

Weighted L^1 and L^0 Regularization Using Proximal Operator Splitting Methods

Zewude A. Berkessa

*Research Unit of Mathematical Sciences
University of Oulu*

zewude.berkessa@oulu.fi

Patrik Waldmann

*Research Unit of Mathematical Sciences
University of Oulu*

patrik.waldmann@oulu.fi

Reviewed on OpenReview: <https://openreview.net/forum?id=9m2k96cDMK>

Abstract

This paper develops a joint weighted L^1 - and L^0 -norm (WL1L0) regularization method by leveraging proximal operators and translation mapping techniques to mitigate the bias introduced by the L^1 -norm in applications to high-dimensional data. A weighting parameter α is incorporated to control the influence of both regularizers. Our broadly applicable model is nonconvex and nonsmooth, but we show convergence for the alternating direction method of multipliers (ADMM) and the strictly contractive Peaceman–Rachford splitting method (SCPRSM). Moreover, we evaluate the effectiveness of our model on both simulated and real high-dimensional genomic datasets by comparing with adaptive versions of the least absolute shrinkage and selection operator (LASSO), elastic net (EN), smoothly clipped absolute deviation (SCAD) and minimax concave penalty (MCP). The results show that WL1L0 outperforms the LASSO, EN, SCAD and MCP by consistently achieving the lowest mean squared error (MSE) across all datasets, indicating its superior ability to handling large high-dimensional data. Furthermore, the WL1L0-SCPRSM also achieves the sparsest solution. Julia code for the WL1L0-ADMM and WL1L0-SCPRSM is available at <https://github.com/ZewAB/WL1L0-ADMM-and-SCPRSM>.

1 Introduction

High-dimensional statistics is a rapidly growing field of research that focuses on statistical analysis in the presence of a large number of variables or predictors (p), often much larger than the sample size (n). For example, high-throughput measurements in genomics contain thousands or millions of variables, such as single nucleotide polymorphism (SNP) markers and gene expression data for each individual. In such settings, traditional statistical methods often fail due to issues like overfitting, multicollinearity and computational complexity. In recent years, a number of regularization methods have been developed that impose a penalty on the size of the regression coefficients, which encourages sparsity and reduces the number of variables in the model (Fan et al., 2011; Fan & Lv, 2010; Heinze et al., 2018). Sparse learning techniques are essential in analyzing high-dimensional data for increased prediction accuracy, reduced computational complexity and enhanced interpretability of the results (Bühlmann & Van De Geer, 2011; Giraud, 2015; Wainwright, 2019).

Among various sparsity-inducing methods, the L^1 regularizer (also known as the least absolute shrinkage and selection operator (LASSO)) stands out for its convex nature and computational efficiency (Tibshirani, 1996). It adds a penalty term to the loss function proportional to the absolute value of the regression coefficients (L^1 -norm), which tends to shrink the coefficients towards zero and force some coefficients to exactly zero. Shrinking these coefficients helps to avoid overfitting, which can happen when a model memorizes the training data too well and does not perform well on new data. The LASSO improves the accuracy of predictions on

unseen data by simultaneously selecting features and mitigating overfitting. However, when coefficients are being shrunk, bias is introduced due to the bias-variance trade-off that is unavoidable in statistical learning. Furthermore, the LASSO tends to favor keeping larger coefficients over smaller ones which can lead to a bias towards larger coefficients in the model estimation process (Hastie et al., 2015).

In the specific context of genomic data, where the goal is often to identify genes associated with certain traits or diseases, this inaccurate selection can lead to the inclusion of incorrect genes in the model. This also poses a risk in terms of impaired prediction, as the estimated coefficients of the selected genes contribute to predicting the trait of interest (Fan et al., 2014b; Fan & Li, 2001; Johnstone & Titterington, 2009; Toloşı & Lengauer, 2011). Hence, the LASSO requires the fulfillment of the irrepresentable condition to obtain valid estimations (Zhao & Yu, 2006). In cases where the underlying datasets fail to meet this condition, the LASSO method may not accurately select the appropriate variables, leading to incorrect discoveries and wrong conclusions. In practice, implementing the irrepresentable condition can be challenging. Studies show that nonconvex regularizers such as SCAD and MCP reduce bias and have better prediction properties than the L^1 regularizer (Bertsimas et al., 2020; Fan & Li, 2001; Zhang, 2010).

On the other hand, L^0 regularization, which is also known as best subset selection (Hocking & Leslie, 1967), directly penalizes the number of non-zero coefficients in the model. It encourages sparsity, meaning it tends to produce models with fewer non-zero coefficients without any shrinkage. This results in a model that only includes the most relevant variables, simplifying the model and potentially improving its predictive performance by reducing overfitting. However, finding the optimal subset of variables using the L^0 -norm is computationally expensive because the L^0 -norm is nonconvex. While L^1 regularization is commonly used because of its convex nature, the L^0 -norm is computationally expensive and often intractable, and hence not frequently used on large data sets (Hastie et al., 2020).

Another regularization method is L^2 , also known as ridge regularization, which shrinks the coefficients towards zero without eliminating any of them completely (Hoerl & Kennard, 1970). Unlike the LASSO, ridge regression produces dense estimated regression coefficients, which means it does not perform feature selection. It reduces the size of the coefficients but does not drive any of them to exactly zero. Furthermore, the elastic net (EN) is another regularization technique that combines the properties of ridge regression and LASSO regression (Zou & Hastie, 2005). It is particularly useful for datasets with many features, especially when some are highly correlated. The LASSO may struggle with grouped variable selection, often picking just one from correlated variables, whereas EN improves feature selection in such cases (Hastie et al., 2015).

In this paper, we propose combining L^1 and L^0 regularization into a method denoted WL1L0 for improved prediction and variable selection. L^1 regularization encourages sparsity by shrinking some coefficients to zero, which helps reduce overfitting and is computationally efficient due to its convex properties. On the other hand, L^0 regularization enforces strict sparsity, directly penalizing the number of non-zero coefficients (i.e., eliminating variables with negligible impact), leading to more interpretable models. This synergy offers a better balance between interpretability and predictive accuracy, particularly in high-dimensional settings like genomics. Furthermore, since L^0 is an unbiased estimator and L^1 often introduces biases in estimation, L^0 can be regarded as debiasing L^1 in this setting. We achieve this goal by using a common regularization parameter and introducing a weight parameter that balances the importance of the two regularization methods.

We address the computational challenges that arise from optimizing the L^0 -norm by using proximal splitting methods, translation mapping and the efficient optimization algorithms ADMM and SCP-RSM. The prediction and sparsity properties of the WL1L0 method is evaluated on one simulated and two real genomic datasets and compared with the popular LASSO, EN, SCAD and MCP regularizers.

2 Related Work

In the rapidly evolving landscape of technology and data, prediction has become a cornerstone for making informed decisions across various domains. Regularization techniques are pivotal in enhancing the performance and generalizability of predictive models, particularly when dealing with complex datasets and high-dimensional data. By imposing penalties on the model parameters, regularization helps prevent over-

fitting, ensuring that the model captures the underlying patterns in the data. In this section, we will review key related works on regularization methods, highlighting significant advancements and methodologies.

We start by introducing a standard regression model

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \boldsymbol{\epsilon}, \quad (1)$$

where $\mathbf{y} \in \mathbb{R}^n$ is the response vector, $\mathbf{X} \in \mathbb{R}^{n \times p}$ is the predictor matrix, $\mathbf{b} \in \mathbb{R}^p$ is the vector of regression coefficients, and $\boldsymbol{\epsilon} \in \mathbb{R}^n$ is a noise (error) vector. For a vector \mathbf{b} , we write the q -norm notation as

$$\|\mathbf{b}\|_q = \begin{cases} \sum_i \mathbf{1}(b_i \neq 0), & \text{if } q = 0, \\ (\sum_i |b_i|^q)^{1/q}, & \text{if } 0 < q < \infty, \\ \max_i |b_i|, & \text{if } q = \infty, \end{cases}$$

where $i = 1, \dots, p$. Here, the $\|\mathbf{b}\|_0$ is the L^0 -norm that is the number of nonzero elements in \mathbf{b} . It is noteworthy that the L^0 -norm does not meet the criteria of a norm, specifically lacking the homogeneity property (Beck, 2017). Despite this, the term is widely used in the literature, and for the sake of consistency, we will retain its adoption.

Ridge regression, also known as Tikhonov regularization was introduced by Hoerl & Kennard (1970), uses an L^2 penalty term that shrinks all the coefficients and reduces their magnitudes. The ridge regression can be formulated as

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_2^2, \quad (2)$$

where $\lambda > 0$ is regularization parameter need to be tuned. This method is particularly effective in addressing multicollinearity in linear regression models. However, it does not necessarily set any coefficients to zero. Hence, ridge regression does not produce a sparse solution of estimated coefficients. On the other hand, LASSO regression

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda \|\mathbf{b}\|_1, \quad (3)$$

incorporates an L^1 regularization penalty, which encourages sparsity in the solution by setting some coefficients exactly to zero (Tibshirani, 1996).

Other penalty functions are introduced to provide a balance between inducing sparsity and reducing estimation bias, aiming to solve the optimization problem as

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + P_\lambda(\mathbf{b}). \quad (4)$$

For example, the smoothly clipped absolute deviation (SCAD) penalty function was introduced by Fan & Li (2001) as an improvement over LASSO regularization, particularly for bias reduction. The SCAD penalty function is defined as

$$P_\lambda^{SCAD}(\mathbf{b}) = \begin{cases} \lambda |\mathbf{b}| & \text{if } |\mathbf{b}| \leq \lambda, \\ \frac{-|\mathbf{b}|^2 + 2a\lambda|\mathbf{b}| - \lambda^2}{2(a-1)} & \text{if } \lambda < |\mathbf{b}| \leq a\lambda, \\ \frac{(a+1)\lambda^2}{2} & \text{if } |\mathbf{b}| > a\lambda, \end{cases} \quad (5)$$

where $\lambda > 0$ and $a > 0$ are unknown parameters. Fan & Li (2001) suggested that $a = 3.7$ is a good choice for various problems, and λ needs to be tuned.

The minimax concave penalty (MCP) is another type of penalty function introduced by Zhang (2010). The MCP penalty function is defined as

$$P_\lambda^{MCP}(\mathbf{b}) = \begin{cases} \lambda |\mathbf{b}| - \frac{\mathbf{b}^2}{2a} & \text{if } |\mathbf{b}| \leq \lambda a, \\ \frac{a\lambda^2}{2} & \text{if } |\mathbf{b}| > \lambda a. \end{cases} \quad (6)$$

According to the estimation theorems of Zhang (2010), $a = 3$ is a good choice for MCP, and λ still needs to be tuned. MCP was developed to address the estimation bias of the LASSO and is generally easier to optimize computationally compared to SCAD.

Both SCAD and MCP aim to eliminate unimportant variables while preserving important ones, achieving the ‘oracle property’ as the sample size grows ($n \rightarrow \infty$). They both asymptotically select the correct model and produce normal, accurate coefficient estimates. MCP is effective with many sparse predictor groups but struggles with tightly clustered non-zero coefficients while SCAD has weaker grouping behavior compared to MCP (Ogutu & Piepho, 2014). We maintain the use of $a = 3.7$ for SCAD and $a = 3$ for MCP throughout the paper.

For the L^0 regularization (best subset selection (BSS)), $P_\lambda(\mathbf{b})$ can be written as

$$P_\lambda^{BSS}(\mathbf{b}) = \lambda \sum_i^p \mathbf{1}(b_i \neq 0). \quad (7)$$

Exact optimization of problem (4) with the L^0 -norm, as defined in (7), is challenging because incorporating (7) into the objective function results in a non-differentiable and non-convex problem. For example, Louizou et al. (2017) propose a method for optimizing a relaxed version of the L^0 norm for parametric models using a distribution called the hard concrete distribution (Maddison et al., 2016), which facilitates gradient-based optimization.

Yun et al. (2019) use a family of M -estimators with trimmed regularization for general high-dimensional problems. The trimmed regularization problem can be formulated for LASSO as

$$\begin{aligned} \hat{\mathbf{b}}, \hat{\boldsymbol{\pi}} &= \underset{\mathbf{b}, \boldsymbol{\pi}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda \sum_{i=1}^p \pi_i |b_i| \\ \text{subject to } & \mathbf{1}^\top \boldsymbol{\pi} \geq p - h, \\ & \boldsymbol{\pi} \in [0, 1], \end{aligned} \quad (8)$$

where h denotes the trimming parameter, which must be appropriately tuned, for instance, through cross-validation.

Another example is the elastic net (EN) regression which combines both L^1 and L^2 regularization penalties, providing a balanced approach to prediction accuracy on future data and model interpretation in linear regression models. It is formulated as

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda_1 \|\mathbf{b}\|_1 + \lambda_2 \|\mathbf{b}\|_2^2, \quad (9)$$

which has two regularization parameters λ_1 and λ_2 to tune (Zou & Hastie, 2005). The LAVA regression model is based on the splitting of the regression component into one sparse and one dense part $\mathbf{b} = \mathbf{c} + \mathbf{d}$ and thereby obtaining the following optimization problem

$$\hat{\mathbf{c}}, \hat{\mathbf{d}} = \underset{\mathbf{c}, \mathbf{d}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}(\mathbf{c} + \mathbf{d})\|_2^2 + \lambda_1 \|\mathbf{c}\|_1 + \lambda_2 \|\mathbf{d}\|_2^2, \quad (10)$$

where the resulting estimator $\hat{\mathbf{b}} = \hat{\mathbf{c}} + \hat{\mathbf{d}}$ (Chernozhukov et al., 2017). The key difference between EN and LAVA is that EN performs variable selection (i.e., is dominated by the L^1 -norm), whereas LAVA is always dense (i.e., is dominated by the L^2 -norm). Waldmann (2021) developed a proximal operator algorithm based on the LAVA regularization method that jointly performs L^1 - and L^2 -norm regularization.

Ziyin & Wang (2023) propose a method called *spread*, for optimizing generic differentiable objectives with an L^1 constraint using a reparametrization. The method is proposed to effectively bridge the gap between sparsity in deep learning and conventional statistical learning by providing a principled way to optimize L^1 constraints in complex nonlinear settings. For example, one can apply the *spread* parametrization to the sparse component of \mathbf{b}_s given the LASSO loss $\|\mathbf{y} - \mathbf{X}\mathbf{b}_s\|_2^2 + 2\kappa \|\mathbf{b}_s\|_1$. The equivalent *spread* loss is then

$$\hat{\mathbf{c}}, \hat{\mathbf{d}} = \underset{\mathbf{c}, \mathbf{d}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}(\mathbf{c} \odot \mathbf{d})\|_2^2 + \kappa (\|\mathbf{d}\|_2^2 + \|\mathbf{c}\|_2^2), \quad (11)$$

where $\hat{\mathbf{b}} = \hat{\mathbf{c}} \odot \hat{\mathbf{d}}$, \odot denotes the element-wise product, and κ is the L^1 regularization strength parameter that needs to be tuned to achieve the best sparsity-performance trade-off.

3 Theoretical Background

Large parts of the theory behind our approach follows from (Bertsekas, 2016) and (Beck, 2017). For an extended real-valued function $f : \mathbb{R}^p \rightarrow [-\infty, \infty]$, we define the following:

- (a) The domain of f is the set

$$\text{dom}(f) = \{\mathbf{b} \in \mathbb{R}^p : f(\mathbf{b}) < \infty\}.$$

- (b) f is proper if $\text{dom}(f) \neq \emptyset$ and f is never $-\infty$.

- (c) The epigraph of f is defined by

$$\text{epi}(f) = \{(\mathbf{b}, a) \in \mathbb{R}^p \times \mathbb{R} : f(\mathbf{b}) \leq a\}.$$

- (d) The function f is closed if its epigraph is closed.

- (e) f is called lower semicontinuous at $\mathbf{b} \in \mathbb{R}^p$ if

$$f(\mathbf{b}) \leq \liminf_{k \rightarrow \infty} f(\mathbf{b}^{(k)})$$

for any sequence $\{\mathbf{b}^{(k)}\}_{k \geq 1} \subseteq \mathbb{R}^p$ for which $\mathbf{b}^{(k)} \rightarrow \mathbf{b}$ as $k \rightarrow \infty$.

- (f) For any $\eta \in \mathbb{R}$, the η -level set of a function f is the set

$$\text{Lev}(f, \eta) = \{\mathbf{b} \in \mathbb{R}^p : f(\mathbf{b}) \leq \eta\}.$$

- (g) A proper function f is called coercive if

$$\lim_{\|\mathbf{b}\| \rightarrow \infty} f(\mathbf{b}) = \infty.$$

For any set $\mathbb{S} \subseteq \mathbb{R}^p$ and any point $\mathbf{b} \in \mathbb{R}^p$, the distance from \mathbf{b} to \mathbb{S} is defined as $D(\mathbf{b}, \mathbb{S}) := \inf\{\|\mathbf{m} - \mathbf{b}\|, \mathbf{m} \in \mathbb{S}\}$, and $D(\mathbf{b}, \mathbb{S}) = \infty$ for all \mathbf{b} when $\mathbb{S} = \emptyset$.

A proper closed and coercive function f attains its minimal value over \mathbb{S} for a nonempty closed set satisfying $\mathbb{S} \cap \text{dom}(f) = \emptyset$. Moreover, a closed coercive function possesses a minimizer on any closed set that has a nonempty intersection with the domain of the function (Beck, 2017). For an extended real-valued function $f : \mathbb{R}^p \rightarrow [-\infty, \infty]$, the following three claims are equivalent:

- i f is lower semicontinuous.
- ii f is closed.
- iii For any $\eta \in \mathbb{R}$, the level set

$$\text{Lev}(f, \eta) = \{\mathbf{b} \in \mathbb{R}^p : f(\mathbf{b}) \leq \eta\}$$

is closed.

The proof of these claims can be found in (Beck, 2017), see Theorem 2.6.

3.1 Subdifferentials of Nonconvex and Nonsmooth Functions

Subdifferentials are important in analyzing complex functions, especially when dealing with nonsmooth and nonconvex functions. Following Clarke et al. (2008) and Mordukhovich (2006), we explore subdifferentiability.

Let $g : \mathbb{R}^p \rightarrow (-\infty, +\infty]$ be a proper and lower semicontinuous function. Then

- (i) For a given $\mathbf{b} \in \text{dom } g$, the Fréchet subdifferential of g at \mathbf{b} , denoted by $\hat{\partial}g(\mathbf{b})$, is the set of all vectors $\mathbf{u} \in \mathbb{R}^p$ which satisfy

$$\liminf_{\mathbf{m} \rightarrow \mathbf{b}} \frac{g(\mathbf{m}) - g(\mathbf{b}) - \langle \mathbf{u}, \mathbf{m} - \mathbf{b} \rangle}{\|\mathbf{m} - \mathbf{b}\|} \geq 0,$$

and we set $\hat{\partial}g = \emptyset$ when $\mathbf{b} \notin \text{dom } g$.

- (ii) The limiting-subdifferential, or simply the subdifferential, of g at \mathbf{b} , written by $g(\mathbf{b})$, is defined by

$$\partial g(\mathbf{b}) := \{\mathbf{u} \in \mathbb{R}^p : \exists \mathbf{b}^{(k)} \rightarrow \mathbf{b}, g(\mathbf{b}^{(k)}) \rightarrow g(\mathbf{b}) \text{ and } \mathbf{u}^{(k)} \in \hat{\partial}g(\mathbf{b}^{(k)}) \xrightarrow{k \rightarrow \infty} \mathbf{u}\},$$

where $\hat{\partial}g(\mathbf{b}) \subset \partial g(\mathbf{b})$ for all $\mathbf{b} \in \mathbb{R}^p$.

- (iii) A point \mathbf{b}^* is called critical point or stationary point of g if it satisfies $0 \in \partial g(\mathbf{b}^*)$.

Please refer to Wu et al. (2021) for generalized subdifferentials of the L^0 , with its regular subdifferentials provided in Le (2013).

3.2 The Kurdyka–Łojasiewicz Inequality and its Property

The Kurdyka–Łojasiewicz (KŁ) inequality deals with the behavior of certain functions near their critical points. It is an important tool for analyzing the convergence of nonconvex nonsmooth optimization problems (Attouch et al., 2010; 2013; Bolte et al., 2014). We now review the KŁ property.

Let $g : \mathbb{R}^p \rightarrow (-\infty, +\infty]$ be a proper lower semicontinuous function. Then,

- (a) The function $g : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{+\infty\}$ is said to have the KŁ property at $\mathbf{b}^* \in \text{dom } \partial g$ if there exist $\eta \in (0, +\infty]$, a neighborhood U of \mathbf{b}^* , and a continuous concave function $\phi : [0, \eta) \rightarrow \mathbb{R}^+$ such that
- (i) $\phi(0) = 0$,
 - (ii) ϕ is continuously differentiable on $(0, \eta)$,
 - (iii) $\forall a \in (0, +\infty]$, $\phi'(a) > 0$,
- (IV) For all $\mathbf{b} \in U \cap \{g(\mathbf{b}^*) < g(\mathbf{b}) < g(\mathbf{b}^*) + \eta\}$, the KŁ property holds:

$$\phi'(g(\mathbf{b}) - g(\mathbf{b}^*))d(0, \partial g(\mathbf{b})) \geq 1.$$

- (b) Proper lower semicontinuous functions which satisfy the KŁ inequality at each point of $\text{dom } \partial g$ are called KŁ functions. Examples of KŁ functions include $\|\mathbf{b}\|_1$, $\|\mathbf{b}\|_0$, and $\|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2$. For more examples, please refer to (Attouch et al., 2010; 2013; Bolte et al., 2014; Yashtini, 2022).

3.3 Proximal Operators

Proximal operators are a fundamental concept in optimization, especially for problems involving non-smooth or non-convex functions, which are increasingly common in a wide range of real-world applications (Fukushima & Mine, 1981; Kaplan & Tichatschke, 1998; Parikh & Boyd, 2013). A proximal operator, denoted as $\text{prox}_f(\mathbf{u})$, aims to find a point closer to \mathbf{u} that also minimizes a specific objective function, $f(\mathbf{v})$ in a specific optimization subproblem. This subproblem is assumed to be more manageable to solve than the original problem. The proximal operator can be mathematically expressed as

$$\text{prox}_f(\mathbf{u}) = \underset{\mathbf{v}}{\text{argmin}} \{f(\mathbf{v}) + (1/2)\|\mathbf{v} - \mathbf{u}\|_2^2\}, \quad (12)$$

where \mathbf{u} and \mathbf{v} are vectors of length p . Here, $\text{prox}_f(\mathbf{u})$ is a point that compromises between minimizing f and being close to \mathbf{u} . Note that the right-hand side of (12) is strongly convex, hence there is a unique minimizer for every $\mathbf{u} \in \mathbb{R}^p$. Introducing the parameter $\gamma > 0$ that represents a trade-off parameter between the two terms \mathbf{v} and \mathbf{u} yields a scaled version of (12), in which $\frac{1}{2}$ is replaced by $\frac{1}{2\gamma}$. The proximal operator has useful properties (Beck, 2017), one of which is its behavior when applied to affine functions, as shown below.

Lemma 1 For any affine function $f(\mathbf{u}) = \langle \mathbf{m}, \mathbf{u} \rangle + a$, where $\mathbf{m} \in \mathbb{R}^p$ is a fixed vector and $a \in \mathbb{R}$, then for any $\mathbf{u} \in \mathbb{R}^p$, the proximal operator defined in (12) reduces to a simple translation of the vector \mathbf{u} by \mathbf{m} . Specifically,

$$\text{prox}_f(\mathbf{u}) = \mathbf{u} - \mathbf{m}, \quad (13)$$

which represents a translation mapping.

The proof of Lemma 1 is provided in Appendix B.1. In accordance with Lemma 1, one defines a translation function as a function that incorporates a standard additive term which is expressed as $\mathcal{T}_{\mathbf{m}}(\mathbf{u}) = f(\mathbf{u} + \mathbf{m}) - \mathbf{m}$. Another important property arises in the context of separable sum functions $f(\mathbf{u}, \mathbf{m}) = g(\mathbf{u}) + h(\mathbf{m})$, where the proximal operator is written as $\text{prox}_f(\mathbf{u}, \mathbf{m}) = \text{prox}_g(\mathbf{u}) + \text{prox}_h(\mathbf{m})$. For proximal operators in the framework of L^0 , please refer to (Attouch et al., 2013; Bolte et al., 2014; Beck, 2017).

4 Methodological Framework

To achieve higher sparsity than the EN, one can use the least squares loss function with ℓ_1 and ℓ_0 norm constraints which can be formulated as

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\text{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 \quad \text{s.t.} \quad \|\mathbf{b}\|_1 \leq t, \quad \|\mathbf{b}\|_0 \leq s, \quad (14)$$

where $\hat{\mathbf{b}}$ represents the estimate of the vector of regression coefficients, t is a constant threshold and s is the desired level of sparsity (i.e., the maximum number of nonzero coefficients). The optimization problem (14) can be expressed in the Lagrangian form as

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\text{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda_1 \|\mathbf{b}\|_1 + \lambda_2 \|\mathbf{b}\|_0, \quad (15)$$

where $\lambda_1 > 0$ and $\lambda_2 > 0$ are the regularization parameters.

We formulate our proposed method by building upon the flexible, penalty-based framework introduced in (15) that uses the parameters λ_1 and λ_2 to control both the size of the coefficients and the sparsity in a more flexible and nuanced manner, allowing for a broader range of model behaviors compared to the constant threshold and the strict subset selection enforced by formulation (14). Here, the exact relationship between t and λ_1 and between s and λ_2 is data-dependent.

We now extend (15) by introducing a weight parameter $\alpha \in (0, 1)$, a common regularization parameter $\lambda > 0$ and reformulate the problem as

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\text{argmin}} \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda (\alpha \|\mathbf{b}\|_1 + (1 - \alpha) \|\mathbf{b}\|_0). \quad (16)$$

Furthermore, we formulate the problem (16) following the LAVA method. Hence, we split the regression component \mathbf{b} into the sparse components \mathbf{c} and \mathbf{d} and separately assign penalties L^1 and L^0 to them and therefore obtain

$$\text{WL1L0: } \hat{\mathbf{c}}, \hat{\mathbf{d}} = \underset{\mathbf{c}, \mathbf{d}}{\text{argmin}} \|\mathbf{y} - \mathbf{X}(\mathbf{c} + \mathbf{d})\|_2^2 + \lambda (\alpha \|\mathbf{c}\|_1 + (1 - \alpha) \|\mathbf{d}\|_0). \quad (17)$$

It is here useful to point out that EN combines L^1 -norm and L^2 -norm regularization which leads to variable selection while being less sensitive to correlated predictors than the LASSO. In contrast, LAVA is dominated by the L^2 -norm, leading to dense models without feature selection. WL1L0 combines L^1 and L^0 regularization, providing stricter feature selection by explicitly controlling the number of non-zero coefficients, resulting in more precise sparsity than the EN. See Appendix A for additional insights into the problem and discussion.

5 Optimization Algorithms

Alternating direction method of multipliers (ADMM) is often used as a benchmark algorithm for splitting problems due to its efficiency and flexibility (Boyd et al., 2011). One of its notable strengths is its ability to

handle large-scale optimization problems by decomposing them into smaller, more manageable subproblems. This decomposition not only simplifies the problem-solving process but also allows for parallel processing of these subproblems. The strictly contractive Peaceman–Rachford splitting method (SCPRSM) is a variant of the classical Peaceman–Rachford splitting method (PRSM) (He et al., 2014). Similar to ADMM, PRSM is an operator splitting technique used to solve optimization problems. SCPRSM is a further extension that ensures convergence by imposing a strict contraction condition. Although ADMM is more widely used in practice, SCPRSM is a more specialized method that guarantees convergence through strict contraction. While SCPRSM emphasizes strict contraction, ADMM may exhibit slower convergence under certain conditions. Therefore, we implement our proposed method using both ADMM and SCPRSM frameworks based on the augmented Lagrangian method that combines the original objective function with the constraints of the optimization problem into a single function. Here, the augmented Lagrangian’s advantage lies in enabling the study of convergence for the proposed methods without requiring assumptions like strict convexity (Boyd et al., 2011).

The introduction of two variables (\mathbf{c} and \mathbf{d}) instead of one (\mathbf{b}) in (17) increases the dimensionality of the optimization problem, adding complexity to the theoretical analysis. Therefore, we first study the convergence properties of the problem in (16). We then reformulate problem (16) to establish its convergence in the ADMM and SCPRSM frameworks. The optimization model for (16) can be formulated as

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\operatorname{argmin}} \{f(\mathbf{b}) + g(\mathbf{b})\} \iff \hat{\mathbf{b}} = \underset{\mathbf{b}}{\operatorname{argmin}} \{f(\mathbf{b}) + g(\mathbf{u})\} \quad (18)$$

subject to $\mathbf{b} = \mathbf{u}$,

where $f(\mathbf{b}) = \|\mathbf{y} - \mathbf{X}(\mathbf{b})\|_2^2$ is the loss function and $g(\mathbf{b}) = \lambda(\alpha\|\mathbf{b}\|_1 + (1 - \alpha)\|\mathbf{b}\|_0)$ is the penalty function. We now write the augmented Lagrangian function corresponding to (18) as

$$L_\gamma(\mathbf{b}, \mathbf{u}, \mathbf{z}) = f(\mathbf{b}) + g(\mathbf{u}) + \mathbf{z}^T(\mathbf{b} - \mathbf{u}) + \frac{\gamma}{2}\|\mathbf{b} - \mathbf{u}\|_2^2, \quad (19)$$

where \mathbf{z} is a dual variable or Lagrange multiplier and $\gamma > 0$ is a learning rate. Here, \mathbf{b} and \mathbf{u} are the primal variables.

5.1 Method of Multipliers and ADMM Framework

The method of multipliers jointly minimizes the two primal variables whereas the ADMM efficiently solves optimization problems by alternately updating primal and dual variables, effectively decomposing complex problems into manageable subproblems (Boyd et al., 2011). A more convenient scaled form of (19) can be obtained by completing the square with the dual variable \mathbf{z} and the residual $\mathbf{b} - \mathbf{u}$ in the augmented Lagrangian. This allows the term $\mathbf{z}^T(\mathbf{b} - \mathbf{u}) + \frac{\gamma}{2}\|\mathbf{b} - \mathbf{u}\|_2^2$ to be rewritten as $\frac{\gamma}{2}\|\mathbf{b} - \mathbf{u} + \frac{1}{\gamma}\mathbf{z}\|_2^2 - \frac{\gamma}{2}\|\mathbf{z}\|_2^2$. Introducing the scaled dual variable $\mathbf{m} = \frac{1}{\gamma}\mathbf{z}$, the scaled form of the augmented Lagrangian becomes

$$L_\gamma(\mathbf{b}, \mathbf{u}, \mathbf{m}) = f(\mathbf{b}) + g(\mathbf{u}) + \frac{\gamma}{2}\|\mathbf{b} - \mathbf{u} + \mathbf{m}\|_2^2 - \frac{\gamma}{2}\|\mathbf{m}\|_2^2. \quad (20)$$

This scaled form is better suited for implementing ADMM and SCPRSM schemes with proximal operators (Parikh & Boyd, 2013). The method of multipliers for (20) can be written as

$$(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}) := \underset{\mathbf{b}, \mathbf{u}}{\operatorname{argmin}} L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}), \quad (21)$$

$$\mathbf{m}^{(k+1)} := \mathbf{m}^{(k)} + \mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)}. \quad (22)$$

The method of multipliers is generally not an implementable method since the primal update step (21) can be as hard to solve as the original problem (Beck, 2017; Boyd et al., 2011). To overcome this challenge, ADMM employs an iterative approach in the primal update step. In this approach, \mathbf{b} and \mathbf{u} are updated sequentially in an alternating fashion, which is why the method is called the alternating direction method of multipliers.

An iterative scheme for the ADMM associated with (20) becomes

$$\mathbf{b}^{(k+1)} := \underset{\mathbf{b}}{\operatorname{argmin}} L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}), \quad (23a)$$

$$\mathbf{u}^{(k+1)} := \underset{\mathbf{u}}{\operatorname{argmin}} L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}), \quad (23b)$$

$$\mathbf{m}^{(k+1)} := \mathbf{m}^{(k)} + \mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)}. \quad (23c)$$

5.2 SCPRSM Framework

The difference between ADMM and PRSM in terms of convergence can be explained through the contraction properties of their iterative sequences. The iterative sequence generated by ADMM is strictly contractive with respect to a given solution set, whereas the sequence generated by PRSM is contractive, but not strictly contractive (He et al., 2014; Corman & Yuan, 2014; He et al., 2002). To address the lack of strict contraction in PRSM, He et al. (2014) proposed incorporating a relaxation factor $r > 0$ into the Lagrange multiplier update steps, thus developing a strictly contractive Peaceman-Rachford splitting method (SCPRSM). This modification ensures that the iterative sequence becomes strictly contractive, improving convergence properties. The studies show that SCPRSM outperforms that ADMM generally leads to faster convergence compared to ADMM (Li & Yuan, 2015; Li et al., 2021). The iterative scheme of the SCPRSM associated with the augmented Lagrangian function (20) is written as

$$\mathbf{b}^{(k+1)} := \underset{\mathbf{b}}{\operatorname{argmin}} L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}), \quad (24a)$$

$$\mathbf{m}^{(k+\frac{1}{2})} := \mathbf{m}^{(k)} + r(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k)}), \quad (24b)$$

$$\mathbf{u}^{(k+1)} := \underset{\mathbf{u}}{\operatorname{argmin}} L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k+\frac{1}{2})}), \quad (24c)$$

$$\mathbf{m}^{(k+1)} := \mathbf{m}^{(k+\frac{1}{2})} + r(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)}), \quad (24d)$$

where the parameter $r \in (0, 1)$ is a relaxation factor. In addition to a relaxation factor r , an important distinction between SCPRSM and ADMM is the presence of an intermediate update for the multipliers, represented as $\mathbf{m}^{(k+\frac{1}{2})}$ in the SCPRSM scheme. This step ensures a balanced handling of the vectors \mathbf{b} and \mathbf{u} , thereby leading to a contractive iteration sequence that guarantees convergence to the solution of the original optimization problem (He et al., 2014; Li & Yuan, 2015; Peaceman & Rachford, 1955). The iterative scheme of ADMM (23a - 23c) can be further simplified as

$$\begin{aligned} \mathbf{b}^{(k+1)} &:= \underset{\mathbf{b}}{\operatorname{argmin}} \left\{ f(\mathbf{b}^{(k)}) + \frac{\gamma}{2} \|\mathbf{b}^{(k)} - (\mathbf{u}^{(k)} - \mathbf{m}^{(k)})\|_2^2 \right\}, \\ \mathbf{u}^{(k+1)} &:= \underset{\mathbf{u}}{\operatorname{argmin}} \left\{ g(\mathbf{u}^{(k)}) + \frac{\gamma}{2} \|\mathbf{u}^{(k)} - (\mathbf{b}^{(k+1)} + \mathbf{m}^{(k)})\|_2^2 \right\}, \\ \mathbf{m}^{(k+1)} &:= \mathbf{m}^{(k)} + \mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)}. \end{aligned} \quad (25)$$

With proximal operators, we can now rewrite (25) as

$$\begin{aligned} \mathbf{b}^{(k+1)} &:= \operatorname{prox}_{f_\gamma}(\mathbf{u}^{(k)} - \mathbf{m}^{(k)}), \\ \mathbf{u}^{(k+1)} &:= \operatorname{prox}_{g_\gamma}(\mathbf{b}^{(k+1)} + \mathbf{m}^{(k)}), \\ \mathbf{m}^{(k+1)} &:= \mathbf{m}^{(k)} + \mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)}. \end{aligned} \quad (26)$$

Similarly, the proximal version of Equations (24a - 24d) can be written as

$$\begin{aligned} \mathbf{b}^{(k+1)} &:= \operatorname{prox}_{f_\gamma}(\mathbf{u}^{(k)} - \mathbf{m}^{(k)}), \\ \mathbf{m}^{(k+\frac{1}{2})} &:= \mathbf{m}^{(k)} + r(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k)}), \\ \mathbf{u}^{(k+1)} &:= \operatorname{prox}_{g_\gamma}(\mathbf{b}^{(k+1)} + \mathbf{m}^{(k+\frac{1}{2})}), \\ \mathbf{m}^{(k+1)} &:= \mathbf{m}^{(k+\frac{1}{2})} + r(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)}). \end{aligned} \quad (27)$$

5.3 Convergence Analysis

Now, we delve into the convergence properties of the iterative schemes of ADMM (23) and SCPISM (24). This analysis will elucidate the conditions under which these methods converge and the nature of the solutions they yield. Specifically, the convergence of the proposed method is established through the following theorems.

Theorem 1 *Let the sequences $\{\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}\}_{k=0}^{\infty}$ be generated by the ADMM scheme (23a - 23c) and its Lagrangian is given by (20). Then the following three conditions hold:*

(a) **Sufficient decrease condition:** *For each iteration step k , $\exists \delta_1 > 0$ such that*

$$L_{\gamma}(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k+1)}) - L_{\gamma}(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) \leq -\delta_1 \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2.$$

(b) **Boundedness condition:** *The sequences $\{\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}\}_{k=0}^{\infty}$ are bounded and its Lagrangian $L_{\gamma}(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)})$ is lower bounded.*

(c) **Convergence:** *The Lagrangian in (20) is a Kurdyka-Łojasiewicz (KL) function, then the corresponding sequence $\{\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}\}$ converges to a unique stationary point $\{\mathbf{b}^{(*)}, \mathbf{u}^{(*)}, \mathbf{m}^{(*)}\}$.*

Note that the function $f(\mathbf{b}) = \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2$ is a continuously differentiable function with respect to \mathbf{b} . Its gradient is computed as $\nabla f(\mathbf{b}) = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}\mathbf{b}$. Then the Lipschitz constant for the gradient of the function $f(\mathbf{b})$ can be computed as

$$\begin{aligned} \|\nabla f(\mathbf{b}^{(k)}) - \nabla f(\mathbf{b}^{(k+1)})\| &\leq 2\|\mathbf{X}^T \mathbf{X}\| \|\mathbf{b}^{(k)} - \mathbf{b}^{(k+1)}\| \\ &\leq 2\lambda_{\max}(\mathbf{X}^T \mathbf{X}) \|\mathbf{b}^{(k)} - \mathbf{b}^{(k+1)}\| \\ &= l_f \|\mathbf{b}^{(k)} - \mathbf{b}^{(k+1)}\|, \end{aligned} \quad (28)$$

where $\|\mathbf{X}^T \mathbf{X}\|$ is the largest eigenvalue of $\mathbf{X}^T \mathbf{X}$ computed as $\lambda_{\max}(\mathbf{X}^T \mathbf{X})$. We denote the Lipschitz gradient constant as $l_f = 2\lambda_{\max}(\mathbf{X}^T \mathbf{X})$. The partial derivative of the Lagrangian (20) with respect to \mathbf{b} is given by

$$\partial_{\mathbf{b}} L_{\gamma}(\mathbf{b}, \mathbf{u}, \mathbf{m}) = \nabla f(\mathbf{b}) + \gamma(\mathbf{b} - \mathbf{u} + \mathbf{m}), \quad (29)$$

and the second partial derivative is $\frac{\partial^2 L_{\gamma}(\mathbf{b}, \mathbf{u}, \mathbf{m})}{\partial \mathbf{b}^2} = 2\mathbf{X}^T \mathbf{X} + \gamma \mathbf{I}$ which is positive definite. This implies that the Lagrangian is strongly convex with respect to \mathbf{b} . We will frequently use the properties of the Lipschitz gradient constant of $f(\mathbf{b})$ and the strong convexity of $L_{\gamma}(\mathbf{b}, \mathbf{u}, \mathbf{m})$ with respect to \mathbf{b} in the proof (see Appendix B.2).

Theorem 2 *Let the sequences $\{\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}\}_{k=0}^{\infty}$ be generated by the scheme (24a - 24d) and its Lagrangian is given by (20). Then conditions (a)-(c) in Theorem 1 hold.*

The proof of Theorem 2 is found in Appendix B.3.

In conclusion, the sufficient decreasing and boundedness conditions are satisfied when the learning rate $\gamma > \max\{\frac{2l_f^2}{\rho}, l_f\}$ in both Theorems 1 and 2. In practice, choosing γ involves a trade-off that requires careful consideration (see Section 5.5). The sufficient decreasing condition can be verified for the mean squared error (MSE) loss because the update step inherently minimizes the loss, ensuring it decreases as the number of iterations increases.

5.4 Implementation

The key idea in the WL1L0 optimization problem is to split the variable $\mathbf{b} = \mathbf{c} + \mathbf{d}$ into two components: \mathbf{c} , subject to L^1 -norm regularization, and \mathbf{d} , subject to L^0 -norm regularization. This decouples the optimization of \mathbf{c} and \mathbf{d} within the loss function $f(\mathbf{c} + \mathbf{d}) = \|\mathbf{y} - \mathbf{X}(\mathbf{c} + \mathbf{d})\|_2^2$ allowing separate updates for each part.

The decoupling is achieved through defining two translation functions, which manage the updates efficiently. The translation functions are defined as $\mathcal{T}_v(\mathbf{u}) = f(\mathbf{u} + \mathbf{v}) - \mathbf{v}$ and $\mathcal{T}_u(\mathbf{v}) = f(\mathbf{v} + \mathbf{u}) - \mathbf{u}$. These enable alternating updates between \mathbf{c} and \mathbf{d} , using the current estimates of the other variable. By leveraging these translations, the loss function $f(\mathbf{c} + \mathbf{d})$ is effectively split, allowing the proximal operators to handle both L^1 - and L^0 -norm regularizations. Hence, for WL1L0-ADMM, the updates are made in six steps, alternating between the two primal variables \mathbf{u} and \mathbf{v} , with corresponding dual variables \mathbf{m} and \mathbf{w} . The steps are:

$$\begin{aligned}
\mathbf{c}^{(k+1)} &:= \text{prox}_{\mathcal{T}_v(\mathbf{u})\gamma}(\mathbf{u}^{(k)} - \mathbf{m}^{(k)}), \\
\mathbf{u}^{(k+1)} &:= \text{prox}_{g\gamma}(\mathbf{c}^{(k+1)} + \mathbf{m}^{(k)}), \\
\mathbf{m}^{(k+1)} &:= \mathbf{m}^{(k)} + \mathbf{c}^{(k+1)} - \mathbf{u}^{(k+1)}, \\
\mathbf{d}^{(k+1)} &:= \text{prox}_{\mathcal{T}_u(\mathbf{v})\delta}(\mathbf{v}^{(k)} - \mathbf{w}^{(k)}), \\
\mathbf{v}^{(k+1)} &:= \text{prox}_{h\delta}(\mathbf{d}^{(k+1)} + \mathbf{w}^{(k)}), \\
\mathbf{w}^{(k+1)} &:= \mathbf{w}^{(k)} + \mathbf{d}^{(k+1)} - \mathbf{v}^{(k+1)}.
\end{aligned} \tag{30}$$

For WL1L0-SCPRSM, the updates are made in eight steps as

$$\begin{aligned}
\mathbf{c}^{(k+1)} &:= \text{prox}_{\mathcal{T}_v(\mathbf{u})\gamma}(\mathbf{u}^{(k)} - \mathbf{m}^{(k)}), \\
\mathbf{m}^{(k+\frac{1}{2})} &:= \mathbf{m}^{(k)} + r(\mathbf{c}^{(k+1)} - \mathbf{u}^{(k)}), \\
\mathbf{u}^{(k+1)} &:= \text{prox}_{g\gamma}(\mathbf{c}^{(k+1)} + \mathbf{m}^{(k+\frac{1}{2})}), \\
\mathbf{m}^{(k+1)} &:= \mathbf{m}^{(k+\frac{1}{2})} + r(\mathbf{c}^{(k+1)} - \mathbf{u}^{(k+1)}), \\
\mathbf{d}^{(k+1)} &:= \text{prox}_{\mathcal{T}_u(\mathbf{v})\delta}(\mathbf{v}^{(k)} - \mathbf{w}^{(k)}), \\
\mathbf{w}^{(k+\frac{1}{2})} &:= \mathbf{w}^{(k)} + r(\mathbf{d}^{(k+1)} - \mathbf{v}^{(k)}), \\
\mathbf{v}^{(k+1)} &:= \text{prox}_{h\delta}(\mathbf{d}^{(k+1)} + \mathbf{w}^{(k+\frac{1}{2})}), \\
\mathbf{w}^{(k+1)} &:= \mathbf{w}^{(k+\frac{1}{2})} + r(\mathbf{d}^{(k+1)} - \mathbf{v}^{(k+1)}).
\end{aligned} \tag{31}$$

Here, $\text{prox}_{g\gamma}(\mathbf{c} + \mathbf{m})$ is the proximal operator for L^1 , which is the soft-thresholding function with learning rate γ defined as

$$\text{prox}_{g\gamma}(\mathbf{c} + \mathbf{m}) = \mathcal{S}_\gamma(\mathbf{c} + \mathbf{m}) = \max(0, |\mathbf{c} + \mathbf{m}| - \gamma) \text{sgn}(\mathbf{c} + \mathbf{m}), \tag{32}$$

and $\text{prox}_{h\delta}(\mathbf{d} + \mathbf{w}) = \mathcal{H}_{\sqrt{2\delta}}(\mathbf{d} + \mathbf{w})$ is the proximal operator for L^0 , which is hard thresholding operator defined as

$$\mathcal{H}_{\sqrt{2\delta}}(\mathbf{d} + \mathbf{w}) = \begin{cases} 0, & \text{if } |\mathbf{d} + \mathbf{w}| < \sqrt{2\delta}, \\ \mathbf{d} + \mathbf{w}, & \text{if } |\mathbf{d} + \mathbf{w}| > \sqrt{2\delta}, \\ \{0, \mathbf{d} + \mathbf{w}\}, & \text{if } |\mathbf{d} + \mathbf{w}| = \sqrt{2\delta}. \end{cases} \tag{33}$$

The iterations are terminated when convergence is reached according to $\|(\mathbf{c}^{(k)} + \mathbf{d}^{(k)}) - (\mathbf{u}^{(k)} + \mathbf{v}^{(k)})\|_\infty \leq \beta(1 + \|\mathbf{m}^{(k)} + \mathbf{w}^{(k)}\|_\infty)$ for tolerance parameter β which was set to 10^{-5} .

For comparison purposes, we also implement the LASSO, SCAD, MCP and EN methods using the proximal ADMM and SCPRSM schemes (see Appendix D). Note that the convergence analysis of (30) and (31) is straightforward from (23) and (24). However, the introduction of translation functions and the variables \mathbf{c} and \mathbf{d} increases the dimensionality of the optimization problem, making the theoretical analysis very extensive. Therefore, we omit the detailed proof of the algorithm that involves the translation functions.

5.5 Determining the Learning Rate

Choosing the learning rate (step size) is crucial for efficiency and proper convergence of optimization algorithms. There are two main methods for determining the learning rates γ and δ (Beck, 2017; Bertsekas, 2016; Boyd & Vandenberghe, 2004): 1. Backtracking line-search: This method adjusts the learning rate iteratively based on specific criteria. However, it is both computationally expensive and time-consuming, as it requires

multiple evaluations of the objective function and its gradient during the search process. These repeated evaluations increase the overall computational load, particularly in high-dimensional problems. 2. Constant learning rate: In contrast, this method uses a fixed learning rate throughout the whole optimization process. It is computationally simpler and avoids the time overhead associated with frequent adjustments, making it more efficient in many scenarios.

We adopt the constant learning rate approach, using $\gamma^k = \frac{1}{\|\mathbf{X}\|_2}$ for all k , where $\|\mathbf{X}\|_2$ is the operator norm on the training set defined as

$$\|\mathbf{X}\|_2 := \max_{\|\mathbf{b}\|_2=1} \|\mathbf{X}\mathbf{b}\|_2. \quad (34)$$

Equivalently, $\|\mathbf{X}\|_2$ is the maximum singular value of \mathbf{X} ($\sigma_{\max}(\mathbf{X})$), which measures the maximum amount by which the matrix \mathbf{X} can stretch a vector \mathbf{b} relative to its original length (Horn & Johnson, 2013). The supremum-based definition can be applied in (34), particularly in infinite-dimensional contexts (Bhatia, 1997). However, in finite dimensions, the maximum and supremum coincide for the operator norm defined in (34). By setting the learning rate to $\frac{1}{\|\mathbf{X}\|_2}$, we ensure that the learning rate is scaled appropriately relative to the maximum possible stretch of \mathbf{X} . We use the same formula for δ .

5.6 Bayesian Optimization for Hyperparameter Tuning

Tuning the regularization parameter λ , the weight parameter α and the relaxation factor r via cross-validation or grid search can be computationally expensive. Bayesian Optimization (BO) is a more advanced, data-driven approach which offers a probabilistic model-based method for hyperparameter tuning (Gao et al., 2021; Shahriari et al., 2015). For the latest advancements, see Wang et al. (2023) and Yang et al. (2024).

BO uses a surrogate model, often a Gaussian Processes (GP), to approximate the true objective function. Hyperparameters are collected in $\vartheta = [\alpha, \lambda, r]$ and the objective function $\iota[\vartheta]$ is modeled as $\iota[\vartheta] \sim \mathcal{GP}(m[\vartheta], k[\vartheta, \vartheta'])$, where $m[\vartheta]$ is its mean and $k[\vartheta, \vartheta']$ the kernel (variance) function. The objective function is evaluated at j sequential points $\text{MSE}^{(j)} = \iota(\vartheta^{(j)})$, with $\text{MSE}^{(j)} \sim N(\iota(\vartheta^{(j)}), \sigma^2)$. This process induces a posterior over the acquisition function, guiding the selection of the next hyperparameters. Common acquisition functions include probability of improvement (PI), expected improvement (EI), upper confidence bound (UCB), and mutual information (MI) (Snoek et al., 2012). BO starts with an initial set of hyperparameters and objective function values to train the surrogate model. The acquisition function balances the posterior mean ($\varpi(\vartheta)$) for exploitation and variance ($v(\vartheta)$) for exploration. The GP-UCB is given by

$$\vartheta^{(j+1)} = \underset{\vartheta}{\operatorname{argmax}} \{ \varpi(\vartheta) + \varkappa v(\vartheta) \},$$

where $\varpi(\vartheta)$ is driven by the mean function $m(\vartheta)$, $v(\vartheta)$ by the variance function $k(\vartheta)$, and \varkappa determines the trade-off between exploitation and exploration. Contal et al. (2014) improved GP-UCB with the Gaussian Process Mutual Information algorithm (GP-MI) as $\vartheta^{(j+1)} = \underset{\vartheta}{\operatorname{argmax}} \{ \mu(\vartheta^{(j)}) + \sqrt{\log(2/\varrho)} (\sqrt{\Sigma(\vartheta^{(j)}) + \varsigma^{(j-1)}} - \sqrt{\varsigma^{(j-1)}}) \}$, where ς controls exploration, $0 < \varrho < 1$, and $\Sigma(\vartheta^{(j)})$ is the variance function at $\vartheta^{(j)}$.

6 Numerical Experiments

6.1 Materials

We evaluate our proposed method using one simulated genomic dataset as well as two real-world genomic datasets. Specifically, single-nucleotide polymorphisms (SNPs), which is a type of genetic variation that represent differences in one of the two nucleotides that make up an individual’s DNA at a specific location compared to the most common nucleotide pair found in a population. SNPs are typically represented by a count of 0, 1, or 2, where 0 means both nucleotides at the SNP location match the most common pair, 1 indicates that one of the two nucleotides differs from the common pair and 2 means both nucleotides at the location differ from the most common pair. This count system helps quantify genetic variation at a specific SNP site. A detailed explanation of these datasets is provided in Appendix C.

Simulated QTLMAS 2010 Dataset (Szydlowski & Paczyńska, 2011): This dataset comprises 3226 individuals, with genomic single nucleotide polymorphism (SNP) data organized in a matrix \mathbf{X} of size 3226×9723 , with two observed traits (response variables): a quantitative and a binary trait. In this study, the quantitative trait was chosen as the phenotype which is represented as a continuous response vector (\mathbf{y}) of length 3226.

Real Pig Dataset (Cleveland et al., 2012): This dataset contains genomic SNP data from 3534 individuals, organized in a matrix of size 3534×52842 , along with phenotypic data for five traits. We used trait 4, which had a heritability of 0.58, as the phenotype that constitutes a vector (\mathbf{y}) of length 3534.

Real Mice Dataset (Pérez & de Los Campos, 2014): This dataset contains data from 1814 individuals, with a genomic SNP data matrix of size 1814×10346 , along with two traits: body length (BL) and body mass index (BMI). In this study, the continuous trait BL was chosen to be our response vector (\mathbf{y}) that has a length of 1814.

6.2 Results

The WL1L0-ADMM, WL1L0-SCPRSM, EN-ADMM, EN-SCPRSM, LASSO-ADMM and LASSO-SCPRSM methods were implemented in Julia 1.10.1 (Bezanson et al., 2017) using the ProximalOperators package (Antonello et al., 2018). For SCAD-ADMM, SCAD-SCPRSM, MCP-ADMM and MCP-SCPRSM, we wrote our own code manually in Julia. For all methods, the BO was performed with the BayesianOptimization package using an ElasticGPE model and the squared exponential automatic relevance determination (SEArD) kernel (Fairbrother et al., 2018). The initial values of $\hat{\mathbf{b}}$, $\hat{\mathbf{c}}$ and $\hat{\mathbf{d}}$ were set to the marginal covariances between \mathbf{y} and \mathbf{X} , multiplied by 0.0001. By conducting preliminary runs for each set of hyperparameters using BO, we identified the optimal range of parameters. BO with the MI acquisition function was executed for hyperparameter tuning of all methods. The regression coefficients of the model are obtained from the training dataset, and once the model is trained, it predicts outcomes on the test dataset. The MSE is then calculated on the test dataset to assess the model’s generalization performance. The test MSE was monitored during the BO process to ensure convergence, which was indicated by no further decrease in MSE. All analyses were executed on a Linux computing platform equipped with an AMD EPYC 7302P 16-Core Processor and 32GB of system memory.

6.2.1 Simulated QTLMAS 2010 Dataset

BO was executed for 250 iterations with 4 GP function evaluations per iteration across all methods. The lower and upper bounds for λ_1 were set to 0.001 and 1000.0, 0.1 and 1800.0, and 0.001 and 600.0 for LASSO-ADMM, SCAD-ADMM, and MCP-ADMM, respectively. The lower and upper bounds for r were set to 0.01 and 0.999, 0.001 and 1.0, 0.001 and 1.0, and for λ_1 were set to 0.001 and 1000.0, 0.1 and 1800.0, and 0.001 and 600.0 for LASSO-SCPRSM, SCAD-SCPRSM, and MCP-SCPRSM, respectively. For EN-ADMM the lower and upper bounds for λ_1 were set to 10.0 and 600.0 and for λ_2 , they were set to 0.001 and 1.0, respectively. For EN-SCPRSM the lower and upper bounds for λ_1 were set to 10.0 and 500.0, for λ_2 , they were set to 0.001 and 200.0, for r , they were set to 0.01 and 0.99 respectively. For WL1L0-ADMM the lower and upper bounds for α were set to 0.01 and 0.99, and for λ_1 , they were set to 0.001 and 500.0, respectively. For WL1L0-SCPRSM, the lower and upper bounds for α were set to 0.0001 and 0.999, for r they were set to 0.0001 and 1.0, and for λ_1 they were set to 0.0001 and 500.0, respectively. The best result, with a minimum test MSE of 64.55, was found with WL1L0-SCPRSM at $\lambda_1 = 391.55$, $\alpha = 0.90$, and $r = 0.48$ (Table 1). Timing of the last evaluation with optimized parameters showed that SCAD-ADMM executed most quickly in only 6.68 seconds. It should be noted that those methods with one regularization parameter tend to be faster to train compared to other methods with two or three hyperparameters.

6.2.2 Real Pig Dataset

For the Pig dataset, we employed 5-fold cross-validation with random allocations into training and test data to obtain the minimum test MSE on the test data set, with the results averaged over the folds. Here, for all methods, BO was executed for 100 iterations with 3 GP function evaluations per iteration due to the large dataset size. The lower and upper bounds for λ_1 were set to 50.0 and 600.0, 0.1 and 1000.0, and 0.1 and 20.0 for LASSO-ADMM, SCAD-ADMM, and MCP-ADMM, respectively. The lower and upper bounds for

Method	min MSE	λ_1	λ_2	α	r	Time (s)	Number of non-zeros
LASSO-ADMM	66.55	294.22	-	-	-	9.31	417
LASSO-SCPRSM	65.95	312.51	-	-	0.19	21.51	334
SCAD-ADMM	66.50	309.65	-	-	-	6.68	386
SCAD-SCPRSM	65.91	299.79	-	-	0.42	19.68	352
MCP-ADMM	69.89	362.57	-	-	-	10.10	3898
MCP-SCPRSM	68.12	268.17	-	-	0.99	17.93	3883
EN-ADMM	66.52	307.65	0.79	-	-	11.87	390
EN-SCPRSM	65.92	288.15	0.001	-	0.99	22.18	375
WL1L0-ADMM	64.77	370.46	-	0.86	-	28.63	324
WL1L0-SCPRSM	64.55	391.55	-	0.90	0.48	25.39	275

Table 1: Performance evaluation of various regularization methods with optimal parameters on simulated QTLMAS data. The best-performing test MSE and most sparse model are highlighted in bold.

λ_1 were set to 50.0 and 400.0, 50.0 and 400.0, and 0.01 and 20.0, and lower and upper bounds for r were set to 0.01 and 1.0, 0.01 and 1.0, and 0.01 and 1.0 for LASSO-SCPRSM, SCAD-SCPRSM, and MCP-SCPRSM, respectively.

For EN-ADMM the lower and upper bounds for λ_1 were set to 10.0 and 600.0 and for λ_2 , they were set to 0.001 and 1.0 respectively. For EN-SCPRSM the lower and upper bounds for λ_1 were set to 0.1 and 200.0, for λ_2 , they were set to 0.01 and 100.0, for r , they were set to 0.001 and 1.0, respectively. For WL1L0-ADMM the lower and upper bounds for α were set to 0.001 and 0.99, and for λ_1 , they were set to 0.001 and 100.0, respectively. For WL1L0-SCPRSM, the lower and upper bounds for α were set to 0.001 and 0.99, for r they were set to 0.001 and 1.0, and for λ_1 they were set to 0.0001 and 200.0, respectively. We observed little variability in the minimum test MSE across the CV-folds for all methods. Hence, we report the mean minimum test MSE using the average estimates of the respective parameters for all methods. The best result, with a mean minimum test MSE of 4.48, was found with WL1L0-SCPRSM with mean estimates $\lambda_1 = 240.63$, $\alpha = 0.54$, and $r = 0.41$ (Table 2). The average timing over the folds of the last evaluation with optimized regularization parameters showed that SCAD-ADMM was fastest, taking only 25.8 seconds.

Method	min MSE	λ_1	λ_2	α	r	Time (s)	Number of non-zeros
LASSO-ADMM	4.53	118.75	-	-	-	26.2	1515
LASSO-SCPRSM	4.50	115.63	-	-	0.32	48.21	1200
SCAD-ADMM	4.64	127.80	-	-	-	25.8	1272
SCAD-SCPRSM	4.50	115.63	-	-	0.32	48.68	1200
MCP-ADMM	6.13	100.0	-	-	-	27.87	22909
MCP-SCPRSM	6.12	106.25	-	-	0.94	25.58	21627
EN-ADMM	4.53	120.63	0.31	-	-	26.86	1447
EN-SCPRSM	4.51	118.79	96.88	-	0.97	32.42	1433
WL1L0-ADMM	4.49	43.75	-	0.56	-	130.82	1093
WL1L0-SCPRSM	4.48	240.63	-	0.54	0.41	124.89	852

Table 2: Evaluation of the performance of various regularization methods with optimal parameters, averaged across five CV-folds on the pig data. The best-performing test MSE and most sparse model are highlighted in bold.

6.2.3 Real Mice Dataset

Similar to the Pig dataset, we employed 5-fold cross-validation also for this data. BO was executed for 100 iterations with 4 GP function evaluations per iteration across all methods. The lower and upper bounds for λ_1 were set to 0.0001 and 20.0, 0.001 and 20.0, and 0.1 and 20.0 for LASSO-ADMM, SCAD-ADMM, and MCP-ADMM, respectively. The lower and upper bounds for λ_1 were set to 0.001 and 35.0, 0.001 and 35.0, and 0.01 and 20.0, and lower and upper bounds for r were set to 0.001 and 1.0, 0.001 and 1.0, and 0.01 and

1.0 for LASSO-SCPRSM, SCAD-SCPRSM, and MCP-SCPRSM, respectively. For EN-ADMM the lower and upper bounds for λ_1 were set to 0.01 and 18.0 and for λ_2 , they were set to 0.001 and 1.0 respectively. For EN-SCPRSM the lower and upper bounds for λ_1 were set to 0.1 and 42.0, for λ_2 , they were set to 0.01 and 40.0, for r , they were set to 0.01 and 1.0, respectively.

For WL1L0-ADMM, the lower and upper bounds for α were set to 0.001 and 0.99, and for λ_1 , they were set to 0.001 and 100.0, respectively. For WL1L0-SCPRSM, the lower and upper bounds for α were set to 0.001 and 0.99, for r they were set to 0.001 and 1.0, and for λ_1 they were set to 0.0001 and 200.0, respectively. The best result, with a mean minimum test MSE of 0.259 was found with WL1L0-SCPRSM at the average estimates $\lambda_1 = 56.25$, $\alpha = 0.28$, and $r = 0.16$ (Table 3).

The average timing over the folds of the last evaluation with optimized regularization parameters showed that SCAD-ADMM once again was fastest with a time of only 2.08 seconds.

Method	min MSE	λ_1	λ_2	α	r	Time (s)	Number of non-zeros
LASSO-ADMM	0.273	20.0	-	-	-	2.10	319
LASSO-SCPRSM	0.267	24.06	-	-	0.81	3.14	280
SCAD-ADMM	0.274	24.0	-	-	-	2.08	319
SCAD-SCPRSM	0.267	24.06	-	-	0.81	2.78	280
MCP-ADMM	0.273	20.0	-	-	-	2.14	432
MCP-SCPRSM	0.273	20.0	-	-	0.01	26.62	387
EN-ADMM	0.276	18.0	0.99	-	-	2.15	539
EN-SCPRSM	0.267	24.98	38.75	-	0.97	2.51	256
WL1L0-ADMM	0.265	43.75	-	0.56	-	19.97	216
WL1L0-SCPRSM	0.259	56.25	-	0.28	0.16	22.83	188

Table 3: Evaluation of the performance of various regularization methods with optimal parameters, averaged across five CV-folds on the mice dataset. The best-performing test MSE and most sparse model are highlighted in bold.

7 Discussion

The WL1L0 method demonstrates superior performance across all datasets by achieving the lowest MSE and the fewest non-zero coefficients. This highlights its effectiveness and efficiency as a regularization technique in high-dimensional data analysis, making it a valuable alternative to the LASSO, SCAD, MCP and EN. The weighting parameter α in WL1L0 provides flexibility in tuning the regularization effect, making the method adaptable to different datasets and problem settings. This adaptability enhances its robustness and applicability across diverse scenarios.

It has been demonstrated several times that the SCAD and MCP often outperform the LASSO (Fan et al., 2014a; Fan & Li, 2001; Zhang, 2010). However, while the LASSO, SCAD, MCP and EN also offer competitive approaches to regularization, the WL1L0 method consistently outperforms them, providing enhanced model sparsity and interpretability without compromising predictive accuracy. The joint sparsity induced by the L^1 and the L^0 components make the resulting model more interpretable. This is crucial in many scientific and industrial applications, where understanding the model is as important as its predictive power.

The use of the SCPRSM algorithm introduces an additional parameter r , which allows for finer control over the optimization process and potentially leads to better convergence properties and more precise model fitting. Across all datasets used, the SCPRSM variants demonstrate strong performance by achieving the smallest minimum MSEs while maintaining a manageable number of non-zero coefficients. Specifically, WL1L0-SCPRSM consistently achieves the lowest MSE across all our datasets, demonstrating its superior ability to minimize prediction errors. It is likely that it will be highly effective in terms of both accuracy and reliability across other types of data. Several other studies have shown that SCPRSM outperforms ADMM (Li & Yuan, 2015; Li et al., 2021).

In this paper, we have mostly focused on the regularization part and note that there certainly is room for computational advancements. ADMM can be improved using techniques such as accelerated ADMM (Zhang et al., 2019; Zeng et al., 2024, and references therein) as well as stochastic distributed ADMM (Chen et al., 2021, and references therein). Similarly, for SCPRSM, further computational improvements can be achieved via stochastic SCPRSM (Na et al., 2017) and indefinite-proximal SCPRSM (Gu et al., 2022; Bai et al., 2023).

8 Conclusion

This paper introduces a novel joint weighted L^1 - and L^0 -norm method denoted WL1L0 based on proximal mappings and translation functions, aiming to debias the bias introduced by the L^1 -norm when applied to high-dimensional data. Our model introduces a weighting parameter α , allowing for the adjustment of the influence of both regularizers. The convergence of ADMM and SCPRSM for the developed method is shown under reasonable assumptions. All hyper-parameters are optimized using Bayesian optimization. The WL1L0-SCPRSM method consistently achieves the lowest MSE across all datasets when compared to all other tested regularization methods (LASSO, EN, SCAD and MCP). Hence, the WL1L0-SCPRSM's superior performance across different genomic high-dimensional datasets demonstrates its versatility. Our current paper focuses primarily on prediction. In future work, we plan to specifically address the properties of variable selection.

Acknowledgements

We acknowledge funding from the University of Oulu & the Academy of Finland Profi 326291.

References

- Niccolò Antonello, Lorenzo Stella, Panagiotis Patrinos, and Toon Van Waterschoot. Proximal gradient algorithms: Applications in signal processing. *arXiv preprint arXiv:1803.01621*, 2018.
- Hédy Attouch, Jérôme Bolte, Patrick Redont, and Antoine Soubeyran. Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Łojasiewicz inequality. *Mathematics of Operations Research*, 35(2):438–457, 2010.
- Hedy Attouch, Jérôme Bolte, and Benar Fux Svaiter. Convergence of descent methods for semi-algebraic and tame problems: Proximal algorithms, forward–backward splitting, and regularized Gauss–Seidel methods. *Mathematical Programming, Series A*, 137(1):91–129, 2013.
- Jian-Chao Bai, Feng-Miao Bian, Xiao-Kai Chang, and Lin Du. Accelerated stochastic peaceman–rachford method for empirical risk minimization. *Journal of the Operations Research Society of China*, 11(4): 783–807, 2023.
- Amir Beck. *First-order Methods in Optimization*. SIAM, 2017.
- Dimitri P. Bertsekas. *Nonlinear Programming, 3rd ed.* Athena Scientific, Nashua, NH, 2016.
- Dimitris Bertsimas, Jean Pauphilet, and Bart Van Parys. Sparse regression: Scalable algorithms and empirical performance. *Statistical Science*, 35(4):555–578, 2020.
- Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral B Shah. Julia: A fresh approach to numerical computing. *SIAM Review*, 59(1):65–98, 2017.
- Rajendra Bhatia. *Matrix Analysis*, volume 169. Springer-Verlag New York, 1997.
- Jérôme Bolte, Shoham Sabach, and Marc Teboulle. Proximal alternating linearized minimization for non-convex and nonsmooth problems. *Mathematical Programming*, 146(1):459–494, 2014.
- Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122, 2011.

- Stephen P Boyd and Lieven Vandenbergh. *Convex Optimization*. Cambridge University Press, 2004.
- Peter Bühlmann and Sara Van De Geer. *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer-Verlag Berlin Heidelberg, Berlin, Germany, 2011.
- Hao Chen, Yu Ye, Ming Xiao, Mikael Skoglund, and H Vincent Poor. Coded stochastic admm for decentralized consensus optimization with edge computing. *IEEE Internet of Things Journal*, 8(7):5360–5373, 2021.
- Victor Chernozhukov, Christian Hansen, and Yuan Liao. A lava attack on the recovery of sums of dense and sparse signals. *The Annals of Statistics*, 45(1):39–76, 2017.
- Francis H Clarke, Yuri S Ledyae, Ronald J Stern, and Peter R Wolenski. *Nonsmooth Analysis and Control Theory*, volume 178. Springer Science & Business Media, 2008.
- Matthew A Cleveland, John M Hickey, and Selma Forni. A common dataset for genomic analysis of livestock populations. *G3: Genes/ Genomes/ Genetics*, 2(4):429–435, 2012.
- Emile Contal, Vianney Perchet, and Nicolas Vayatis. Gaussian process optimization with mutual information. In *International Conference on Machine Learning*, volume 32, pp. 253–261. PMLR, 2014.
- Etienne Corman and Xiaoming Yuan. A generalized proximal point algorithm and its convergence rate. *SIAM Journal on Optimization*, 24(4):1614–1638, 2014.
- Jamie Fairbrother, Christopher Nemeth, Maxime Rischard, Johanni Brea, and Thomas Pinder. Gaussian-processes. jl: A nonparametric Bayes package for the Julia language. *arXiv preprint arXiv:1812.09064*, 2018.
- Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2001.
- Jianqing Fan and Jinchi Lv. A selective overview of variable selection in high dimensional feature space. *Statistica Sinica*, 20(1):101–148, 2010.
- Jianqing Fan, Jinchi Lv, and Lei Qi. Sparse high dimensional models in economics. *Annual Review of Economics*, 3(1):291–317, 2011.
- Jianqing Fan, Yingying Fan, and Emre Barut. Adaptive robust variable selection. *Annals of Statistics*, 42(1):324–351, 2014a.
- Jianqing Fan, Fang Han, and Han Liu. Challenges of big data analysis. *National Science Review*, 1(2):293–314, 2014b.
- Masao Fukushima and Hisashi Mine. A generalized proximal point algorithm for certain non-convex minimization problems. *International Journal of Systems Science*, 12(8):989–1000, 1981.
- Haiping Gao, Shifa Zhong, Wenlong Zhang, Thomas Igou, Eli Berger, Elliot Reid, Yangying Zhao, Dylan Lambeth, Lan Gan, Moyosore A. Afolabi, et al. Revolutionizing membrane design using machine learning-Bayesian optimization. *Environmental Science & Technology*, 56(4):2572–2581, 2021.
- Christophe Giraud. *Introduction to High-Dimensional Statistics*. CRC Press, Boca Raton, FL, 2015.
- Yan Gu, Bo Jiang, and DR Han. An indefinite-proximal-based strictly contractive peaceman-rachford splitting method. *Journal of Computational Mathematics*, 41:1017–1040, 2022.
- Trevor Hastie, Robert Tibshirani, and Martin Wainwright. *Statistical Learning with Sparsity: The Lasso and Generalizations*. CRC Press, New York, 2015.
- Trevor Hastie, Robert Tibshirani, and Ryan Tibshirani. Best subset, forward stepwise or Lasso? Analysis and recommendations based on extensive comparisons. *Statistical Science*, 35(4):579–592, 2020.

- Bingsheng He, Li-Zhi Liao, Deren Han, and Hai Yang. A new inexact alternating directions method for monotone variational inequalities. *Mathematical Programming*, 92:103–118, 2002.
- Bingsheng He, Han Liu, Zhaoran Wang, and Xiaoming Yuan. A strictly contractive Peaceman–Rachford splitting method for convex programming. *SIAM Journal on Optimization*, 24(3):1011–1040, 2014.
- Georg Heinze, Christine Wallisch, and Daniela Dunkler. Variable selection—a review and recommendations for the practicing statistician. *Biometrical Journal*, 60(3):431–449, 2018.
- R R Hocking and R N Leslie. Selection of the best subset in regression analysis. *Technometrics*, 9:531–540, 1967.
- Arthur E Hoerl and Robert W Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- Roger A Horn and Charles R Johnson. *Matrix Analysis*. Cambridge University Press, 2013.
- Iain M Johnstone and D Michael Titterton. Statistical challenges of high-dimensional data. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 367(1906):4237–4253, 2009.
- Alexander Kaplan and Rainer Tichatschke. Proximal point methods and nonconvex optimization. *Journal of Global Optimization*, 13:389–406, 1998.
- Hai Yen Le. Generalized subdifferentials of the rank function. *Optimization Letters*, 7:731–743, 2013.
- Peixuan Li, Yuan Shen, Suhong Jiang, Zehua Liu, and Caihua Chen. Convergence study on strictly contractive peaceman–rachford splitting method for nonseparable convex minimization models with quadratic coupling terms. *Computational Optimization and Applications*, 78:87–124, 2021.
- Xinxin Li and Xiaoming Yuan. A proximal strictly contractive Peaceman-Rachford splitting method for convex programming with applications to imaging. *SIAM Journal on Imaging Sciences*, 8(2):1332–1365, 2015.
- Xingran Liao, Xuekai Wei, and Mingliang Zhou. Minimax concave penalty regression for superresolution image reconstruction. *IEEE Transactions on Consumer Electronics*, 70(1):2999–3007, 2023.
- Christos Louizos, Max Welling, and Diederik P Kingma. Learning sparse neural networks through l_0 regularization. *arXiv preprint arXiv:1712.01312*, 2017.
- Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- Boris S Mordukhovich. *Variational Analysis and Generalized Differentiation II: Applications*, volume 331. Springer, 2006.
- Sen Na, Mingyuan Ma, Shuming Ma, and Guangju Peng. Stochastic strictly contractive peaceman-rachford splitting method. *arXiv preprint arXiv:1711.04955*, 2017.
- Joseph O Ogotu and Hans-Peter Piepho. Regularized group regression methods for genomic prediction: Bridge, MCP, SCAD, group bridge, group lasso, sparse group lasso, group MCP and group SCAD. In *BMC proceedings*, volume 8, pp. 1–9. Springer, 2014.
- Neal Parikh and Stephen Boyd. Proximal algorithms. *Foundations and Trends® in Optimization*, 1(3):123–231, 2013.
- Donald W Peaceman and Henry H Rachford, Jr. The numerical solution of parabolic and elliptic differential equations. *Journal of the Society for Industrial and Applied Mathematics*, 3(1):28–41, 1955.
- Paulino Pérez and Gustavo de Los Campos. Genome-wide regression and prediction with the BGLR statistical package. *Genetics*, 198(2):483–495, 2014.

- Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. Taking the human out of the loop: A review of Bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2015.
- Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical Bayesian optimization of machine learning algorithms. *Advances in Neural Information Processing Systems*, 25, 2012.
- Maciej Szydlowski and Paulina Paczyńska. QTLMAS 2010: Simulated dataset. In *BMC proceedings*, volume 5, pp. 1–3. Springer, 2011.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- Laura Tološi and Thomas Lengauer. Classification with correlated features: unreliability of feature ranking and solutions. *Bioinformatics*, 27(14):1986–1994, 2011.
- Martin J Wainwright. *High-dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press, Cambridge, United Kingdom, 2019.
- Patrik Waldmann. A proximal LAVA method for genome-wide association and prediction of traits with mixed inheritance patterns. *BMC Bioinformatics*, 22(1):1–16, 2021.
- Hao Wang, Zhanglei Shi, Chi-Sing Leung, and Hing Cheung So. ADMM-MCP framework for sparse recovery with global convergence. *IEEE Transactions on Signal Processing*, 2018.
- Ting Wang and Hongwei Liu. A class of modified accelerated proximal gradient methods for nonsmooth and nonconvex minimization problems. *Numerical Algorithms*, 95(1):207–241, 2024.
- Xilu Wang, Yaochu Jin, Sebastian Schmitt, and Markus Olhofer. Recent advances in Bayesian optimization. *ACM Computing Surveys*, 55(13s):1–36, 2023.
- Yu Wang, Wotao Yin, and Jinshan Zeng. Global convergence of ADMM in nonconvex nonsmooth optimization. *Journal of Scientific Computing*, 78:29–63, 2019.
- Yuqia Wu, Shaohua Pan, and Shujun Bi. Kurdyka–łojasiewicz property of zero-norm composite functions. *Journal of Optimization Theory and Applications*, 188:94–112, 2021.
- Kaixin Yang, Long Liu, and Yalu Wen. The impact of Bayesian optimization on feature selection. *Scientific Reports*, 14(1):3948, 2024.
- Maryam Yashtini. Convergence and rate analysis of a proximal linearized admm for nonconvex nonsmooth optimization. *Journal of Global Optimization*, 84(4):913–939, 2022.
- Jihun Yun, Peng Zheng, Eunho Yang, Aurelie Lozano, and Aleksandr Aravkin. Trimming the ℓ_1 regularizer: Statistical analysis, optimization, and applications to deep learning. In *International Conference on Machine Learning*, pp. 7242–7251. PMLR, 2019.
- Jihun Yun, Aurélie C Lozano, and Eunho Yang. Adaptive proximal gradient methods for structured neural networks. *Advances in Neural Information Processing Systems*, 34:24365–24378, 2021.
- Yuxuan Zeng, Zhiguo Wang, Jianchao Bai, and Xiaojing Shen. An accelerated stochastic admm for nonconvex and nonsmooth finite-sum optimization. *Automatica*, 163:111554, 2024.
- Cun-Hui Zhang. Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, 38(2):894–942, 2010.
- Juyong Zhang, Yue Peng, Wenqing Ouyang, and Bailin Deng. Accelerating admm for efficient simulation and optimization. *ACM Transactions on Graphics (TOG)*, 38(6):1–21, 2019.
- Peng Zhao and Bin Yu. On model selection consistency of Lasso. *The Journal of Machine Learning Research*, 7:2541–2563, 2006.

Tuo Zhao, Han Liu, and Tong Zhang. Pathwise coordinate optimization for sparse learning: Algorithm and theory. *The Annals of Statistics*, 46(1):180–218, 2018.

Liu Ziyin and Zihao Wang. spread: Solving L1 penalty with SGD. In *International Conference on Machine Learning*, pp. 43407–43422. PMLR, 2023.

Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.

Appendix

A Related Problems

Consider the following formulation

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\operatorname{argmin}} F(\mathbf{b}) := \Gamma(\mathbf{b}) + \lambda(1 - \alpha)\|\mathbf{b}\|_0, \quad (35)$$

where $\Gamma(\mathbf{b}) := \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2 + \lambda\alpha\|\mathbf{b}\|_1$. Since $\Gamma(\mathbf{b})$ is a convex function, it has a global minimum value. By the Weierstrass theorem, a continuous function over a nonempty compact set attains a minimum. The existence of an optimal solution is guaranteed if a function is continuous over a closed set and coercive over the set (Bertsekas, 2016). Beck (2017) demonstrates that the latter extends to closed functions, i.e. a closed and coercive function over a closed set attains an optimal solution.

For the case $\|\mathbf{b}\|_0 = \sum_{i=1}^p \mathbf{1}(b_i \neq 0)$ with $\lambda(1 - \alpha) > 0$, we need to show it is a closed function. Let

$g(\mathbf{b}) = \sum_{i=1}^p I(b_i)$, where $I : \mathbb{R} \rightarrow \{0, 1\}$ is defined as

$$I(b_i) = \begin{cases} \lambda(1 - \alpha), & b_i \neq 0, \\ 0, & b_i = 0. \end{cases}$$

The function $I(\cdot)$ is closed since its level sets, given by

$$\operatorname{Lev}(I, \eta) = \begin{cases} \emptyset, & \eta < 0, \\ \{0\}, & \eta \in [0, 1), \\ \mathbb{R}, & \eta \geq 1, \end{cases} \quad (36)$$

are closed sets. Here, g is a closed function. Furthermore, using Theorem 2.6 in (Beck, 2017), the closedness of $\|\mathbf{b}\|_0$ implies its lower semi-continuity.

A vector \mathbf{b}^* is a local minimum of the function F , if there exists $\varepsilon > 0$ such that $F(\mathbf{b}^*) \leq F(\mathbf{b})$ for all $\mathbf{b} \in \mathbb{R}^p$ with $\|\mathbf{b} - \mathbf{b}^*\| < \varepsilon$. A vector \mathbf{b}^* is a global minimum if $F(\mathbf{b}^*) \leq F(\mathbf{b})$ for all $\mathbf{b} \in \mathbb{R}^p$.

For illustration purposes, we generated a random design matrix \mathbf{X} with dimension 100×500 by simulating 100 samples and 500 features. For each value in the range between -5 and 5, the outcomes of a function $F(b)$ that combines the squared error loss with the regularization term that consists of the weighted sum of the L^1 and L^0 -norms were produced. The regularization parameters were set to $\lambda = 2$ and $\alpha = 0.6$. Therefore, the plotted $F(b)$ includes both the error and regularization terms, rather than solely the penalty norms. One can see that $F(b)$ is nonconvex because any point between the endpoints A and B , as indicated by the dashed red line in Figure 1, lies outside the domain of $F(b)$. In fact, the shape of the function $F(b)$ is similar to that of nonconvex regularization methods such as SCAD and MCP (Fan & Li, 2001; Zhang, 2010; Zhao et al., 2018).

Various shapes of the function $F(b)$ for different values of α and $\lambda = 1$ are depicted in Figure 2 to illustrate the impact of α on the function’s behavior.

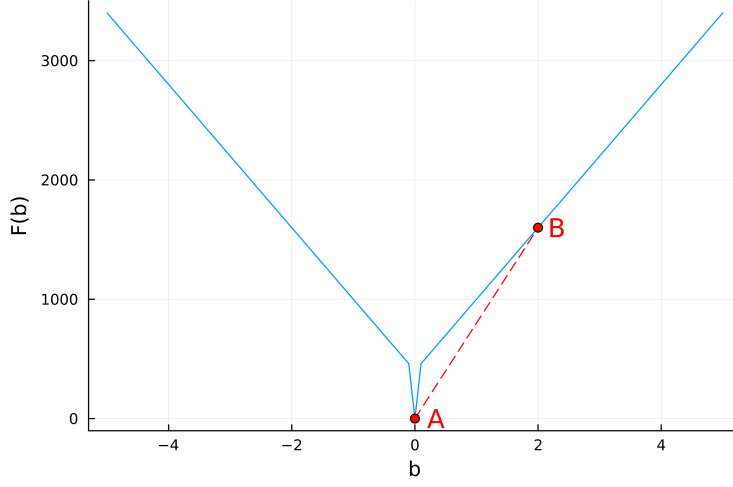


Figure 1: Illustration of the nonconvexity of the function F

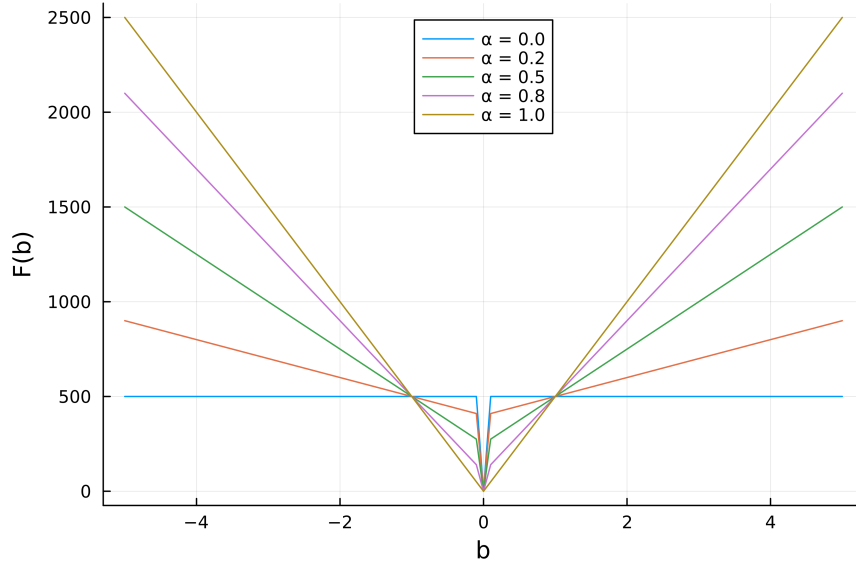


Figure 2: Different function values of F with regularizer $\lambda = 1$ and $\alpha = 0 (L^0), 0.2, 0.5, 0.8, 1 (L^1)$.

B Proofs

B.1 Proof of Lemma 1

Given $f(\mathbf{u}) = \langle \mathbf{m}, \mathbf{u} \rangle + a$, the proximal operator of the function $f(\mathbf{u})$ is defined as $\text{prox}_f(\mathbf{u}) = \underset{\mathbf{v}}{\text{argmin}} \{f(\mathbf{v}) + \frac{1}{2}\|\mathbf{v} - \mathbf{u}\|_2^2\}$. Substituting the affine function $f(\mathbf{v}) = \langle \mathbf{m}, \mathbf{v} \rangle + a$ gives us $\text{prox}_f(\mathbf{u}) = \underset{\mathbf{v}}{\text{argmin}} \{\langle \mathbf{m}, \mathbf{v} \rangle + a + \frac{1}{2}\|\mathbf{v} - \mathbf{u}\|_2^2\}$. The term a is constant and does not affect the minimization (it doesn't depend on \mathbf{v}), so we can ignore it. Next, the squared norm term $\frac{1}{2}\|\mathbf{v} - \mathbf{u}\|_2^2$ can be expanded as $\frac{1}{2}\|\mathbf{v} - \mathbf{u}\|_2^2 = \frac{1}{2}(\|\mathbf{v}\|_2^2 - 2\langle \mathbf{v}, \mathbf{u} \rangle + \|\mathbf{u}\|_2^2)$. Since $\|\mathbf{u}\|_2^2$ is constant with respect to \mathbf{v} , it can also be ignored in the minimization. Now, combining the linear terms involving \mathbf{v} , we have $\underset{\mathbf{v}}{\text{argmin}} \{\frac{1}{2}\|\mathbf{v}\|_2^2 + \langle \mathbf{m} - \mathbf{u}, \mathbf{v} \rangle\}$. This expression is a quadratic function in \mathbf{v} . To minimize it, we set the gradient of the objective function

with respect to \mathbf{v} equal to zero. The gradient of $\frac{1}{2}\|\mathbf{v}\|_2^2$ is \mathbf{v} . The gradient of $\langle \mathbf{m} - \mathbf{u}, \mathbf{v} \rangle$ is $\mathbf{m} - \mathbf{u}$. Setting the total gradient to zero gives us: $\mathbf{v} + (\mathbf{m} - \mathbf{u}) = 0$. Solving for \mathbf{v} , we find $\mathbf{v} = \mathbf{u} - \mathbf{m}$. Therefore, $\text{prox}_f(\mathbf{u}) = \mathbf{u} - \mathbf{m}$.

B.2 Proof of Theorem 1

While our model differs from that of Wang et al. (2019; 2018), we adopt a similar proof framework.

Proof (a): From (23a), $\mathbf{b}^{(k+1)}$ minimizes $L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)})$ and since $L_\gamma(\mathbf{b}, \mathbf{u}, \mathbf{m})$ is strongly convex with respect to \mathbf{b} , the Lagrangian function satisfies the following inequality (Beck, 2017):

$$L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) - L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) \leq -\frac{\rho}{2}\|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2, \quad (37)$$

where $L_\gamma(\cdot)$ is a ρ -strongly convex function ($\rho > 0$). From the augmented Lagrangian function in (20), we have

$$L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k+1)}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k)}) = \gamma \left(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)} \right)^\top \left(\mathbf{m}^{(k+1)} - \mathbf{m}^{(k)} \right). \quad (38)$$

Now we rewrite (23c) as

$$\mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)} = \mathbf{m}^{(k+1)} - \mathbf{m}^{(k)}. \quad (39)$$

From (23b) and (29), we get

$$\nabla f(\mathbf{b}^{(k+1)}) + \gamma(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)} + \mathbf{m}^{(k)}) = 0. \quad (40)$$

Substituting (23c) into (40), we obtain

$$\nabla f(\mathbf{b}^{(k+1)}) = -\gamma \mathbf{m}^{(k+1)}. \quad (41)$$

Using (39) and (41), (38) becomes

$$\gamma \left(\|\mathbf{m}^{(k+1)} - \mathbf{m}^{(k)}\|^2 \right) = \gamma \left(\left\| -\frac{1}{\gamma} \nabla f(\mathbf{b}^{(k+1)}) + \frac{1}{\gamma} \nabla f(\mathbf{b}^{(k)}) \right\|^2 \right) \leq \frac{l_f^2}{\gamma} \left(\|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|^2 \right). \quad (42)$$

Hence,

$$L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k+1)}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k)}) \leq \frac{l_f^2}{\gamma} \left(\|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|^2 \right). \quad (43)$$

Here, the term $l_f \geq 0$ denotes a Lipschitz gradient of the function $f(\mathbf{b})$. From (23b) we have that

$$L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k)}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) \leq 0. \quad (44)$$

Finally, combining (37), (43) and (44), we obtain the desired inequality as

$$\begin{aligned} & L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k+1)}) - L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) \\ &= L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k+1)}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k)}) \\ &\quad + L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k)}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) \\ &\quad + L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) - L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) \\ &\leq \left(\frac{l_f^2}{\gamma} - \frac{\rho}{2} \right) \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2 \\ &= -\delta_1 \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2, \end{aligned} \quad (45)$$

where $\sigma_1 = \frac{\rho}{2} - \frac{l_f^2}{\gamma}$ and $\gamma > \frac{2l_f^2}{\rho}$. Hence, the sufficient decreasing condition is met.

Proof (b): We utilize the descent lemma to prove that $L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)})$ is lower bounded for any k .

Lemma 2 (*Descent lemma*) *Let the function f belong to the class of continuously differentiable functions with constant l_f Lipschitz continuous gradients. Then for any two points $\mathbf{b}^{(k)}$ and $\mathbf{u}^{(k)}$,*

$$f(\mathbf{u}^{(k)}) \leq f(\mathbf{b}^{(k)}) + \nabla f(\mathbf{b}^{(k)})^\top (\mathbf{u}^{(k)} - \mathbf{b}^{(k)}) + \frac{l_f}{2} \|\mathbf{u}^{(k)} - \mathbf{b}^{(k)}\|_2^2. \quad (46)$$

The proof of the descent lemma can be found in (Beck, 2017), see Lemma 5.7.

As a result of the Descent lemma, the sequence is lower bounded as

$$\begin{aligned} L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) &= f(\mathbf{b}) + g(\mathbf{u}) + \frac{\gamma}{2} \|\mathbf{b} - \mathbf{u} + \frac{1}{\gamma} \mathbf{m}\|_2^2 - \frac{\gamma}{2} \|\frac{1}{\gamma} \mathbf{m}\|_2^2 \\ &= f(\mathbf{b}^{(k)}) + g(\mathbf{u}^{(k)}) + \mathbf{m}^T (\mathbf{b}^{(k)} - \mathbf{u}^{(k)}) + (\gamma/2) \|\mathbf{b}^{(k)} - \mathbf{u}^{(k)}\|_2^2 \\ &\geq f(\mathbf{u}^{(k)}) + g(\mathbf{u}^{(k)}) + \left(\frac{\gamma}{2} - \frac{l_f}{2} \right) \|\mathbf{u}^{(k)} - \mathbf{m}^{(k)}\|_2^2 \\ &\geq -\infty \text{ for } \gamma \geq l_f. \end{aligned} \quad (47)$$

Hence, from (47), $L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)})$ is lower bounded.

As established in the proof (a), the sufficient descent property implies that $L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)})$ is upper-bounded by $L_\gamma(\mathbf{b}^0, \mathbf{u}^0, \mathbf{m}^0)$. To prove that the sequence $\{\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}\}$ is bounded, we start by rewriting (45) as

$$\begin{aligned} \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2 &\leq \frac{1}{\delta_1} (L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k+1)})) \\ \sum_{k=0}^l \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2 &\leq \frac{1}{\delta_1} (L_\gamma(\mathbf{b}^0, \mathbf{u}^0, \mathbf{m}^0) - L_\gamma(\mathbf{b}^{l+1}, \mathbf{u}^{l+1}, \mathbf{m}^{l+1})) \\ &< \infty. \end{aligned} \quad (48)$$

Equation (48) also holds as $l \rightarrow \infty$. Hence, $\mathbf{b}^{(k)}$ is bounded.

From (42), we obtain

$$\begin{aligned} \|\mathbf{m}^{(k+1)} - \mathbf{m}^{(k)}\|_2^2 &\leq \frac{l_f^2}{\gamma^2} \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2 \\ \sum_{k=0}^l \|\mathbf{m}^{(k+1)} - \mathbf{m}^{(k)}\|_2^2 &< \infty. \end{aligned} \quad (49)$$

This implies that $\mathbf{m}^{(k)}$ is bounded.

Finally, from (39) we obtain $\mathbf{u}^{(k+1)} = \mathbf{b}^{(k+1)} - \mathbf{m}^{(k+1)} + \mathbf{m}^{(k)}$ and $\mathbf{u}^{(k)} = \mathbf{b}^{(k)} - \mathbf{m}^{(k)} + \mathbf{m}^{(k-1)}$. Then

$$\begin{aligned} \|\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}\|_2^2 &= \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)} + \mathbf{m}^{(k)} - \mathbf{m}^{(k+1)} + \mathbf{m}^{(k-1)} - \mathbf{m}^{(k)}\|_2^2 \\ &\leq \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2 + \|\mathbf{m}^{(k+1)} - \mathbf{m}^{(k)}\|_2^2 + \|\mathbf{m}^{(k)} - \mathbf{m}^{(k-1)}\|_2^2. \end{aligned}$$

Consequently, we obtain

$$\sum_{k=1}^{\infty} \|\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}\|_2^2 < \infty. \quad (50)$$

Hence, the sequence $\{\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}\}$ is bounded.

Proof (c): The augmented Lagrangian function $L_\gamma(\mathbf{b}, \mathbf{u}, \mathbf{m}) = f(\mathbf{b}) + g(\mathbf{u}) + \frac{\gamma}{2} \|\mathbf{b} - \mathbf{u} + \mathbf{m}\|_2^2 - \frac{\gamma}{2} \|\mathbf{m}\|_2^2$ defined as $L_\gamma : \mathbb{R}^n \rightarrow (-\infty, \infty]$ is proper and lower semi-continuous, where $f(\mathbf{b}) = \|\mathbf{y} - \mathbf{X}\mathbf{b}\|_2^2$, $g(\mathbf{u}) = \|\mathbf{u}\|_1 + \|\mathbf{u}\|_0$ and $h(\mathbf{b}, \mathbf{u}, \mathbf{m}) = \frac{\gamma}{2} \|\mathbf{b} - \mathbf{u} + \mathbf{m}\|_2^2 - \frac{\gamma}{2} \|\mathbf{m}\|_2^2$. If $L_\gamma(\mathbf{b}, \mathbf{u}, \mathbf{m})$ is semi-algebraic, then it satisfies the KL property

at any point of its domain. Note that both f and h are real polynomial functions, which are semi-algebraic functions (Attouch et al., 2013; Bolte et al., 2014). Both $\|\mathbf{u}\|_0$ and $\|\mathbf{u}\|_1$ have piecewise linear graphs and are therefore semi-algebraic (see Example 3 and 4 in (Bolte et al., 2014), respectively).

Furthermore, consider that $g_1(\mathbf{u}) = \lambda\alpha\|\mathbf{u}\|_1$ and $g_2 = \lambda(1-\alpha)\|\mathbf{u}\|_0$. Their proximal operators have piecewise linear graphs and are perfectly known objects (Attouch et al., 2013; Beck, 2017). The proximal operator for $g_1(\mathbf{u}) = \lambda\alpha\|\mathbf{u}\|_1$, $\text{prox}_{g_1(\lambda\alpha)}(\mathbf{u}) = [|\mathbf{u}| - \lambda\alpha]_+ \text{sgn}(\mathbf{u})$ (the so-called soft thresholding function) is defined as

$$[|\mathbf{u}| - \lambda\alpha]_+ \text{sgn}(\mathbf{u}) = \begin{cases} \mathbf{u} - \lambda\alpha, & \text{if } \mathbf{u} \geq \lambda\alpha, \\ 0, & \text{if } |\mathbf{u}| < \lambda\alpha, \\ \mathbf{u} + \lambda\alpha, & \text{if } \mathbf{u} \leq -\lambda\alpha. \end{cases}$$

Hence, $\text{prox}_{g_1(\lambda\alpha)}(\mathbf{u})$ has a piecewise-linear graph and is semi-algebraic. The proximal operator for g_2 can be written as

$$\text{prox}_{g_2(\lambda(1-\alpha))}(\mathbf{u}) = \begin{cases} 0, & \text{if } |\mathbf{u}| < \sqrt{2\lambda(1-\alpha)}, \\ \mathbf{u}, & \text{if } |\mathbf{u}| > \sqrt{2\lambda(1-\alpha)}, \\ \{0, \mathbf{u}\}, & \text{if } |\mathbf{u}| = \sqrt{2\lambda(1-\alpha)}. \end{cases}$$

Clearly, $\text{prox}_{g_2(\lambda(1-\alpha))}(\mathbf{u})$ is also piecewise linear and semi-algebraic. Note that $\text{prox}_{g_2(\lambda(1-\alpha))}(\mathbf{u}) = \mathcal{H}_\nu(\mathbf{u})$ the so-called hard thresholding operator, is defined as

$$\mathcal{H}_\nu(\mathbf{u}) \equiv \begin{cases} 0, & \text{if } |\mathbf{u}| < \nu, \\ \mathbf{u}, & \text{if } |\mathbf{u}| > \nu, \\ \{0, \mathbf{u}\}, & \text{if } |\mathbf{u}| = \nu, \end{cases}$$

where $\nu = \sqrt{2\lambda(1-\alpha)}$. Here, $g_1 + g_2$ is also semi-algebraic.

Consequently, for any nonnegative real numbers λ and α , the function $f(\mathbf{b}) + \lambda\alpha\|\mathbf{u}\|_1 + \lambda(1-\alpha)\|\mathbf{u}\|_0 + \frac{\gamma}{2}\|\mathbf{b} - \mathbf{u} + \mathbf{m}\|_2^2 - \frac{\gamma}{2}\|\mathbf{m}\|_2^2$ is semi-algebraic. Hence, we conclude that the Lagrangian function in (20) is a KŁ function.

Since $\{\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}\}$ is bounded, there exists a subsequence $\{\mathbf{b}^{kl}, \mathbf{u}^{kl}, \mathbf{m}^{kl}\}$ converging to a stationary point $\{\mathbf{b}^*, \mathbf{u}^*, \mathbf{m}^*\}$, where $l \in \mathbb{N}$. Since the Lagrangian function in (20) is a KŁ function (using the lower semicontinuous property), we have

$$L_\gamma(\mathbf{b}^*, \mathbf{u}^*, \mathbf{m}^*) \leq \lim_{l \rightarrow \infty} L_\gamma(\mathbf{b}^{kl}, \mathbf{u}^{kl}, \mathbf{m}^{kl}). \quad (51)$$

In conclusion, all the conditions (a)-(c) in Theorem 1 hold.

B.3 Proof of Theorem 2

Starting with $r = 1$ (PRSM), we update \mathbf{b} , \mathbf{m} , and \mathbf{u} iteratively according to (24a - 24d)

$$\mathbf{b}^{(k+1)} := \underset{\mathbf{b}}{\text{argmin}} L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}), \quad (52a)$$

$$\mathbf{m}^{(k+\frac{1}{2})} := \mathbf{m}^{(k)} + \mathbf{b}^{(k+1)} - \mathbf{u}^{(k)}, \quad (52b)$$

$$\mathbf{u}^{(k+1)} := \underset{\mathbf{u}}{\text{argmin}} L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k+\frac{1}{2})}), \quad (52c)$$

$$\mathbf{m}^{(k+1)} := \mathbf{m}^{(k+\frac{1}{2})} + \mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)}. \quad (52d)$$

Proof (a): From (52a), since $\mathbf{b}^{(k+1)}$ minimizes $L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)})$ and Lagrangian is strongly convex with respect to the variable \mathbf{b} , (37) holds.

Next, using the augmented Lagrangian function in (20), we compute

$$L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k+1)}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{k+\frac{1}{2}}) = \gamma \left(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)} \right)^\top \left(\mathbf{m}^{(k+1)} - \mathbf{m}^{k+\frac{1}{2}} \right). \quad (53)$$

Now we rewrite (52d) as

$$\mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)} = \mathbf{m}^{(k+1)} - \mathbf{m}^{k+\frac{1}{2}}. \quad (54)$$

From (52c) and (29), we obtain

$$\nabla f(\mathbf{b}^{(k+1)}) + \gamma(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)} + \mathbf{m}^{k+\frac{1}{2}}) = 0. \quad (55)$$

Substituting (52d) into (55), we obtain (41). Again, from (52a) and (29), we obtain

$$\nabla f(\mathbf{b}^{(k+1)}) + \gamma(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k)} + \mathbf{m}^{(k)}) = 0. \quad (56)$$

Substituting (52b) into (56), we obtain

$$\nabla f(\mathbf{b}^{(k+1)}) = -\gamma \mathbf{m}^{k+\frac{1}{2}}. \quad (57)$$

Using (41),(54) and (57), (53) becomes

$$\gamma \left(\|\mathbf{m}^{(k+1)} - \mathbf{m}^{k+\frac{1}{2}}\|^2 \right) = \gamma \left(\left\| -\frac{1}{\gamma} \nabla f(\mathbf{b}^{(k+1)}) + \frac{1}{\gamma} \nabla f(\mathbf{b}^{(k+1)}) \right\|^2 \right) = 0. \quad (58)$$

Hence,

$$L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k+1)}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{k+\frac{1}{2}}) \leq 0. \quad (59)$$

Using (20), we have

$$L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{k+\frac{1}{2}}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) = \gamma \left(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k)} \right)^\top \left(\mathbf{m}^{k+\frac{1}{2}} - \mathbf{m}^{(k)} \right). \quad (60)$$

Next, we reformulate (52b) as

$$\mathbf{b}^{(k+1)} - \mathbf{u}^{(k)} = \mathbf{m}^{k+\frac{1}{2}} - \mathbf{m}^{(k)}. \quad (61)$$

Using (57) and (61), (60) becomes

$$\gamma \left(\|\mathbf{m}^{k+\frac{1}{2}} - \mathbf{m}^{(k)}\|^2 \right) = \gamma \left(\left\| -\frac{1}{\gamma} \nabla f(\mathbf{b}^{(k+1)}) + \frac{1}{\gamma} \nabla f(\mathbf{b}^{(k)}) \right\|^2 \right). \quad (62)$$

Therefore,

$$L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{k+\frac{1}{2}}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) \leq \frac{l_f^2}{\gamma} \left(\|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|^2 \right), \quad (63)$$

where $l_f \geq 0$ is a Lipschitz gradient of the function $f(\mathbf{b})$. From (52c) we have that

$$L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{k+\frac{1}{2}}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{k+\frac{1}{2}}) \leq 0. \quad (64)$$

Finally, combining (37), (59), (63) and (64), we get the desired inequality as follows

$$\begin{aligned} & L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k+1)}) - L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) \\ &= L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{(k+1)}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{k+\frac{1}{2}}) \\ &\quad + L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k+1)}, \mathbf{m}^{k+\frac{1}{2}}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{k+\frac{1}{2}}) \\ &\quad + L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{k+\frac{1}{2}}) - L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) \\ &\quad + L_\gamma(\mathbf{b}^{(k+1)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) - L_\gamma(\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}) \\ &\leq \left(\frac{l_f^2}{\gamma} - \frac{\rho}{2} \right) \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2 \\ &= -\delta_1 \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2, \end{aligned}$$

which is the sufficient decreasing condition (45).

Proof (b): The difference lies in some steps to show the boundedness of $\mathbf{u}^{(k)}$ and $\mathbf{m}^{(k)}$. The rest is the same as in the proof of Theorem 1(b). From (54) and (61) we obtain $\mathbf{u}^{(k+1)} = \mathbf{b}^{(k+1)} - \mathbf{m}^{(k+1)} + \mathbf{m}^{k+\frac{1}{2}}$ and $\mathbf{u}^{(k)} = \mathbf{b}^{(k+1)} - \mathbf{m}^{k+\frac{1}{2}} + \mathbf{m}^{(k)}$, respectively. Then

$$\begin{aligned} \|\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}\|_2^2 &= \|\mathbf{m}^{k+\frac{1}{2}} - \mathbf{m}^{(k+1)} + \mathbf{m}^{k+\frac{1}{2}} - \mathbf{m}^{(k)}\|_2^2 \\ &\leq \|\mathbf{m}^{(k+1)} - \mathbf{m}^{k+\frac{1}{2}}\|_2^2 + \|\mathbf{m}^{k+\frac{1}{2}} - \mathbf{m}^{(k)}\|_2^2, \end{aligned}$$

This inequality can be rewritten using (58) and (62) as

$$\|\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}\|_2^2 \leq \frac{l_f^2}{\gamma^2} \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2.$$

Consequently, we obtain

$$\sum_{k=0}^{\infty} \|\mathbf{u}^{(k+1)} - \mathbf{u}^{(k)}\|_2^2 < \infty. \quad (65)$$

Equation (65) implies that $\mathbf{u}^{(k)}$ is bounded. To show that $\mathbf{m}^{(k)}$ is bounded, we analyze the difference $\mathbf{m}^{(k+1)} - \mathbf{m}^{(k)}$ as follows $\mathbf{m}^{(k+1)} - \mathbf{m}^{(k)} = \mathbf{m}^{(k+1)} - \mathbf{m}^{k+\frac{1}{2}} + \mathbf{m}^{k+\frac{1}{2}} - \mathbf{m}^{(k)}$. Then, we obtain

$$\begin{aligned} \|\mathbf{m}^{(k+1)} - \mathbf{m}^{(k)}\|_2^2 &\leq \frac{l_f^2}{\gamma^2} \|\mathbf{b}^{(k+1)} - \mathbf{b}^{(k)}\|_2^2, \\ \sum_{k=0}^{\infty} \|\mathbf{m}^{(k+1)} - \mathbf{m}^{(k)}\|_2^2 &< \infty. \end{aligned} \quad (66)$$

Hence, $\mathbf{m}^{(k)}$ is bounded.

Proof (c): See the proof of Theorem 1 (c).

For the case $r \in (0, 1)$, all conditions are valid. Hence, for the sequences $\{\mathbf{b}^{(k)}, \mathbf{u}^{(k)}, \mathbf{m}^{(k)}\}_{k=0}^t$ generated by the SCPISM scheme (24a - 24d), and its Lagrangian given by (20), the three conditions from Theorem 1 (a)-(c) hold, achieving a worst-case convergence rate of $O(\frac{1}{k})$. Here, a worst-case $O(\frac{1}{k})$ convergence rate indicates that the solution's accuracy, based on specific criteria, improves gradually at a rate proportional to one divided by the number of iterations (k) within an iterative algorithm (He et al., 2014).

C Datasets

C.1 Simulated QTLMAS 2010 Dataset

The dataset comprises 3226 individuals across 5 generations, including 20 founders (5 males and 15 females), with two observed traits (responses): a quantitative trait and a binary trait (Szydlowski & Paczyńska, 2011). Each female mates once, producing approximately 30 progeny per birth. SNP data were simulated using a coalescent model on five autosomal chromosomes, each 100 Mbp long. A total of 10031 markers were generated, including 263 monomorphic SNPs and 9768 biallelic SNPs. The continuous quantitative trait is controlled by 9 major QTLs at fixed positions, including two pairs of epistatic genes, 3 maternally imprinted genes, and two additive major genes with phenotypic effects of -3 and 3. The additive genes are positioned at SNP indices 4354 and 5327, whereas the major epistatic locus is at SNP 931. Additionally, a dominance locus was positioned at SNP number 9212, with an effect of 5.00 assigned to the heterozygote and 5.01 to the upper homozygote. Moreover, an over-dominance locus was placed at SNP 9404, with an effect of 5.00 assigned to the heterozygote, -0.01 to the lower homozygote, and 0.01 to the upper homozygote. After filtering SNPs with $MAF < 0.01$, 9723 markers were retained and transformed into one-hot encoding, resulting in 29169 genomic markers. We used the quantitative trait in our study. Generations 1 to 4 (individuals 1 to 2326) were used for training, and generation 5 (individuals 2327 to 3226) served as test data.

C.2 Real Pig Dataset

The Pig dataset contains data from 3534 individuals, with high-density genotypes and phenotypes for five traits (Cleveland et al., 2012). Using the PorcineSNP60 chip, 52842 SNPs were assessed and filtered to 50282 based on a minor allele frequency threshold of < 0.01 . The chosen trait had a heritability of 0.58. After adjusting the phenotypic data and excluding individuals with missing data, the final dataset included 3152 individuals and was transformed into one-hot encoding, resulting in 150840 genomic markers.

C.3 Real Mice Dataset

This dataset comes from an experiment aimed at identifying and locating quantitative trait loci (QTLs) associated with various complex traits in a population of mice. The dataset contains 1814 individuals who were genotyped for 10346 polymorphic markers and two traits: body length (BL) and body mass index (BMI). In this study, we used BL trait. After transforming the data into one-hot encoding, the dataset resulted in 31038 genomic markers. This dataset is from the Wellcome Trust and is available in the R package BGLR (Pérez & de Los Campos, 2014).

D Implementation of Baseline Methods

For comparison purposes, we implement the LASSO (3) using the proximal ADMM and SCPRSM schemes (referred to as LASSO-ADMM and LASSO-SCPRSM, respectively) as

$$\begin{aligned}\mathbf{b}^{(k+1)} &:= \text{prox}_{f_\gamma}(\mathbf{u}^{(k)} - \mathbf{m}^{(k)}), \\ \mathbf{u}^{(k+1)} &:= \text{prox}_{g_\gamma}(\mathbf{b}^{(k+1)} + \mathbf{m}^{(k)}), \\ \mathbf{m}^{(k+1)} &:= \mathbf{m}^{(k)} + \mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)}\end{aligned}\tag{67}$$

and

$$\begin{aligned}\mathbf{b}^{(k+1)} &:= \text{prox}_{f_\gamma}(\mathbf{u}^{(k)} - \mathbf{m}^{(k)}), \\ \mathbf{m}^{(k+\frac{1}{2})} &:= \mathbf{m}^{(k)} + r(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k)}), \\ \mathbf{u}^{(k+1)} &:= \text{prox}_{g_\gamma}(\mathbf{b}^{(k+1)} + \mathbf{m}^{(k+\frac{1}{2})}), \\ \mathbf{m}^{(k+1)} &:= \mathbf{m}^{(k+\frac{1}{2})} + r(\mathbf{b}^{(k+1)} - \mathbf{u}^{(k+1)}),\end{aligned}\tag{68}$$

where $\text{prox}_{g_\gamma}(\mathbf{b}) = \mathcal{S}_\gamma(\mathbf{b})$.

Similarly, for EN-ADMM and EN-SCPRSM, $\text{prox}_{g_\gamma}(\mathbf{b}) = \frac{1}{1+\gamma\xi}\mathcal{S}_\gamma(\mathbf{b})$, where $\xi > 0$ is a linear combination of the L^1 and L^2 penalties (Parikh & Boyd, 2013). The closed-form proximal mappings of the SCAD (5) and MCP (6) penalty functions can be found in (Fan & Li, 2001; Liao et al., 2023; Wang & Liu, 2024; Yun et al., 2021). Here, we utilize the scaled versions

$$\text{prox}_{g_\gamma}(\mathbf{b}) = \text{prox}_{\text{scad}_\gamma}(\mathbf{b}) = \begin{cases} \mathcal{S}_{\gamma\lambda}(\mathbf{b}) & \text{if } |\mathbf{b}| \leq (1+\gamma)\lambda, \\ \frac{(a-1)(\mathbf{b}) - \text{sign}(\mathbf{b})a\lambda\gamma}{a-1-\gamma} & \text{if } (1+\gamma)\lambda < |\mathbf{b}| \leq a\lambda, \\ \mathbf{b} & \text{if } |\mathbf{b}| > a\lambda, \end{cases}\tag{69}$$

$$\text{prox}_{g_\gamma}(\mathbf{b}) = \text{prox}_{\text{mcp}_\gamma}(\mathbf{b}) = \begin{cases} \frac{a\gamma}{a\gamma-1}\mathcal{S}_{\gamma\lambda}(\mathbf{b}) & \text{if } |\mathbf{b}| \leq a\gamma\lambda, \\ \mathbf{b} & \text{otherwise,} \end{cases}\tag{70}$$

with respect to SCAD and MCP, respectively. All iterations of LASSO-ADMM, LASSO-SCPRSM, EN-ADMM, EN-SCPRSM, SCAD-ADMM, SCAD-SCPRSM, MCP-ADMM and MCP-SCPRSM terminate upon achieving convergence, defined by the condition $\|\mathbf{b}^{(k)} - \mathbf{u}^{(k)}\|_\infty \leq \beta(1 + \|\mathbf{m}^{(k)}\|_\infty)$, where the tolerance parameter β is set to 10^{-5} .