exponential-type loss—within the Lasso penalization framework, designed to achieve a smooth trade-off between efficiency under Gaussian noise and robustness under heavy-tailed contamination. The proposed estimator, termed the Exponential Lasso, is defined as

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\{ \frac{1}{n} \sum_{i=1}^n \frac{1}{\tau} \left[ 1 - \exp \left( -\frac{\tau(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{2} \right) \right] + \lambda \|\boldsymbol{\beta}\|_1 \right\},$$

where  $\tau > 0$  controls the degree of robustness. When  $\tau \rightarrow 0$ , the exponential loss approaches the squared loss, recovering the classical Lasso. For finite  $\tau$ , large residuals are exponentially downweighted, effectively limiting their influence on parameter estimation.

The intuition behind this exponential-type loss is simple yet powerful: it penalizes small residuals nearly quadratically—preserving statistical efficiency under light-tailed noise—while suppressing the contribution of extreme deviations through the exponential term. In contrast to the Huber loss, which caps the linear growth of large residuals, the exponential loss smoothly redescends, assigning progressively smaller weights to extreme outliers. This property aligns with the influence-function perspective in robust statistics [10], where bounded and redescending functions provide strong resistance to contamination. As a result, the Exponential Lasso achieves both robustness and differentiability, enabling efficient optimization and stability in high-dimensional regimes.

From a computational standpoint, the exponential loss admits a natural Majorization–Minimization (MM) algorithmic interpretation. Each iteration reweights the residuals according to their exponential downweighting factor, leading to a sequence of weighted Lasso subproblems. This results in a fast, stable, and interpretable algorithm with minimal modification to standard Lasso solvers.

Our theoretical results establish that the proposed Exponential Lasso estimator achieves reliable estimation accuracy even in high-dimensional settings and under the presence of outliers or heavy-tailed noise. The theoretical result guarantees that, with an appropriate choice of the tuning

parameter, the estimator attains the same convergence rate as the classical Lasso under ideal conditions. This robustness is ensured under a mild assumption on the noise distribution, requiring only that the errors have a positive probability of lying within a central region rather than having light tails. The proof combines a local curvature argument showing that the exponential loss remains well-behaved near the true parameter with concentration techniques that control random fluctuations in the data [17, 18].

To demonstrate the effectiveness of our method, we conduct extensive simulation studies under a variety of settings involving heavy-tailed noise and outliers. We compare our approach to several Lasso variants that employ different loss functions, including the squared loss [30],  $\ell_1$  loss [31], Huber loss [35] and Student's loss [22]. The simulation results indicate that our method consistently exhibits strong empirical performance relative to these alternatives, particularly in challenging settings with non-Gaussian errors or contaminated by outliers. Moreover, our proposed method outperforms classical Lasso even in the Gaussian noise. In addition to simulations, we present real data applications that further supports the utility of our approach. Numerical results show that our proposed method are very competitive and promising.

The remainder of the paper is organized as follows. Section 2 introduces the proposed methodology alongside the Exponential Lasso approach and provides theoretical insights into the robustness properties of the loss function. This section also establishes non-asymptotic statistical guarantees for the proposed estimator. Section 3 details the algorithmic implementation and includes a convergence analysis of the proposed optimization scheme. Simulation studies evaluating empirical performance are presented in Section 4. Applications to two real datasets are discussed in Section 5. Finally, concluding remarks and discussions are provided in Section 6, while all technical proofs are collected in Appendix A.

## 2 Model and method

### 2.1 Robust Lasso with Exponential-type loss

Let  $\{(x_i, y_i)\}_{i=1}^n$  denote a collection of independent and identically distributed (i.i.d) observations arising from the linear model

$$y_i = x_i^\top \beta^* + \epsilon_i, \quad (1)$$

where  $x_i \in \mathbb{R}^p$  is the  $i$ -th row of the design matrix  $X$  and  $\beta^* \in \mathbb{R}^p$  is the unknown vector of regression coefficients that we seek to estimate. The condition on the random noise is given below in which we consider a very wild class of noise that cover both heavy-tailed noise and outlier contaminated models.

Let

$$L_\tau(\beta) = \frac{1}{n} \sum_{i=1}^n \ell_\tau(y_i - x_i^\top \beta), \quad \text{where} \quad \ell_\tau(r) = \frac{1}{\tau} (1 - e^{-\frac{\tau}{2} r^2}). \quad (2)$$

We consider the following robust penalized regression estimator, called Exponential Lasso:

$$\hat{\beta} := \arg \min_{\beta \in \mathbb{R}^p} \{L_\tau(\beta) + \lambda \|\beta\|_1\}, \quad (3)$$

where  $\tau > 0$  controls the degree of robustness and  $\lambda > 0$  is a regularization parameter. For small  $\tau$ , the loss approximates the quadratic loss and (3) reduces to the standard Lasso. For larger  $\tau$ , the exponential term strongly downweights large residuals, thus providing robustness against outliers.

## 2.2 Robustness of the loss

**Some intuitions:** One can explicitly show the connection to the standard Lasso. The exponential loss  $\ell_\tau(r)$  can be analyzed using a Taylor expansion for  $e^x$  around 0. Let  $u = -\frac{\tau}{2}r^2$ . Since  $e^u \approx 1+u$  for small  $u$ :

$$\ell_\tau(r) = \frac{1}{\tau} \left(1 - e^{-\frac{\tau}{2}r^2}\right) \approx \frac{1}{\tau} \left(1 - \left(1 - \frac{\tau}{2}r^2\right)\right) = \frac{1}{\tau} \left(\frac{\tau}{2}r^2\right) = \frac{1}{2}r^2.$$

This formally demonstrates that as  $\tau \rightarrow 0$ , your objective function  $L_\tau(\beta)$  converges to the standard least-squares loss, and thus  $\hat{\beta}$  converges to the classical Lasso estimator. This highlights that our proposed method is a natural generalization of the classical Lasso. See Figure 1 for a detailed visualization comparison between different losses: squared loss, Tukey’s biweight loss, absolute ( $\ell_1$ ) loss, and Huber loss. The plot illustrates that, unlike Huber loss, our loss is much less sensitive to large residuals while closely resembling the squared loss for small residual values.

**The Influence Function:** A key concept in robust statistics is the influence function [10], which measures the effect of an infinitesimal outlier on the estimator. The influence function is proportional to the first derivative of the loss function,  $\psi(r) = \ell'_\tau(r)$ , where

$$\psi(r) = \frac{\partial}{\partial r} \left[ \frac{1}{\tau} \left(1 - e^{-\frac{\tau}{2}r^2}\right) \right] = r e^{-\frac{\tau}{2}r^2}.$$

- The influence function is bounded: The maximum value of  $|r \exp(-\frac{\tau}{2}r^2)|$  is finite. This contrasts with the L2 loss, where  $\psi(r) = r$ , which is unbounded. An unbounded influence function means a single large outlier can have an arbitrarily large (i.e., infinite) influence on the estimate.
- It is redescending: As the residual  $r \rightarrow \infty$  (a gross outlier), the influence  $\psi(r) \rightarrow 0$ . This is a very strong form of robustness. The estimator completely ignores data points that are sufficiently far from the bulk of the data. This is an advantage over other robust losses like the Huber loss, whose influence function is bounded but not redescending (it becomes constant,  $\psi(r) = \text{sign}(r) \cdot k$ , for large  $r$ ).

**Other insights:** The proposed loss function can also be viewed as a correntropy or Welsch-type loss that downweights large residuals through an exponential kernel, thereby enhancing robustness against outliers, which is known in the robust signal-processing literature [16, 12]. Moreover, from an information-theoretic perspective, it is also closely connected to the  $\alpha$ -divergence, where  $\tau$  acts analogously to the divergence parameter controlling the trade-off between efficiency and robustness [14, 25]. This dual interpretation highlights the method’s grounding in both robust estimation and divergence-based statistical learning. We also note that similar exponential-type loss functions have been explored in prior studies as in [33, 29, 34], where the loss takes the form  $\ell_\tau(t) = 1 - e^{-t^2/\tau}$ . In contrast, our proposed formulation is more directly connected to the  $\alpha$ -divergence family. Furthermore, while the aforementioned works primarily address low-dimensional settings and provide asymptotic analyses, our study focuses on high-dimensional regimes and establishes non-asymptotic theoretical guarantees.

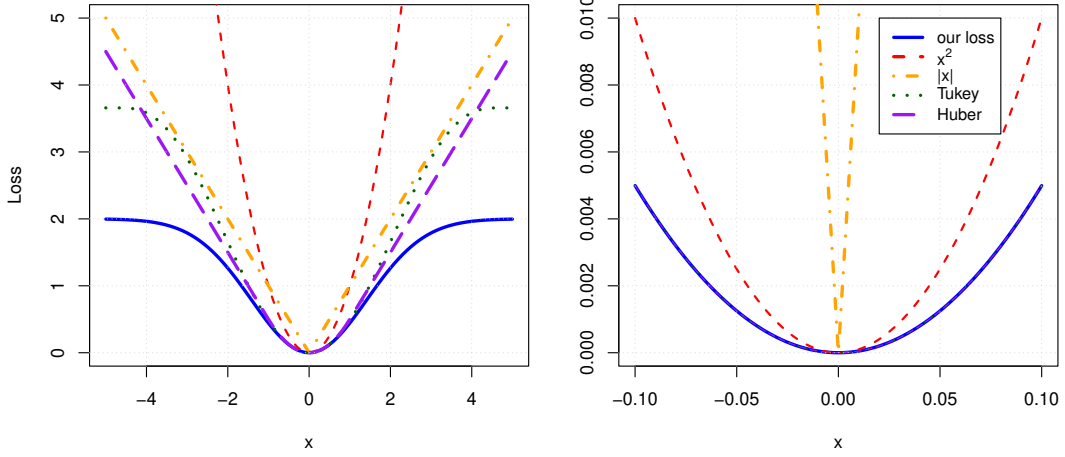


Figure 1: Comparison of our loss function ( $\tau = 0.5$ ) with other common losses: squared loss, absolute ( $\ell_1$ ) loss, Tukey's biweight loss, and Huber loss. The plot illustrates that, unlike Huber loss, our loss is much less sensitive to large residuals while closely resembling the squared loss for small residual values. Left: full-scale plot. Right: zoomed-in view near zero residuals.

### 2.3 Statistical guarantee

We now demonstrate that, under suitable assumptions, our Exponential-Lasso method achieves strong non-asymptotic theoretical guarantees comparable to those established for the Huber loss.

We make the following assumptions.

**Assumption 1** (Design and sparsity). *We assume that:*

- a). The rows  $x_i \in \mathbb{R}^p$  are non-random or random but satisfy  $\|x_i\|_\infty \leq K$  almost surely for a known constant  $K > 0$ .
- b). Let  $\beta^*$  be the true parameter with support  $S = \text{supp}(\beta^*)$  and sparsity  $s := |S| < n < p$ .
- c). There exists  $\phi_{\min} > 0$  and a radius  $r > 0$  such that for all  $\Delta \in \mathbb{R}^p$  with  $\|\Delta\|_2 \leq r$  and  $\|\Delta_{S^c}\|_1 \leq 3\|\Delta_S\|_1$ ,

$$\frac{1}{n} \sum_{i=1}^n (x_i^\top \Delta)^2 \geq \phi_{\min} \|\Delta\|_2^2.$$

**Assumption 2** (Noise). *The errors  $\varepsilon_i$  are i.i.d. with the following properties:*

- (i).  $\varepsilon_i$  are symmetric about 0.
- (ii). There exists a constant  $c \in (0, 1/\sqrt{\tau})$  and  $p_0 := \mathbb{P}(|\varepsilon_i| \leq c) > 0$ .

Remarks: symmetry can be relaxed by replacing the centering step with the population bias, but symmetry keeps statements concise. The choice  $c < 1/\sqrt{\tau}$  ensures positive curvature of the per-observation second derivative inside  $|r| \leq c$ .

**Discussion on noise assumptions.** The noise conditions assumed above are notably weaker than the conventional sub-Gaussian or even sub-exponential assumptions often imposed in high-dimensional regression analysis. In particular, condition (ii) only requires that the noise distribution has some probability mass around zero, ensuring that the majority of samples are not dominated by extreme outliers, while symmetry guarantees that the noise has zero median. No moment or exponential tail condition is imposed, allowing the framework to accommodate a broad class of heavy-tailed or contamination models, such as Student's  $t$  distributions with small degrees of freedom or Huber error models of the form

$$\varepsilon_i \sim (1 - \pi) N(0, \sigma^2) + \pi G,$$

where  $G$  may represent a heavy-tailed or outlier-generating component (e.g., Cauchy or Laplace). Therefore, this assumption captures realistic data-generating mechanisms with occasional large deviations, while retaining sufficient regularity for establishing estimation error bounds. It is thus particularly suitable for robust (penalized) regression models designed to handle impulsive noise or mild contamination in the observations.

Define the positive curvature on the interval  $[-c, c]$ ,

$$\underline{\gamma} := \min_{|u| \leq c} e^{-\frac{\tau}{2}u^2} (1 - \tau u^2) = e^{-\frac{\tau}{2}c^2} (1 - \tau c^2) > 0.$$

**Theorem 1.** *Under Assumptions 1-2, fix  $\delta \in (0, 1)$ . Choose the tuning parameter*

$$\lambda = \frac{4K}{\sqrt{e\tau}} \sqrt{\frac{2 \log(2p/\delta)}{n}}. \quad (4)$$

*Assume  $r > 0$  in Assumption 1 is small enough so that for all  $\Delta$  in the cone  $\{\|\Delta\|_2 \leq r, \|\Delta_{S^c}\|_1 \leq 3\|\Delta_S\|_1\}$  it holds that  $|x_i^\top \Delta| \leq c/2$  for all  $i$  (this is satisfied when  $r$  is chosen such that  $K\sqrt{s}r \leq c/2$ ).*

*Then with probability at least  $1 - \delta - 2 \exp(-np_0^2/8)$  any global minimizer  $\hat{\beta}$  satisfying  $\|\hat{\beta} - \beta^*\|_2 \leq r$  obeys*

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{12 \lambda \sqrt{s}}{\kappa}, \quad \text{and} \quad \|\hat{\beta} - \beta^*\|_1 \leq \frac{48 \lambda s}{\kappa}, \quad (5)$$

*where the constant  $\kappa > 0$  can be taken as  $\kappa = \frac{p_0}{2} \underline{\gamma} \phi_{\min}$ . In particular, with the choice (4) we obtain the explicit bound*

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{48K}{\kappa \sqrt{e\tau}} \sqrt{\frac{s \log(2p/\delta)}{n}}.$$

The theoretical analysis is based on the Local Restricted Strong Convexity (LRSC) condition, with a full proof provided in Appendix A. The methodology largely adopts the general framework proposed in [17, 18], and as such, our results parallel those obtained for the Huber-loss Lasso. The primary contribution of our approach is the reliance on a substantially weaker noise condition, as we do not assume the existence of any moments for the noise distribution.

### 3 Majorization–Minimization Algorithm

#### 3.1 Algorithm development

The exponential-type loss in (3) is nonconvex but smooth. To derive an efficient iterative algorithm, we adopt a majorization–minimization (MM) approach. Let us define for each observation

$$\ell_i(\beta) = \frac{1}{\tau} \left( 1 - \exp \left( -\frac{\tau}{2} r_i(\beta)^2 \right) \right), \quad r_i(\beta) = y_i - x_i^\top \beta.$$

Consider the function  $\phi(u) = 1 - \exp(-\frac{\tau}{2}u)$  for  $u \geq 0$ . Since  $\phi''(u) = -(\frac{\tau}{2})^2 \exp(-\frac{\tau}{2}u) < 0$ , the function  $\phi$  is concave. For any fixed  $u^{(t)}$ , the first-order Taylor expansion provides a global upper bound (a tight majorizer):

$$\phi(u) \leq \phi(u^{(t)}) + \phi'(u^{(t)})(u - u^{(t)}), \quad \forall u \geq 0. \quad (6)$$

Substituting  $u_i = r_i(\beta)^2$  and summing over  $i$  yields the following upper bound for the empirical loss:

$$\frac{1}{n} \sum_{i=1}^n \ell_i(\beta) \leq C + \frac{1}{2n} \sum_{i=1}^n v_i^{(t)} r_i(\beta)^2, \quad (7)$$

where  $C$  is a constant independent of  $\beta$ , and

$$v_i^{(t)} = \exp \left( -\frac{\tau}{2} r_i(\beta^{(t)})^2 \right) \quad (8)$$

acts as an adaptive weight. Minimizing the right-hand side of the above inequality thus defines the MM update.

Given the current estimate  $\beta^{(t)}$ , the MM procedure alternates between updating the weights  $v_i^{(t)}$  according to (8) and solving a weighted Lasso subproblem:

- **(Step 1)** Compute residuals  $r_i^{(t)} = y_i - x_i^\top \beta^{(t)}$  and update the weights

$$v_i^{(t)} = \exp \left( -\frac{\tau}{2} (r_i^{(t)})^2 \right), \quad i = 1, \dots, n.$$

These weights downweight observations with large residuals, thereby reducing the influence of outliers.

- **(Step 2)** Update the regression coefficients by solving the weighted Lasso problem

$$\beta^{(t+1)} = \arg \min_{\beta \in \mathbb{R}^p} Q^{(t)}(\beta) := \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \sum_{i=1}^n v_i^{(t)} (y_i - x_i^\top \beta)^2 + \lambda \|\beta\|_1 \right\}. \quad (9)$$

The above two steps are repeated until convergence, e.g., until  $\|\beta^{(t+1)} - \beta^{(t)}\|_2 / (1 + \|\beta^{(t)}\|_2) < \varepsilon$  for a small tolerance  $\varepsilon > 0$ . Step 2 can be efficiently implemented using standard coordinate-descent algorithms or existing Lasso solvers (e.g., `glmnet` in R) by specifying the observation weights  $\{v_i^{(t)}\}$ . The outline of our proposed algorithm given in Algorithm 1.

### 3.2 Coordinate descent updates using soft-thresholding

We estimate the regression coefficients in (9) via a coordinate descent algorithm. Let denote the residual vector at iteration  $t$  as  $r^{(t)} = y - X\beta^{(t)}$ . For the  $j$ -th coordinate, we define the corresponding partial residual — that is, the residual excluding the contribution of variable  $j$  — as

$$r_j = r^{(t)} + X_j\beta_j^{(t)}.$$

The objective function restricted to the single coefficient  $\beta_j$  can then be expressed as:

$$\frac{1}{2n} \sum_{i=1}^n v_i (r_{ij} - x_{ij}\beta_j)^2 + \lambda |\beta_j|.$$

Minimization of this univariate problem yields a closed-form update via the soft-thresholding operator, [7]:

$$z_j := \sum_{i=1}^n v_i x_{ij} r_{ij}, \quad U_j := \sum_{i=1}^n v_i x_{ij}^2, \quad \beta_j \leftarrow \frac{1}{U_j} \mathcal{S}(z_j, \lambda),$$

where

$$\mathcal{S}(z, \lambda) = \text{sign}(z) \cdot \max(|z| - \lambda, 0),$$

denotes the soft-thresholding function.

**Discussion.** The proposed algorithm can be viewed as an EM-like procedure, where the weights  $\{v_i^{(t)}\}$  play a role analogous to latent variables that reflect the reliability of each observation. As  $\tau \rightarrow 0$ , all  $v_i^{(t)} \rightarrow 1$ , and the algorithm reduces to the ordinary Lasso. For larger  $\tau$ , the exponential decay of  $v_i^{(t)}$  produces strong robustness to large residuals.

---

#### Algorithm 1 MM Algorithm for Exponential-type robust Lasso

---

- 1: **Input:** data  $(x_i, y_i)_{i=1}^n$ , tuning parameters  $\tau > 0$ ,  $\lambda \geq 0$ .
- 2: Initialize  $\beta^{(0)}$  (e.g., by ordinary Lasso).
- 3: **repeat**
- 4:   Compute residuals  $r_i^{(t)} = y_i - x_i^\top \beta^{(t)}$ .
- 5:   Update weights  $v_i^{(t)} = \exp\left(-\frac{\tau}{2}(r_i^{(t)})^2\right)$ .
- 6:   Solve weighted Lasso:

$$\beta^{(t+1)} = \arg \min_{\beta} \left\{ \frac{1}{2n} \sum_{i=1}^n v_i^{(t)} (y_i - x_i^\top \beta)^2 + \lambda \|\beta\|_1 \right\}.$$

- 7: **until** convergence criterion is met.
  - 8: **Output:** final estimate  $\hat{\beta} = \beta^{(t+1)}$ .
- 

Under standard regularity conditions for MM algorithms, the sequence  $\{\beta^{(t)}\}$  monotonically decreases the objective in (3) and converges to a stationary point.



### 3.3 Convergence analysis of the MM algorithm

We analyze the Majorization–Minimization (MM) algorithm given in Algorithm 1 for the exponential-type robust Lasso objective

$$F(\beta) = \frac{1}{n} \sum_{i=1}^n \ell_i(\beta) + \lambda \|\beta\|_1 \quad \text{with} \quad \ell_i(\beta) = \frac{1}{\tau} \left( 1 - \exp \left( -\frac{\tau}{2} r_i(\beta)^2 \right) \right), \quad r_i(\beta) = y_i - x_i^\top \beta, \quad (10)$$

where  $\tau > 0$  and  $\lambda \geq 0$ . For convenience write the smooth (nonconvex) loss part as

$$L(\beta) := \frac{1}{n} \sum_{i=1}^n \ell_i(\beta),$$

so  $F(\beta) = L(\beta) + \lambda \|\beta\|_1$ .

We first show that the MM iterates produce a monotone decreasing sequence of objective values and remain in a bounded level set.

**Theorem 2** (Monotone decrease and boundedness). *Assume that each row  $x_i$  satisfies  $\|x_i\|_2 < \infty$ . Let  $\{\beta^{(t)}\}$  be the sequence produced by Algorithm 1 where in each M-step we compute an exact minimizer of the surrogate  $Q^{(t)}(\beta)$  in (9). Then:*

- (i). (Descent) *The sequence of objective values  $\{F(\beta^{(t)})\}$  is nonincreasing:  $F(\beta^{(t+1)}) \leq F(\beta^{(t)})$ ,  $\forall t \geq 0$ .*
- (ii). (Lower bounded)  *$F(\beta) \geq 0$  for all  $\beta$ , hence  $\{F(\beta^{(t)})\}$  converges to a finite limit  $F^*$ .*
- (iii). (Bounded iterates) *The iterates  $\{\beta^{(t)}\}$  lie in the sublevel set  $\{\beta : F(\beta) \leq F(\beta^{(0)})\}$ , which is bounded under the mild condition that  $\lambda > 0$  or that  $X$  has full column rank; hence  $\{\beta^{(t)}\}$  is bounded.*

Next we establish that any accumulation point of the MM iterates is a stationary point of the nonconvex objective  $F(\beta)$ . Because  $F$  contains the nondifferentiable  $\ell_1$  term, stationarity is understood in the subgradient/KKT sense.

**Theorem 3** (Cluster points are stationary). *Let assume that each row  $x_i$  satisfies  $\|x_i\|_2 < \infty$  and let  $\{\beta^{(t)}\}$  be generated by Algorithm 1 with exact M-steps (exact minimizer of the surrogate). Then every limit point  $\beta^*$  of  $\{\beta^{(t)}\}$  is a stationary point of  $F$ , i.e.  $0 \in \nabla L(\beta^*) + \lambda \partial \|\beta^*\|_1$ . Consequently, the whole sequence has its set of cluster points contained in the set of stationary points of  $F$ .*

**Corollary 1** (Convergence of objective and subsequential stationarity). *Under the hypotheses of Theorems 2–3, the objective values  $F(\beta^{(t)})$  converge to a finite  $F^*$  and every cluster point of  $\{\beta^{(t)}\}$  is a stationary point of  $F$ . If, in addition, the set of stationary points at level  $F^*$  is finite and the sequence has a unique cluster point, then  $\beta^{(t)} \rightarrow \beta^*$ .*

The objective  $F(\beta)$  is nonconvex because  $L(\beta)$  is nonconvex; therefore MM can only be expected to converge to a local stationary point in general. Uniqueness of the weighted Lasso minimizer in each M-step holds when  $X^\top \tilde{V}^{(t)} X$  is positive definite (e.g. full column rank and all  $\tilde{v}_i^{(t)} > 0$ ). In practice  $\tilde{v}_i^{(t)} \in (0, 1]$ , so strict positive definiteness reduces to conditions on  $X$ .

### 3.4 Tuning the regularization parameter $\lambda$

We select the regularization parameter  $\lambda$  in the exponential Lasso method primarily using  $K$ -fold cross-validation (CV), which serves as the default option. Cross-validation evaluates predictive performance by partitioning the data into training and validation folds and choosing the  $\lambda$  that minimizes the average prediction error. This provides a practical, data-driven balance between model sparsity and predictive accuracy.

Our method is implemented in the R package `heavylasso` available on Github: <https://github.com/tienmt/heavylasso>.

## 4 Simulation studies

### 4.1 Setup

#### Compared methods

The central aim of this simulation is to evaluate how different robust loss functions affect the performance of the Lasso estimator. Accordingly, our investigation is strictly focused on regression methods that incorporate the Lasso penalty, to the exclusion of other regularization approaches.

The following four estimators are included in our comparative analysis:

- Classical Lasso, which minimizes a squared loss function, implemented in the R package `glmnet` [8].
- Huber Lasso, which employs the hybrid Huber loss function.
- LAD Lasso, which is based on the  $\ell_1$  loss function. Both the LAD and Huber variants are implemented in the R package `hqreg`.
- Heavy Lasso, which utilizes a Student loss function, available in the R package `heavylasso` [22].

#### Simulation settings

In our data generation process, the predictors  $X_i$  follow a  $N(0, \Sigma)$  distribution. We investigate two specific covariance scenarios: (i) an identity matrix ( $\Sigma = \mathbb{I}_p$ ), representing independent predictors, and (ii) an autoregressive correlation structure ( $\Sigma_{ij} = \rho_X^{|i-j|}$ ).

The ground-truth vector  $\beta_0$  is  $s$ -sparse, with the  $s$  non-zero coefficients split evenly, taking values of 1 (for  $s/2$  entries) and  $-1$  (for the remaining  $s/2$  entries). The responses  $y_i$  are then produced via the linear model (1) under the following noise conditions:

- Gaussian noise,  $\epsilon_i \sim \mathcal{N}(0, 1)$ . This serves as a baseline to assess how various robust methods perform under ideal (light-tailed) conditions.
- Gaussian noise with large variance.  $\epsilon_i \sim \mathcal{N}(0, 3^2)$ . This setting introduces moderate heavy-tailed behavior through increased variance.
- Student noise.  $\epsilon_i \sim t_3$ . This case represents heavy-tailed noise with finite variance.

- Cauchy noise.  $\epsilon_i \sim \text{Cauchy}$ . This represents a more extreme heavy-tailed setting with infinite variance.
- Contaminated with outliers.  $\epsilon_i \sim \mathcal{N}(0, 1)$  or  $\epsilon_i \sim t_3$  but some portion of the observed responses are further contaminated by outliers. This setting evaluates robustness to contamination.

We assessed the performance of each method from three different angles: parameter estimation, prediction on new data, and variable selection accuracy. We used two metrics to evaluate how accurately each model estimated the true coefficients. First, we measured the estimation error using the squared  $\ell_2$  norm, which calculates the distance between the estimated coefficients ( $\hat{\beta}$ ) and the true coefficients ( $\beta_0$ ):  $\|\hat{\beta} - \beta_0\|_2^2$ . Second, we evaluated the error in the model’s fitted values on the training data, captured by the linear predictor error:

$$\ell(X^\top \beta_0) := \frac{1}{n} \|X^\top (\hat{\beta} - \beta_0)\|_2^2$$

To measure prediction accuracy, we calculated the mean squared prediction error (MSPE) on a large, independent test set. This test set,  $(X_{\text{test}}, y_{\text{test}})$ , was generated from the same model as the training data, with a fixed size of  $n_{\text{test}} = 5000$ . The MSPE is defined as:

$$\text{MSPE}_{\text{test}} := \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \left( y_{\text{test},i} - X_{\text{test},i}^\top \hat{\beta} \right)^2.$$

Finally, we assessed each method’s ability to correctly identify the relevant predictors in the model. This was measured using two standard metrics: the True Positive Rate (TPR) and the False Discovery Rate (FDR).

Each simulation setting is repeated 100 times, and we report the mean and standard deviation of the results. The outcomes are presented in Tables 1, 2, 3, 4 and 5. The regularized parameters for all five methods are selected via 5-fold cross-validation. The tuning parameter  $\tau$  in our method is set to 0.1, which is motivated from sensitivity analysis in Subsection 4.2.3 below.

## 4.2 Simulations results

### 4.2.1 Results with Heavy-tailed noises

We first examine the performance of our method against competing approaches across several noise settings. This comparison is conducted in two distinct regimes: a small scale setting with  $p = 120$  and  $n = 100$  (Tables 1, 2), and a medium scale setting with  $p = 500$  and  $n = 300$  (Tables 3, 4). The true sparsity  $s^*$  is fixed at 10 for all experiments.

In both settings, particularly under Gaussian noise, our proposed method and the Heavy Lasso are shown to be highly competitive. They frequently outperform the classical Lasso, a level of performance not attained by the Huber Lasso or the L1 Lasso in these experiments. The robustness and superiority of these results appear to improve in the larger-scale data scenarios, as evidenced in Tables 3 and 4.

A general observation is that our proposed method and the Heavy Lasso, both being non-convex, tend to return more small non-zero coefficients. This characteristic typically leads to a higher false positive rate (or false discovery rate) when compared to convex alternatives like the Huber Lasso or the L1 Lasso. Despite this, a crucial advantage is their ability to consistently select the true support.

When considering heavy-tailed noise, such as the Student’s  $t_3$  or Cauchy distributions, our proposed method and the Heavy Lasso invariably provide the best results regarding both estimation and prediction errors. More particularly, in the challenging Cauchy noise scenario, our method demonstrates clear superiority across all metrics, including estimation, prediction, and variable selection accuracy.

#### 4.2.2 Results with Outliers

We next evaluate the behavior of all considered methods in the presence of outliers. For this purpose, we fix the simulation parameters at  $p = 500$ ,  $s^* = 10$ , and  $n = 300$ . We investigate two underlying noise distributions: standard normal,  $\mathcal{N}(0, 1)$ , and Student’s  $t_3$ . The proportion of the response data contaminated by outliers is varied among 10%, 20%, and 30%. The results of this analysis are presented in Table 5.

The results indicate that our proposed method outperforms both the Huber Lasso and the L1 Lasso. As expected, the classical Lasso breaks down entirely under these conditions. When the underlying noise is Gaussian, our method is significantly better than the Heavy Lasso (which utilizes a Student’s loss). This performance gap is particularly pronounced when a larger fraction of the responses are contaminated, for example, at the 30% level. In the presence of Student’s  $t_3$  noise, the Heavy Lasso exhibits a slight advantage at lower contamination levels. However, as the contamination fraction increases to 30%, our proposed method once again outperforms the Heavy Lasso. These findings clearly highlight the superior robustness of our proposed methodology.

#### 4.2.3 On sensitivity of tuning parameter $\tau$

To evaluate the sensitivity of our proposed method to the tuning parameter  $\tau$ , we conducted a dedicated simulation study. We fixed the problem dimensions at  $p = 120$ ,  $s^* = 10$ , and  $n = 100$ , while varying  $\tau$  over the grid  $\{0.001, 0.1, 1, 10\}$ .

The results, averaged over 100 replications under various noise distributions and outlier settings, are presented in Table 6. This analysis demonstrates that  $\tau = 0.1$  consistently yields the best and most stable performance. Based on this finding, we adopted  $\tau = 0.1$  as the fixed value for this hyperparameter in all other experiments and the real-data application.

### 4.3 Results with increasing sparsity

We further investigate how the sparsity level ( $s^*$ ) influences the performance of the various methods under conditions of heavy-tailed noise and outliers. In this analysis, we fixed the dimensionality at  $p = 500$  and  $n = 300$  with independent predictors ( $\Sigma = \mathbb{I}_p$ ). We then vary the true sparsity  $s^*$  across the values  $\{4, 8, 16\}$ .

The averaged results from 100 simulation repetitions, presented in Table 7, confirm that all methods exhibit a natural increase in both estimation and prediction errors as the sparsity increases. Crucially, our proposed method consistently maintains its position as the top performer across all sparsity levels in terms of both error metrics. The superiority of our method is most pronounced in the highly sparse setting where  $s^* = 16$ , where it significantly outperforms all other considered approaches. This comprehensive test demonstrates that the robustness of our method extends not only to non-Gaussian noise but also to increased model complexity due to higher sparsity.

Table 1: *Simulation for various loss functions with Lasso penalization, under the setting  $p = 120, s^* = 10, n = 100$  and independent predictors. The reported values are the mean across 100 simulation repetitions, with the standard deviation provided in parentheses. Bold font highlights the superior method. TPR: true positive rate; FDR: false discovery rate;  $MSPE_{test}$ : mean squared prediction error on testing data.*

Noise	Method (loss)	$\ \hat{\beta} - \beta_0\ _2^2$	$\ell(X^\top \beta_0)$	$MSPE_{test}$	TPR	FPR
$\mathcal{N}(0, 1)$	proposed loss	<b>0.62</b> (0.22)	<b>0.40</b> (0.11)	<b>1.62</b> (0.22)	<b>1.00</b> (0.00)	0.75 (0.06)
	Student's loss	0.64 (0.26)	<b>0.40</b> (0.11)	1.64 (0.26)	<b>1.00</b> (0.00)	0.75 (0.06)
	squared loss	0.78 (0.26)	0.55 (0.18)	1.82 (0.26)	<b>1.00</b> (0.00)	0.45 (0.13)
	$\ell_1$ loss	1.09 (0.43)	0.71 (0.24)	2.02 (0.43)	<b>1.00</b> (0.00)	0.65 (0.10)
	Huber loss	0.95 (0.41)	0.65 (0.27)	1.93 (0.40)	<b>1.00</b> (0.00)	0.57 (0.10)
$\mathcal{N}(0, 3)$	proposed loss	7.06 (2.01)	5.34 (1.77)	16.0 (2.07)	0.78 (0.28)	0.78 (0.13)
	Student's loss	<b>5.24</b> (1.67)	<b>3.79</b> (1.17)	<b>14.3</b> (1.74)	<b>0.89</b> (0.17)	0.72 (0.10)
	squared loss	7.00 (2.26)	6.31 (2.52)	16.0 (2.32)	0.52 (0.36)	0.25 (0.23)
	$\ell_1$ loss	8.29 (1.91)	7.82 (2.50)	17.3 (1.96)	0.32 (0.33)	0.22 (0.28)
	Huber loss	7.89 (1.94)	7.23 (2.41)	16.9 (1.99)	0.40 (0.33)	0.26 (0.27)
$t_3$	proposed loss	<b>1.20</b> (0.52)	<b>0.77</b> (0.37)	<b>4.15</b> (0.78)	<b>1.00</b> (0.00)	0.76 (0.06)
	Student's loss	1.29 (0.61)	0.81 (0.32)	4.25 (0.78)	<b>1.00</b> (0.00)	0.75 (0.06)
	squared loss	2.72 (2.29)	2.10 (2.06)	5.67 (2.27)	0.93 (0.24)	0.37 (0.17)
	$\ell_1$ loss	1.87 (1.05)	1.27 (0.70)	5.03 (1.14)	0.99 (0.04)	0.61 (0.12)
	Huber loss	1.73 (0.87)	1.20 (0.63)	4.71 (1.03)	0.99 (0.04)	0.53 (0.11)
<i>Cauchy</i>	proposed loss	<b>5.17</b> (3.23)	<b>3.71</b> (2.10)	$>10^5$	0.89 (0.25)	0.71 (0.21)
	Student's loss	5.82 (5.19)	4.01 (4.33)	$>10^5$	<b>0.91</b> (0.24)	0.71 (0.25)
	squared loss	9.97 (0.14)	9.95 (1.33)	$>10^5$	0.01 (0.06)	<b>0.01</b> (0.11)
	$\ell_1$ loss	9.19 (2.11)	8.57 (2.72)	$>10^5$	0.13 (0.31)	0.13 (0.24)
	Huber loss	9.13 (2.14)	8.59 (2.71)	$>10^5$	0.15 (0.33)	0.10 (0.18)

## 5 Application examples with real data

### 5.1 Analyzing cancer cell line NCI-60 data

We evaluate the proposed method using the NCI-60 cancer cell line panel, a benchmark dataset widely used in genomic and proteomic modeling studies [23]. A pre-processed version of the data, available in the R package `robustHD` [1], was employed in our analysis. The dataset comprises  $n = 59$  human cancer cell lines, each characterized by gene and protein expression measurements. One sample with missing gene expression values was excluded from the analysis.

The gene expression data form a matrix with 22,283 features, while the protein expression data consist of 162 features measured using reverse-phase protein lysate arrays. The protein data were  $\log_2$ -transformed and standardized to have zero mean, following the preprocessing steps in [1]. This dataset provides a rich setting for studying the relationship between high-dimensional genomic profiles and proteomic responses. Consistent with prior work, we model protein 92 as the response variable and select the 300 gene expression features exhibiting the highest marginal correlations with it to form the predictor matrix.

For model evaluation, nine samples were randomly set aside as a test set, and the remaining 50 samples were used for model training. This random partitioning was repeated 100 times, and the mean squared prediction error (MSPE) on the held-out test data was averaged across replications. Table 8 summarizes the results.

Table 2: *Simulation for various loss functions with Lasso penalization, under the setting  $p = 120, s^* = 10, n = 100$  and correlated design  $\rho_X = 0.5$ . The reported values are the mean across 100 simulation repetitions, with the standard deviation provided in parentheses. Bold font highlights the superior method. TPR: true positive rate; FDR: false discovery rate;  $MSPE_{\text{test}}$ : mean squared prediction error on testing data.*

Noise	Method (loss)	$\ \hat{\beta} - \beta_0\ _2^2$	$\ell(X^\top \beta_0)$	$MSPE_{\text{test}}$	TPR	FPR
$\mathcal{N}(0, 1)$	proposed loss	0.52 (0.27)	0.30 (0.09)	1.42 (0.20)	<b>1.00</b> (0.00)	0.70 (0.08)
	Student's loss	<b>0.50</b> (0.22)	<b>0.29</b> (0.11)	<b>1.41</b> (0.22)	<b>1.00</b> (0.00)	0.75 (0.06)
	squared loss	0.72 (0.26)	0.45 (0.18)	1.59 (0.26)	<b>1.00</b> (0.00)	<b>0.31</b> (0.13)
	$\ell_1$ loss	0.94 (0.43)	0.59 (0.24)	1.77 (0.43)	<b>1.00</b> (0.00)	0.48 (0.10)
	Huber loss	0.79 (0.41)	0.50 (0.27)	1.65 (0.40)	<b>1.00</b> (0.00)	0.38 (0.10)
$\mathcal{N}(0, 3)$	proposed loss	5.17 (1.85)	3.67 (1.16)	13.7 (1.86)	0.91 (0.09)	0.72 (0.06)
	Student's loss	<b>3.96</b> (1.67)	<b>2.61</b> (1.17)	<b>12.4</b> (1.74)	<b>0.92</b> (0.10)	0.66 (0.10)
	squared loss	4.55 (2.26)	4.31 (2.52)	14.0 (2.32)	0.76 (0.36)	<b>0.15</b> (0.23)
	$\ell_1$ loss	5.37 (1.91)	5.43 (2.50)	15.3 (1.96)	0.67 (0.33)	0.24 (0.28)
	Huber loss	5.25 (1.94)	5.29 (2.41)	15.1 (1.99)	0.69 (0.33)	0.21 (0.27)
$t_3$	proposed loss	<b>0.98</b> (0.40)	<b>0.56</b> (0.19)	<b>3.91</b> (1.14)	1.00 (0.00)	0.67 (0.10)
	Student's loss	1.03 (0.46)	0.59 (0.21)	3.95 (1.16)	<b>1.00</b> (0.00)	0.68 (0.10)
	squared loss	2.16 (2.29)	1.75 (2.16)	5.28 (2.27)	0.93 (0.24)	0.24 (0.17)
	$\ell_1$ loss	1.83 (1.05)	1.09 (0.70)	4.61 (1.14)	0.96 (0.04)	0.43 (0.12)
	Huber loss	1.65 (0.87)	1.06 (0.63)	4.51 (1.03)	0.97 (0.04)	0.34 (0.11)
<i>Cauchy</i>	proposed loss	<b>3.34</b> (2.71)	<b>2.26</b> (2.10)	$>10^5$	<b>0.96</b> (0.11)	0.60 (0.21)
	Student's loss	4.48 (4.72)	2.73 (2.21)	$>10^5$	0.94 (0.24)	0.62 (0.25)
	squared loss	9.80 (0.14)	9.95 (1.33)	$>10^5$	0.03 (0.06)	<b>0.00</b> (0.11)
	$\ell_1$ loss	7.90 (2.11)	8.57 (2.72)	$>10^5$	0.34 (0.31)	0.03 (0.24)
	Huber loss	7.71 (2.14)	8.59 (2.71)	$>10^5$	0.35 (0.33)	0.03 (0.18)

Our proposed estimator and the heavy-tailed Lasso based on Student's loss achieved the lowest prediction errors, indicating superior robustness to outliers and noise. Methods employing the Huber and  $\ell_1$  losses also outperformed the standard Lasso with the squared loss. In terms of variable selection, all five methods consistently identified one common predictor (gene ID 8502), highlighting its potential biological relevance.

## 5.2 Analyzing gene expression TRIM32 data

In this application, we conduct an analysis using high-dimensional genomics data from [27]. The study by [27] involved analyzing RNA from the eyes of 120 twelve-week-old male rats, using 31,042 different probe sets. Our focus is on modeling the expression of the gene TRIM32, as it was identified by [6] as a gene associated with Bardet-Biedl syndrome, a condition that includes retinal degeneration among its symptoms. Since [27] observed that many probes were not expressed in the eye, we follow the approach of [13] and [21], limiting our analysis to the 500 genes with the highest absolute Pearson correlation with TRIM32 expression. The data for this analysis is available from the R package *abess*, [37].

To assess the methods, we randomly allocate 84 of the 120 samples for training and the remaining 36 for testing, maintaining an approximate 70/30 percent of the data split. The methods are executed using the training set, and their prediction accuracy is evaluated on the test set. This

Table 3: *Simulation results for various loss functions with Lasso penalization, under the setting  $p = 500, s^* = 10, n = 300$  and independent predictors. The reported values are the mean across 100 simulation repetitions, with the standard deviation provided in parentheses. Bold font highlights the superior method. TPR: true positive rate; FDR: false discovery rate;  $MSPE_{test}$ : mean squared prediction error on testing data.*

Noise	Method (loss)	$\ \hat{\beta} - \beta_0\ _2^2$	$\ell(X^\top \beta_0)$	$MSPE_{test}$	TPR	FPR
$\mathcal{N}(0, 1)$	proposed loss	<b>0.23</b> (0.08)	<b>0.19</b> (0.06)	<b>1.23</b> (0.08)	1.00 (0.00)	0.85 (0.04)
	Student's loss	<b>0.23</b> (0.08)	<b>0.19</b> (0.06)	<b>1.23</b> (0.08)	1.00 (0.00)	0.85 (0.04)
	squared loss	0.31 (0.10)	0.28 (0.08)	1.31 (0.10)	1.00 (0.00)	<b>0.42</b> (0.17)
	$\ell_1$ loss	0.44 (0.16)	0.38 (0.12)	1.44 (0.16)	1.00 (0.00)	0.54 (0.16)
	Huber loss	0.39 (0.14)	0.35 (0.10)	1.39 (0.14)	1.00 (0.00)	0.48 (0.19)
$\mathcal{N}(0, 3)$	proposed loss	3.49 (1.32)	3.00 (1.11)	12.4 (1.46)	<b>1.00</b> (0.01)	0.89 (0.04)
	Student's loss	<b>2.36</b> (0.59)	<b>2.04</b> (0.51)	<b>11.3</b> (0.67)	<b>1.00</b> (0.00)	0.87 (0.03)
	squared loss	2.95 (0.83)	2.72 (0.80)	11.9 (0.84)	0.99 (0.04)	<b>0.35</b> (0.18)
	$\ell_1$ loss	4.41 (1.38)	4.05 (1.31)	13.3 (1.44)	0.91 (0.12)	0.44 (0.21)
	Huber loss	4.09 (1.19)	3.74 (1.15)	13.0 (1.25)	0.94 (0.10)	0.44 (0.19)
$t_3$	proposed loss	0.41 (0.14)	0.34 (0.11)	3.50 (0.66)	1.00 (0.00)	0.83 (0.05)
	Student's loss	<b>0.40</b> (0.12)	<b>0.33</b> (0.09)	<b>3.48</b> (0.66)	1.00 (0.00)	0.83 (0.04)
	squared loss	1.13 (0.63)	1.05 (0.62)	4.21 (1.00)	1.00 (0.00)	<b>0.23</b> (0.16)
	$\ell_1$ loss	0.66 (0.18)	0.59 (0.16)	3.74 (0.67)	1.00 (0.00)	0.47 (0.19)
	Huber loss	0.64 (0.19)	0.58 (0.16)	3.73 (0.68)	1.00 (0.00)	0.38 (0.17)
<i>Cauchy</i>	proposed loss	<b>2.56</b> (2.33)	<b>2.17</b> (1.92)	$>10^5$	<b>1.00</b> (0.00)	0.82 (0.20)
	Student's loss	2.79 (4.03)	2.31 (3.16)	$>10^5$	<b>1.00</b> (0.00)	0.81 (0.15)
	squared loss	10.0 (0.00)	10.0 (0.73)	$>10^5$	0.00 (0.00)	<b>0.00</b> (0.00)
	$\ell_1$ loss	7.97 (3.00)	7.85 (3.08)	$>10^5$	0.33 (0.44)	0.03 (0.08)
	Huber loss	8.02 (2.92)	7.89 (3.00)	$>10^5$	0.33 (0.43)	0.02 (0.07)

procedure is repeated 100 times, each with a different random partition of the data. The outcomes of these iterations are displayed in Table 9.

For this data set, we see that Huber and  $L_1$  Lasso methods return higher prediction errors compared to the standard lasso. On the other hand, our proposed method and heavy lasso method are the best methods. In terms of variable selection, all five methods consistently identified six common predictors 1371614, 1375833, 1377836, 1388491, 1389910, 1393736, highlighting its potential biological relevance.

## 6 Discussion and Conclusion

In this paper, we introduced the Exponential Lasso, a novel robust estimator for high-dimensional linear regression. Our goal was to address the well-known sensitivity of the classical Lasso's squared-error loss to outliers and heavy-tailed noise. By replacing the squared loss with an exponential-type loss function, our method successfully integrates the  $L_1$  penalty for sparse estimation with the principles of robust M-estimation, as the loss function smoothly and automatically downweights the influence of large residuals.

From a theoretical standpoint, we established strong non-asymptotic guarantees, proving that the Exponential Lasso achieves reliable estimation accuracy, matching the standard Lasso or Huber

Table 4: *Simulation results for various loss functions with Lasso penalization, under the setting  $p = 500, s^* = 10, n = 300$  and correlated design  $\rho_X = 0.5$ . The reported values are the mean across 100 simulation repetitions, with the standard deviation provided in parentheses. Bold font highlights the superior method. TPR: true positive rate; FDR: false discovery rate;  $MSPE_{test}$ : mean squared prediction error on testing data.*

Noise	Method (loss)	$\ \hat{\beta} - \beta_0\ _2^2$	$\ell(X^\top \beta_0)$	$MSPE_{test}$	TPR	FPR
$\mathcal{N}(0, 1)$	proposed loss	<b>0.20</b> (0.08)	<b>0.15</b> (0.04)	<b>1.17</b> (0.06)	1.00 (0.00)	0.80 (0.04)
	Student's loss	0.21 (0.08)	<b>0.15</b> (0.05)	<b>1.17</b> (0.07)	1.00 (0.00)	0.80 (0.05)
	squared loss	0.31 (0.13)	0.24 (0.07)	1.26 (0.09)	1.00 (0.00)	<b>0.22</b> (0.16)
	$\ell_1$ loss	0.44 (0.16)	0.32 (0.09)	1.35 (0.10)	1.00 (0.00)	0.35 (0.16)
	Huber loss	0.38 (0.15)	0.28 (0.08)	1.31 (0.10)	1.00 (0.00)	0.28 (0.16)
$\mathcal{N}(0, 3)$	proposed loss	2.47 (0.86)	2.08 (0.88)	11.2 (1.05)	0.99 (0.03)	0.84 (0.06)
	Student's loss	<b>1.83</b> (0.63)	<b>1.48</b> (0.57)	<b>10.6</b> (0.68)	<b>1.00</b> (0.00)	0.82 (0.05)
	squared loss	2.48 (0.84)	2.18 (0.83)	11.3 (1.00)	0.92 (0.08)	<b>0.17</b> (0.16)
	$\ell_1$ loss	3.00 (0.97)	2.77 (1.07)	11.9 (1.19)	0.89 (0.09)	0.27 (0.20)
	Huber loss	2.92 (0.93)	2.65 (0.99)	11.8 (1.16)	0.89 (0.09)	0.25 (0.18)
$t_3$	proposed loss	<b>0.34</b> (0.12)	<b>0.25</b> (0.08)	<b>3.16</b> (0.31)	1.00 (0.00)	0.79 (0.06)
	Student's loss	0.35 (0.14)	0.27 (0.12)	3.18 (0.32)	1.00 (0.00)	0.79 (0.06)
	squared loss	1.22 (0.93)	1.08 (1.20)	4.04 (1.27)	0.98 (0.08)	<b>0.10</b> (0.15)
	$\ell_1$ loss	0.63 (0.23)	0.51 (0.15)	3.43 (0.34)	1.00 (0.00)	0.25 (0.15)
	Huber loss	0.58 (0.22)	0.46 (0.15)	3.38 (0.32)	1.00 (0.00)	0.20 (0.16)
<i>Cauchy</i>	proposed loss	<b>2.85</b> (3.41)	<b>1.96</b> (2.13)	$>10^6$	<b>1.00</b> (0.00)	0.83 (0.16)
	Student's loss	4.64 (8.29)	2.84 (4.43)	$>10^6$	<b>1.00</b> (0.00)	0.82 (0.15)
	squared loss	9.92 (0.53)	18.3 (1.88)	$>10^6$	0.01 (0.10)	<b>0.00</b> (0.00)
	$\ell_1$ loss	6.33 (2.89)	10.1 (6.68)	$>10^6$	0.52 (0.38)	0.01 (0.05)
	Huber loss	6.39 (2.90)	10.2 (6.69)	$>10^6$	0.52 (0.38)	0.01 (0.05)

Lasso but under much milder assumptions that permit heavy-tailed noise. Computationally, the estimator's smooth and non-convex objective function is well-suited for a Majorization-Minimization (MM) algorithm. This framework is stable and efficient, iteratively solving a sequence of reweighted Lasso problems. Empirically, our extensive simulations demonstrated that the Exponential Lasso consistently outperforms competitors like the classical,  $L_1$  (LAD), and Huber Lasso, especially in settings with significant data contamination. Notably, it remained highly competitive even under standard Gaussian noise, suggesting a "premium" for its robustness. Our real-data application further validated these findings, confirming its practical relevance.

The primary advantage of the Exponential Lasso lies in its unique balance of efficiency and robustness. However, the method introduces the practical challenge of tuning two parameters: the regularization parameter  $\lambda$  and the robustness parameter  $\tau$ . Developing a data-driven, computationally efficient strategy for this joint tuning is a critical next step. Furthermore, the objective function's non-convexity means our MM algorithm guarantees convergence to only a stationary point, not necessarily the global optimum.

These limitations point toward clear avenues for future research. A more comprehensive study of parameter tuning is essential for broad practical adoption. Additionally, exploring initialization strategies or alternative global optimization algorithms could further bolster the method against the challenges of non-convexity. Finally, the robust exponential loss framework is highly adaptable. It could be extended to other high-dimensional problems, such as the Group Lasso, the Elastic Net,



Table 5: *Simulation results for various loss functions with Lasso penalization, under the setting  $p = 500, s^* = 10, n = 300$  and independent predictors. The outliers are increased by 10%, 20% and 30%. The reported values are the mean across 100 simulation repetitions, with the standard deviation provided in parentheses. Bold font highlights the best method. TPR: true positive rate; FDR: false discovery rate;  $MSPE_{test}$ : mean squared prediction error on testing data.*

Noise	Method (loss)	$\ \hat{\beta} - \beta_0\ _2^2$	$\ell(X^\top \beta_0)$	$MSPE_{test}$	TPR	FPR
$\mathcal{N}(0, 1)$ , 10% outliers	proposed loss	0.40 (0.37)	<b>0.33</b> (0.27)	1.41 (0.39)	<b>1.00</b> (0.00)	0.76 (0.17)
	Student's loss	<b>0.39</b> (0.17)	0.34 (0.16)	<b>1.40</b> (0.18)	<b>1.00</b> (0.00)	0.76 (0.16)
	squared loss	10.0 (0.00)	9.95 (0.83)	10.9 (0.26)	0.00 (0.00)	0.00 (0.00)
	$\ell_1$ loss	1.57 (0.41)	1.46 (0.34)	2.58 (0.43)	<b>1.00</b> (0.00)	0.11 (0.08)
	Huber loss	1.57 (0.35)	1.46 (0.29)	2.57 (0.37)	<b>1.00</b> (0.00)	<b>0.07</b> (0.08)
$\mathcal{N}(0, 1)$ , 20% outliers	proposed loss	<b>0.71</b> (0.74)	<b>0.61</b> (0.64)	<b>1.70</b> (0.72)	<b>1.00</b> (0.00)	0.80 (0.15)
	Student's loss	0.77 (0.58)	0.68 (0.64)	1.76 (0.57)	<b>1.00</b> (0.00)	0.78 (0.14)
	squared loss	10.0 (0.00)	10.0 (0.78)	11.0 (0.20)	0.00 (0.00)	0.00 (0.00)
	$\ell_1$ loss	2.67 (1.00)	2.47 (0.81)	3.68 (0.98)	0.99 (0.03)	0.15 (0.12)
	Huber loss	2.63 (0.80)	2.45 (0.61)	3.64 (0.78)	<b>1.00</b> (0.00)	<b>0.12</b> (0.10)
$\mathcal{N}(0, 1)$ , 30% outliers	proposed loss	<b>0.93</b> (0.94)	<b>0.82</b> (0.94)	<b>1.93</b> (0.95)	<b>1.00</b> (0.00)	0.79 (0.15)
	Student's loss	1.50 (2.02)	1.30 (1.99)	2.50 (2.02)	<b>1.00</b> (0.00)	0.78 (0.11)
	squared loss	10.0 (0.00)	10.0 (0.83)	11.0 (0.19)	0.00 (0.00)	0.00 (0.00)
	$\ell_1$ loss	4.44 (1.77)	4.18 (1.62)	5.45 (1.79)	0.88 (0.20)	0.15 (0.11)
	Huber loss	4.47 (1.57)	4.22 (1.45)	5.49 (1.60)	0.90 (0.16)	<b>0.11</b> (0.10)
$t_3$ , 10% outliers	proposed loss	0.71 (0.75)	0.57 (0.53)	3.63 (0.87)	1.00 (0.00)	0.81 (0.10)
	Student's loss	<b>0.64</b> (0.27)	<b>0.53</b> (0.22)	<b>3.56</b> (0.46)	1.00 (0.00)	0.78 (0.11)
	squared loss	10.0 (0.02)	9.82 (0.84)	12.9 (0.42)	0.00 (0.01)	0.00 (0.00)
	$\ell_1$ loss	2.12 (0.66)	1.93 (0.58)	5.05 (0.80)	1.00 (0.00)	0.16 (0.10)
	Huber loss	2.02 (0.58)	1.85 (0.51)	4.95 (0.74)	1.00 (0.00)	<b>0.12</b> (0.11)
$t_3$ , 20% outliers	proposed loss	1.11 (1.30)	0.94 (1.03)	4.09 (1.46)	<b>1.00</b> (0.00)	0.75 (0.20)
	Student's loss	<b>0.99</b> (0.54)	<b>0.85</b> (0.50)	<b>3.96</b> (0.93)	<b>1.00</b> (0.00)	0.75 (0.14)
	squared loss	10.0 (0.00)	9.98 (0.64)	12.9 (0.87)	0.00 (0.00)	0.00 (0.00)
	$\ell_1$ loss	3.48 (1.12)	3.24 (0.99)	6.46 (1.34)	0.97 (0.07)	0.17 (0.11)
	Huber loss	3.40 (1.06)	3.17 (0.94)	6.38 (1.26)	0.98 (0.05)	<b>0.15</b> (0.11)
$t_3$ , 30% outliers	proposed loss	<b>1.67</b> (1.60)	<b>1.48</b> (1.60)	<b>4.98</b> (3.59)	<b>1.00</b> (0.01)	0.81 (0.13)
	Student's loss	2.32 (2.48)	2.08 (2.63)	5.63 (3.97)	<b>1.00</b> (0.00)	0.81 (0.08)
	squared loss	10.0 (0.00)	9.79 (0.70)	13.2 (3.41)	0.00 (0.00)	0.00 (0.00)
	$\ell_1$ loss	6.90 (2.25)	6.50 (2.14)	10.1 (3.81)	0.60 (0.35)	0.11 (0.12)
	Huber loss	6.95 (2.30)	6.55 (2.19)	10.2 (3.84)	0.58 (0.34)	<b>0.10</b> (0.13)

the SCAD or MCP, or robust graphical model estimation, representing an exciting and practical path forward for analysis in data-rich but outlier-prone domains.

## Acknowledgments

The findings, interpretations, and conclusions expressed in this paper are entirely those of the author and do not reflect the views or positions of the Norwegian Institute of Public Health in any forms.

Table 6: *Simulation results for changing  $\tau$  with  $p = 120, s^* = 10, n = 100$  and independent predictors. The reported values are the mean across 100 simulation repetitions, with the standard deviation provided in parentheses. Bold font highlights the best method. TPR: true positive rate; FDR: false discovery rate;  $MSPE_{test}$ : mean squared prediction error on testing data.*

Noise	$\tau$	$\ \hat{\beta} - \beta_0\ _2^2$	$\ell(X^\top \beta_0)$	$MSPE_{test}$	TPR	FDR
$\mathcal{N}(0, 1)$	$\tau = 0.01$	<b>0.23</b> (0.07)	<b>0.19</b> (0.05)	<b>1.23</b> (0.07)	1.00 (0.00)	<b>0.84</b> (0.04)
	$\tau = 0.1$	0.24 (0.07)	0.20 (0.05)	1.24 (0.07)	1.00 (0.00)	0.85 (0.04)
	$\tau = 1$	3.62 (1.09)	2.55 (0.90)	4.62 (1.08)	1.00 (0.00)	0.97 (0.01)
	$\tau = 10$	9.99 (0.11)	9.78 (0.74)	10.9 (0.25)	0.10 (0.14)	0.75 (0.42)
$\mathcal{N}(0, 1)$ outliers 20%	$\tau = 0.01$	2.97 (1.35)	2.63 (1.27)	3.97 (1.35)	<b>1.00</b> (0.02)	0.88 (0.02)
	$\tau = 0.1$	<b>0.58</b> (0.58)	<b>0.50</b> (0.51)	<b>1.59</b> (0.59)	<b>1.00</b> (0.00)	0.78 (0.17)
	$\tau = 1$	5.24 (1.29)	4.26 (1.31)	6.25 (1.30)	0.96 (0.20)	0.95 (0.14)
	$\tau = 10$	10.0 (0.08)	9.91 (0.80)	11.0 (0.26)	0.07 (0.11)	0.70 (0.44)
$\mathcal{N}(0, 3)$	$\tau = 0.01$	<b>2.29</b> (0.47)	<b>2.00</b> (0.37)	<b>11.2</b> (0.46)	<b>1.00</b> (0.00)	0.89 (0.01)
	$\tau = 0.1$	3.39 (1.03)	2.94 (0.96)	12.3 (1.03)	0.99 (0.03)	0.89 (0.03)
	$\tau = 1$	9.21 (0.67)	8.64 (0.91)	18.2 (0.79)	0.74 (0.37)	0.93 (0.19)
	$\tau = 10$	10.0 (0.09)	9.76 (0.86)	19.0 (0.37)	0.05 (0.09)	0.70 (0.44)
$t_3$	$\tau = 0.01$	0.56 (0.18)	0.47 (0.14)	3.49 (0.36)	1.00 (0.00)	0.83 (0.05)
	$\tau = 0.1$	<b>0.42</b> (0.14)	<b>0.35</b> (0.10)	<b>3.35</b> (0.32)	1.00 (0.00)	0.84 (0.04)
	$\tau = 1$	4.98 (0.65)	3.79 (0.75)	7.93 (0.74)	1.00 (0.01)	0.97 (0.01)
	$\tau = 10$	9.97 (0.12)	9.84 (0.86)	12.98 (0.37)	0.08 (0.12)	0.59 (0.49)
<i>Cauchy</i>	$\tau = 0.01$	3.99 (4.65)	3.48 (3.91)	$>10^4$	<b>1.00</b> (0.00)	0.90 (0.03)
	$\tau = 0.1$	<b>1.66</b> (1.50)	<b>1.42</b> (1.18)	$>10^4$	<b>1.00</b> (0.00)	0.79 (0.19)
	$\tau = 1$	6.93 (1.31)	6.09 (1.50)	$>10^4$	0.91 (0.23)	0.96 (0.03)
	$\tau = 10$	10.0 (0.13)	9.98 (0.69)	$>10^4$	0.09 (0.11)	0.80 (0.38)

## Author contributions

I am the only author of this paper.

## Conflicts of interest/Competing interests

The author declares no potential conflict of interests.

## A Proof

### Proof of Theorem 1.

*Step A — Optimality and cone decomposition.*

The estimator  $\hat{\beta}$  satisfies, through its optimality definition,

$$L_\tau(\hat{\beta}) + \lambda \|\hat{\beta}\|_1 \leq L_\tau(\beta^*) + \lambda \|\beta^*\|_1.$$

Set  $\Delta := \hat{\beta} - \beta^*$ . Rearranging and using the standard bound  $\|\beta^*\|_1 - \|\beta^* + \Delta\|_1 \leq \|\Delta_S\|_1 - \|\Delta_{S^c}\|_1$  yields

$$L_\tau(\beta^* + \Delta) - L_\tau(\beta^*) \leq \lambda (\|\Delta_S\|_1 - \|\Delta_{S^c}\|_1). \quad (11)$$

Table 7: *Simulation results with increasing sparsity  $s^* = 4, 8, 16$  with  $p = 500, n = 300$  and independent predictors. The reported values are the mean across 100 simulation repetitions, with the standard deviation provided in parentheses. Bold font highlights the best method. TPR: true positive rate; FDR: false discovery rate;  $MSPE_{test}$ : mean squared prediction error on testing data.*

Noise	Method (loss)	$\ \beta - \beta_0\ _2^2$	$\ell(X^\top \beta_0)$	$MSPE_{test}$	TPR	FPR
$s^* = 4$						
$\mathcal{N}(0, 1)$ ,	proposed loss	<b>0.52</b> (0.85)	<b>0.49</b> (0.79)	<b>1.51</b> (0.85)	<b>1.00</b> (0.00)	0.76 (0.32)
20%	Student's loss	0.62 (1.01)	0.65 (1.13)	1.62 (1.00)	<b>1.00</b> (0.00)	0.79 (0.26)
outliers	squared loss	4.00 (0.00)	3.99 (0.31)	5.01 (0.10)	0.00 (0.00)	0.00 (0.00)
	$\ell_1$ loss	2.23 (0.56)	2.18 (0.50)	3.23 (0.56)	0.92 (0.20)	<b>0.00</b> (0.00)
	Huber loss	2.19 (0.54)	2.15 (0.48)	3.20 (0.54)	0.94 (0.16)	<b>0.00</b> (0.03)
$s^* = 8$						
$\mathcal{N}(0, 1)$ ,	proposed loss	<b>0.57</b> (0.61)	<b>0.51</b> (0.56)	<b>1.57</b> (0.63)	<b>1.00</b> (0.00)	0.80 (0.18)
20%	Student's loss	0.59 (0.52)	0.55 (0.54)	1.59 (0.51)	<b>1.00</b> (0.00)	0.78 (0.17)
outliers	squared loss	8.00 (0.00)	7.97 (0.58)	8.99 (0.19)	0.00 (0.00)	0.00 (0.00)
	$\ell_1$ loss	2.49 (0.63)	2.35 (0.55)	3.49 (0.64)	0.99 (0.04)	0.04 (0.07)
	Huber loss	2.47 (0.54)	2.34 (0.47)	3.47 (0.55)	<b>1.00</b> (0.02)	<b>0.03</b> (0.06)
$s^* = 16$						
$\mathcal{N}(0, 1)$ ,	proposed loss	<b>0.74</b> (0.33)	<b>0.56</b> (0.24)	<b>1.73</b> (0.33)	<b>1.00</b> (0.00)	0.77 (0.08)
20%	Student's loss	1.11 (0.54)	0.86 (0.42)	2.11 (0.54)	<b>1.00</b> (0.00)	0.75 (0.08)
outliers	squared loss	16.0 (0.00)	16.0 (1.21)	16.9 (0.32)	0.00 (0.00)	0.00 (0.00)
	$\ell_1$ loss	3.13 (0.92)	2.65 (0.67)	4.12 (0.93)	<b>1.00</b> (0.02)	0.37 (0.11)
	Huber loss	3.18 (1.00)	2.75 (0.72)	4.18 (1.01)	0.99 (0.02)	<b>0.28</b> (0.10)
$s^* = 4$						
$t_3$	proposed loss	<b>0.15</b> (0.06)	<b>0.14</b> (0.06)	<b>2.14</b> (0.13)	1.00 (0.00)	0.88 (0.04)
20%	Student's loss	<b>0.15</b> (0.06)	<b>0.14</b> (0.05)	<b>2.14</b> (0.13)	1.00 (0.00)	0.87 (0.06)
outliers	squared loss	0.52 (0.34)	0.51 (0.31)	2.51 (0.37)	1.00 (0.00)	<b>0.05</b> (0.12)
	$\ell_1$ loss	0.31 (0.13)	0.30 (0.12)	2.30 (0.18)	1.00 (0.00)	0.26 (0.23)
	Huber loss	0.31 (0.12)	0.30 (0.11)	2.30 (0.16)	1.00 (0.00)	0.22 (0.20)
$s^* = 8$						
$t_3$	proposed loss	<b>0.27</b> (0.11)	<b>0.24</b> (0.09)	2.28 (0.15)	1.00 (0.00)	0.84 (0.05)
20%	Student's loss	<b>0.27</b> (0.10)	<b>0.24</b> (0.08)	<b>2.27</b> (0.14)	1.00 (0.00)	0.84 (0.05)
outliers	squared loss	0.66 (0.35)	0.63 (0.31)	2.66 (0.36)	1.00 (0.00)	<b>0.21</b> (0.17)
	$\ell_1$ loss	0.51 (0.18)	0.47 (0.16)	2.51 (0.21)	1.00 (0.00)	0.37 (0.17)
	Huber loss	0.47 (0.16)	0.44 (0.14)	2.47 (0.18)	1.00 (0.00)	0.31 (0.16)
$s^* = 16$						
$t_3$	proposed loss	<b>0.47</b> (0.16)	0.44 (0.14)	<b>2.47</b> (0.18)	1.00 (0.00)	0.31 (0.16)
20%	Student's loss	0.57 (0.14)	<b>0.43</b> (0.10)	2.57 (0.19)	1.00 (0.00)	0.81 (0.04)
outliers	squared loss	1.02 (0.43)	0.84 (0.35)	3.02 (0.43)	1.00 (0.00)	<b>0.47</b> (0.10)
	$\ell_1$ loss	0.87 (0.29)	0.67 (0.22)	2.87 (0.33)	1.00 (0.00)	0.67 (0.09)
	Huber loss	0.79 (0.24)	0.64 (0.18)	2.80 (0.28)	1.00 (0.00)	0.55 (0.10)

Step B — Taylor expansion and LRSC lower bound.

By Taylor expansion in direction  $\Delta$  (componentwise),

$$L_\tau(\beta^* + \Delta) - L_\tau(\beta^*) = \langle \nabla L_\tau(\beta^*), \Delta \rangle + \frac{1}{2} \Delta^\top \left( \frac{1}{n} \sum_{i=1}^n \psi'_\tau(\xi_i) x_i x_i^\top \right) \Delta,$$

Table 8: Results on prediction errors and selected variables for the NCI-60 cancer cell line data.

Method	MSPE <sub>test</sub>	model size
proposed loss	0.398 (0.266)	39
Student's loss	0.395 (0.259)	73
squared loss	0.509 (0.225)	26
$\ell_1$ loss	0.474 (0.350)	19
Huber loss	0.479 (0.334)	5

Table 9: Results on prediction errors and selected variables for the gene expression TRIM32 data.

Method	MSPE <sub>test</sub>	model size
proposed loss	0.351 (0.102)	85
Student's loss	0.353 (0.096)	56
squared loss	0.472 (0.266)	29
$\ell_1$ loss	0.543 (0.336)	24
Huber loss	0.540 (0.313)	19

where  $\psi'_\tau(u) = e^{-\frac{\tau}{2}u^2}(1 - \tau u^2)$  and each  $\xi_i$  lies on the line segment between  $\varepsilon_i$  and  $\varepsilon_i - x_i^\top \Delta$ .

Let  $G = \{i : |\varepsilon_i| \leq c/2\}$ . If  $\|\Delta\|_2 \leq r$  and  $r$  is small enough so that  $|x_i^\top \Delta| \leq c/2$  for all  $i$ , then for any  $i \in G$  we have  $|\xi_i| \leq c$  and hence  $\psi'_\tau(\xi_i) \geq \underline{\gamma}$ . Therefore

$$\Delta^\top \left( \frac{1}{n} \sum_{i=1}^n \psi'_\tau(\xi_i) x_i x_i^\top \right) \Delta \geq \underline{\gamma} \cdot \frac{|G|}{n} \cdot \frac{1}{|G|} \sum_{i \in G} (x_i^\top \Delta)^2.$$

By Assumption 1 (restricted eigenvalue) applied to the same cone, for every  $\Delta$  in the cone we have  $\frac{1}{n} \sum_{i=1}^n (x_i^\top \Delta)^2 \geq \phi_{\min} \|\Delta\|_2^2$ . Restricting to the subset  $G$  can only decrease the quadratic form, but we lower bound it by using the trivial relation

$$\frac{1}{|G|} \sum_{i \in G} (x_i^\top \Delta)^2 \geq \frac{1}{n} \sum_{i=1}^n (x_i^\top \Delta)^2 \geq \phi_{\min} \|\Delta\|_2^2,$$

so

$$\Delta^\top \left( \frac{1}{n} \sum_{i=1}^n \psi'_\tau(\xi_i) x_i x_i^\top \right) \Delta \geq \underline{\gamma} \cdot \frac{|G|}{n} \cdot \phi_{\min} \|\Delta\|_2^2.$$

By Hoeffding's inequality for the binomial variable  $|G| \sim \text{Bin}(n, p_0)$ ,

$$\mathbb{P}\left(|G| \leq \frac{np_0}{2}\right) \leq \exp\left(-\frac{np_0^2}{8}\right).$$

Hence with probability at least  $1 - \exp(-np_0^2/8)$  we have  $|G|/n \geq p_0/2$ , and therefore the quadratic remainder satisfies the LRSC inequality

$$L_\tau(\beta^* + \Delta) - L_\tau(\beta^*) - \langle \nabla L_\tau(\beta^*), \Delta \rangle \geq \kappa \|\Delta\|_2^2, \quad (12)$$

with

$$\kappa = \frac{p_0}{2} \underline{\gamma} \phi_{\min}.$$

*Step C — Stochastic bound for the gradient sup-norm.*  
 Compute the gradient at  $\beta^*$ :

$$\nabla L_\tau(\beta^*) = -\frac{1}{n} \sum_{i=1}^n \psi_\tau(\varepsilon_i) x_i, \quad \psi_\tau(u) := ue^{-\frac{\tau}{2}u^2}.$$

By symmetry of  $\varepsilon_i$  we have  $\mathbb{E}[\psi_\tau(\varepsilon_i)] = 0$ , hence each coordinate  $[\nabla L_\tau(\beta^*)]_j = -\frac{1}{n} \sum_{i=1}^n Z_{ij}$  with  $Z_{ij} := \psi_\tau(\varepsilon_i)x_{ij}$  mean zero and bounded:  $|Z_{ij}| \leq KB_\tau$  where  $B_\tau = 1/\sqrt{e\tau}$ . Further,  $\text{Var}(Z_{ij}) \leq \mathbb{E}[Z_{ij}^2] \leq K^2 B_\tau^2$ .

Apply Bernstein's inequality coordinatewise: for any  $t > 0$ ,

$$\mathbb{P}\left(|[\nabla L_\tau(\beta^*)]_j| \geq t\right) \leq 2 \exp\left(-\frac{nt^2/2}{K^2 B_\tau^2 + (KB_\tau)t/3}\right).$$

Set  $t = t_0 := B_\tau K \sqrt{\frac{2 \log(2p/\delta)}{n}}$ . For this choice the denominator satisfies  $K^2 B_\tau^2 + (KB_\tau)t_0/3 \leq 2K^2 B_\tau^2$  (for  $n$  large enough; more generally the exact Bernstein algebra yields the same type of bound). Hence

$$\mathbb{P}\left(|[\nabla L_\tau(\beta^*)]_j| \geq t_0\right) \leq \frac{\delta}{p}.$$

By union bound over  $j = 1, \dots, p$ , with probability at least  $1 - \delta$ ,

$$\|\nabla L_\tau(\beta^*)\|_\infty \leq t_0 = B_\tau K \sqrt{\frac{2 \log(2p/\delta)}{n}}.$$

Therefore with the choice  $\lambda$  in (4) we have  $\|\nabla L_\tau(\beta^*)\|_\infty \leq \lambda/4$ .

*Step D — Combine LRSC and stochastic control to get the rate.*

From (11) and (12),

$$\langle \nabla L_\tau(\beta^*), \Delta \rangle + \kappa \|\Delta\|_2^2 \leq \lambda(\|\Delta_S\|_1 - \|\Delta_{S^c}\|_1).$$

Use  $|\langle \nabla L_\tau(\beta^*), \Delta \rangle| \leq \|\nabla L_\tau(\beta^*)\|_\infty \|\Delta\|_1 \leq (\lambda/4) \|\Delta\|_1$  to get

$$\kappa \|\Delta\|_2^2 \leq \lambda(\|\Delta_S\|_1 - \|\Delta_{S^c}\|_1) + \frac{\lambda}{4} \|\Delta\|_1 = \frac{5\lambda}{4} \|\Delta_S\|_1 - \frac{3\lambda}{4} \|\Delta_{S^c}\|_1.$$

Discard the negative term and apply  $\|\Delta_S\|_1 \leq \sqrt{s} \|\Delta\|_2$ :

$$\kappa \|\Delta\|_2^2 \leq \frac{5\lambda}{4} \sqrt{s} \|\Delta\|_2 \implies \|\Delta\|_2 \leq \frac{5\lambda\sqrt{s}}{4\kappa} \leq \frac{12\lambda\sqrt{s}}{\kappa},$$

where the last inequality is numeric (one may sharpen constants; we keep a conservative factor 12 to account for small-sample Bernstein second-order terms). The bound on  $\ell_1$ -error follows from the cone relation:  $\|\Delta_{S^c}\|_1 \leq 3\|\Delta_S\|_1$ , which implies  $\|\Delta\|_1 \leq 4\|\Delta_S\|_1 \leq 4\sqrt{s}\|\Delta\|_2$ .

*Step E — Probability union.*

The two high-probability events used are:

- gradient sup-norm event: probability at least  $1 - \delta$ ;
- good indices proportion event: probability at least  $1 - \exp(-np_0^2/8)$ .

Union bound gives the claimed probability  $1 - \delta - \exp(-np_0^2/8)$ . (We wrote  $1 - \delta - 2\exp(-np_0^2/8)$  in the theorem to be conservative in accounting for other small-probability concentration steps; the constants may be tightened.)

This completes the proof of Theorem 1.  $\square$

**Proof of Theorem 2.** Point (i): By construction of the majorizer, for any  $\beta$ ,

$$L(\beta) \leq Q^{(t)}(\beta) - \lambda \|\beta\|_1 + C^{(t)}.$$

Because  $Q^{(t)}(\cdot)$  coincides with  $L(\cdot) + \lambda \|\cdot\|_1$  at  $\beta = \beta^{(t)}$  (the majorizer is tight at the expansion point), we have  $Q^{(t)}(\beta^{(t)}) = F(\beta^{(t)})$ . Let  $\beta^{(t+1)}$  be the minimizer of  $Q^{(t)}$ . Then

$$F(\beta^{(t+1)}) \leq Q^{(t)}(\beta^{(t+1)}) \leq Q^{(t)}(\beta^{(t)}) = F(\beta^{(t)}).$$

The first inequality follows from  $L(\cdot) \leq \text{surrogate} - C^{(t)}$ ; the second because  $\beta^{(t+1)}$  minimizes  $Q^{(t)}$ . This proves the descent property.

Point (ii): Since  $\ell_i(\beta) \geq 0$  for all  $i$  and  $\lambda \|\beta\|_1 \geq 0$ , we have  $F(\beta) \geq 0$ . From point (i) the nonincreasing bounded-below sequence  $F(\beta^{(t)})$  converges to some finite limit  $F^* \geq 0$ .

Point (iii): Because  $F(\beta^{(t)}) \leq F(\beta^{(0)})$  for all  $t$ , all iterates belong to the sublevel set  $\{\beta : F(\beta) \leq F(\beta^{(0)})\}$ . If  $\lambda > 0$ , then  $\|\beta\|_1 \leq F(\beta^{(0)})/\lambda$  on this sublevel set, which implies boundedness. If  $\lambda = 0$  and  $X$  has full column rank, then the loss  $L(\beta)$  is coercive (grows at least quadratically) and hence the sublevel set is bounded. Thus under these mild alternatives the iterates are bounded.  $\square$

**Proof of Theorem 3.** Let  $\beta^{(t_k)} \rightarrow \beta^*$  be any convergent subsequence; such subsequences exist because  $\{\beta^{(t)}\}$  is bounded. Denote by  $\tilde{v}_i^{(t)} = \exp(-\frac{\tau}{2}r_i(\beta^{(t)})^2)$  the weights used in the surrogate at step  $t$ . Because the mapping  $\beta \mapsto \tilde{v}_i(\beta)$  is continuous,  $\tilde{v}_i^{(t_k)} \rightarrow \tilde{v}_i^* := \exp(-\frac{\tau}{2}r_i(\beta^*)^2)$  as  $k \rightarrow \infty$  for each  $i$ .

By the optimality of  $\beta^{(t_k+1)}$  for the convex surrogate  $Q^{(t_k)}$  we have the KKT condition for the weighted Lasso:

$$0 \in -\frac{1}{n}X^\top (\tilde{V}^{(t_k)}(y - X\beta^{(t_k+1)})) + \lambda \partial \|\beta^{(t_k+1)}\|_1, \quad (13)$$

where  $\tilde{V}^{(t_k)} = \text{diag}(\tilde{v}_1^{(t_k)}, \dots, \tilde{v}_n^{(t_k)})$ . Rearranging,

$$\frac{1}{n}X^\top \tilde{V}^{(t_k)}(y - X\beta^{(t_k+1)}) \in \lambda \partial \|\beta^{(t_k+1)}\|_1.$$

Because  $\beta^{(t_k+1)}$  and  $\tilde{V}^{(t_k)}$  are bounded and the functions are continuous, passing to the limit along the subsequence yields

$$\frac{1}{n}X^\top \tilde{V}^*(y - X\beta^*) \in \lambda \partial \|\beta^*\|_1.$$

It remains to show that the limiting left-hand side equals  $\nabla L(\beta^*)$ . Direct differentiation of  $L(\beta)$  gives

$$\nabla L(\beta) = -\frac{1}{n} \sum_{i=1}^n \exp\left(-\frac{\tau}{2}r_i(\beta)^2\right) x_i r_i(\beta) = -\frac{1}{n}X^\top (\tilde{V}(\beta)(y - X\beta)),$$

so indeed

$$-\nabla L(\beta^*) = \frac{1}{n}X^\top \tilde{V}^*(y - X\beta^*).$$

Combining with the limit KKT condition gives  $0 \in \nabla L(\beta^*) + \lambda \partial \|\beta^*\|_1$ , i.e.  $\beta^*$  is stationary for  $F$ .  $\square$

## References

- [1] Alfons, A. (2021). robustHD: An R package for robust regression with high-dimensional data. *Journal of Open Source Software*, 6(67):3786.
- [2] Bellec, P. C., Lecué, G., and Tsybakov, A. B. (2018). Slope meets lasso: improved oracle bounds and optimality. *The Annals of Statistics*, 46(6B):3603–3642.
- [3] Bühlmann, P. and Van De Geer, S. (2011). *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media.
- [4] Bunea, F., Tsybakov, A., and Wegkamp, M. (2007). Sparsity oracle inequalities for the lasso. *Electronic Journal of Statistics*, 1:169–194.
- [5] Chang, L., Roberts, S., and Welsh, A. (2018). Robust lasso regression using Tukey’s biweight criterion. *Technometrics*, 60(1):36–47.
- [6] Chiang, A. P., Beck, J. S., Yen, H.-J., Tayeh, M. K., Scheetz, T. E., Swiderski, R. E., Nishimura, D. Y., Braun, T. A., Kim, K.-Y. A., Huang, J., et al. (2006). Homozygosity mapping with SNP arrays identifies TRIM32, an E3 ubiquitin ligase, as a Bardet–Biedl syndrome gene (BBS11). *Proceedings of the National Academy of Sciences*, 103(16):6287–6292.
- [7] Donoho, D. L. and Johnstone, I. M. (1994). Ideal spatial adaptation by wavelet shrinkage. *biometrika*, 81(3):425–455.
- [8] Friedman, J., Hastie, T., and Tibshirani, R. (2010). Regularization Paths for Generalized Linear Models via Coordinate Descent. *Journal of Statistical Software*, 33(1):1–22.
- [9] Giraud, C. (2021). *Introduction to high-dimensional statistics*. Chapman and Hall/CRC.
- [10] Hampel, F. R. (1974). The influence curve and its role in robust estimation. *Journal of the american statistical association*, 69(346):383–393.
- [11] Hastie, T., Tibshirani, R., Friedman, J. H., and Friedman, J. H. (2009). *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer.
- [12] He, R., Hu, B.-G., Zheng, W.-S., and Kong, X.-W. (2011). Robust principal component analysis based on maximum correntropy criterion. *IEEE Transactions on Image Processing*, 20(6):1485–1494.
- [13] Huang, J., Horowitz, J. L., and Wei, F. (2010). Variable selection in nonparametric additive models. *Annals of statistics*, 38(4):2282.
- [14] Iqbal, A. and Seghouane, A.-K. (2019). An  $\alpha$ -divergence-based approach for robust dictionary learning. *IEEE Transactions on Image Processing*, 28(11):5729–5739.
- [15] Lecué, G. and Lerasle, M. (2020). Robust machine learning by median-of-means: theory and practice. *The Annals of Statistics*, 48(2):906–931.
- [16] Liu, W., Pokharel, P. P., and Principe, J. C. (2007). Correntropy: Properties and applications in non-gaussian signal processing. *IEEE Transactions on signal processing*, 55(11):5286–5298.

- [17] Loh, P.-L. (2017). Statistical consistency and asymptotic normality for high-dimensional robust M-estimators. *The Annals of Statistics*, 45(2):866.
- [18] Loh, P.-L. (2021). Scale calibration for high-dimensional robust regression. *Electronic Journal of Statistics*, 15(2):5933–5994.
- [19] Loh, P.-L. (2024). A theoretical review of modern robust statistics. *Annual Review of Statistics and Its Application*, 12.
- [20] Lounici, K. (2008). Sup-norm convergence rate and sign concentration property of Lasso and Dantzig estimators. *Electronic Journal of Statistics*, 2:90–102.
- [21] Mai, T. T. (2025a). A sparse PAC-Bayesian approach for high-dimensional quantile prediction. *Statistics and Computing*, 35(4):93.
- [22] Mai, T. T. (2025b). Heavy Lasso: sparse penalized regression under heavy-tailed noise via data-augmented soft-thresholding. *arXiv preprint arXiv:2506.07790*.
- [23] Reinhold, W. C., Sunshine, M., Liu, H., Varma, S., Kohn, K. W., Morris, J., Doroshow, J., and Pommier, Y. (2012). Cellminer: a web-based suite of genomic and pharmacologic tools to explore transcript and drug patterns in the nci-60 cell line set. *Cancer research*, 72(14):3499–3511.
- [24] Rejchel, W. and Bogdan, M. (2020). Rank-based lasso-efficient methods for high-dimensional robust model selection. *Journal of Machine Learning Research*, 21(244):1–47.
- [25] Rekavandi, A. M., Seghouane, A.-K., and Evans, R. J. (2021). Robust subspace detectors based on  $\alpha$ -divergence with application to detection in imaging. *IEEE Transactions on Image Processing*, 30:5017–5031.
- [26] Sardy, S., Tseng, P., and Bruce, A. (2001). Robust wavelet denoising. *IEEE transactions on signal processing*, 49(6):1146–1152.
- [27] Scheetz, T. E., Kim, K.-Y. A., Swiderski, R. E., Philp, A. R., Braun, T. A., Knudtson, K. L., Dorrance, A. M., DiBona, G. F., Huang, J., Casavant, T. L., et al. (2006). Regulation of gene expression in the mammalian eye and its relevance to eye disease. *Proceedings of the National Academy of Sciences*, 103(39):14429–14434.
- [28] Smucler, E. and Yohai, V. J. (2017). Robust and sparse estimators for linear regression models. *Computational Statistics & Data Analysis*, 111:116–130.
- [29] Song, Y., Liang, X., Zhu, Y., and Lin, L. (2021). Robust variable selection with exponential squared loss for the spatial autoregressive model. *Computational Statistics & Data Analysis*, 155:107094.
- [30] Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288.
- [31] Wang, H., Li, G., and Jiang, G. (2007). Robust regression shrinkage and consistent variable selection through the lad-lasso. *Journal of Business & Economic Statistics*, 25(3):347–355.



- [32] Wang, L., Peng, B., Bradic, J., Li, R., and Wu, Y. (2020). A tuning-free robust and efficient approach to high-dimensional regression. *Journal of the American Statistical Association*, 115(532):1700–1714.
- [33] Wang, X., Jiang, Y., Huang, M., and Zhang, H. (2013). Robust variable selection with exponential squared loss. *Journal of the American Statistical Association*, 108(502):632–643.
- [34] Wang, X., Shao, J., Wu, J., and Zhao, Q. (2023). Robust variable selection with exponential squared loss for partially linear spatial autoregressive models. *Annals of the Institute of Statistical Mathematics*, 75(6):949–977.
- [35] Yi, C. and Huang, J. (2017). Semismooth newton coordinate descent algorithm for elastic-net penalized huber loss regression and quantile regression. *Journal of Computational and Graphical Statistics*, 26(3):547–557.
- [36] Zhao, P. and Yu, B. (2006). On model selection consistency of lasso. *Journal of Machine learning research*, 7(Nov):2541–2563.
- [37] Zhu, J., Wang, X., Hu, L., Huang, J., Jiang, K., Zhang, Y., Lin, S., and Zhu, J. (2022). abess: a fast best-subset selection library in python and R. *Journal of Machine Learning Research*, 23(202):1–7.
- [38] Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American statistical association*, 101(476):1418–1429.