

COMBSS: Best Subset Selection via Continuous Optimization

Sarat Moka (✉ s.moka@unsw.edu.au)

UNSW Sydney

Benoit Liquet

Macquarie University

Houying Zhu

Macquarie University

Samuel Muller

Macquarie University

Research Article

Keywords: Linear regression, High-dimensional regression, Model selection, Variable selection

Posted Date: June 22nd, 2023

DOI: <https://doi.org/10.21203/rs.3.rs-3077764/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: No competing interests reported.

Version of Record: A version of this preprint was published at Statistics and Computing on February 12th, 2024. See the published version at <https://doi.org/10.1007/s11222-024-10387-8>.

COMBSS: Best Subset Selection via Continuous Optimization

Sarat Moka^{1,2*}, Benoit Liquet^{2,3}, Houying Zhu² and Samuel Muller^{2,1}

^{1*}School of Mathematics and Statistics, The University of New South Wales, Sydney, NSW, Australia.

²School of Mathematical and Physical Sciences, Macquarie University, Sydney, NSW, Australia.

³Laboratoire de Mathématiques et de leurs Applications, Université de Pau et des Pays de l'Adour, Pau, France.

⁴School of Mathematics and Statistics, The University of Sydney, Sydney, NSW, Australia.

*Corresponding author(s). E-mail(s): s.moka@unsw.edu.au;

Contributing authors: benoit.liquet-weiland@mq.edu.au; houying.zhu@mq.edu.au;
samuel.muller@mq.edu.au;

Abstract

The problem of best subset selection in linear regression is considered with the aim to find a fixed size subset of features that best fits the response. This is particularly challenging when the total available number of features is very large compared to the number of data samples. Existing optimal methods for solving this problem tend to be slow while fast methods tend to have low accuracy. Ideally, new methods perform best subset selection faster than existing optimal methods but with comparable accuracy, or, being more accurate than methods of comparable computational speed. Here, we propose a novel continuous optimization method that identifies a solution path, a small set of models of varying size, that consists of candidates for the best subset in linear regression. Our method turns out to be very fast, making the best subset selection possible when the number of features is well in excess of thousands. Because of the outstanding overall performance, framing the best subset selection challenge as a continuous optimization problem opens new research directions for feature extraction for a large variety of regression models.

Keywords: Linear regression, High-dimensional regression, Model selection, Variable selection

1 Introduction

Recent developments in information technology have enabled the collection of high-dimensional complex data in engineering, economics, finance, biology, health sciences and other fields (Fan and Li, 2006). In high-dimensional data, the number of features is large and often far higher than the number of collected data samples. In

many applications, it is desirable to find a parsimonious best subset of predictors so that the resulting model has desirable prediction accuracy (Müller and Welsh, 2010; Fan and Lv, 2010; Miller, 2019). This article is recasting the challenge of best subset selection in linear regression as a novel continuous optimization problem. We show that this reframing has enormous potential and substantially advances research into larger dimensional and exhaustive feature selection in

regression, making available technology that can reliably and exhaustively select variables when the total number of variables is well in excess of thousands.

Here, we aim to develop a method that performs best subset selection and an approach that is faster than existing exhaustive methods while having comparable accuracy, or, that is more accurate than other methods of comparable computational speed.

Consider the linear regression model of the form $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$, where $\mathbf{y} = (y_1, \dots, y_n)^\top$ is an n -dimensional known response vector, X is a known design matrix of dimension $n \times p$ with $x_{i,j}$ indicating the i th observation of the j th explanatory variable, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$ is the p -dimensional vector of unknown regression coefficients, and $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^\top$ is a vector of unknown errors, unless otherwise specified, assumed to be independent and identically distributed. Best subset selection is a classical problem that aims to solve,

$$\begin{aligned} & \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\text{minimize}} && \frac{1}{n} \|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 \\ & \text{subject to} && \|\boldsymbol{\beta}\|_0 = k, \end{aligned} \quad (1)$$

for a given k , where $\|\cdot\|_2$ is the \mathcal{L}_2 -norm, $\|\boldsymbol{\beta}\|_0 = \sum_{j=1}^p I(\beta_j \neq 0)$ is the number of non-zero elements in $\boldsymbol{\beta}$, and $I(\cdot)$ is the indicator function. For ease of presentation, we assume that all columns of X are subject to selection, but generalizations are immediate (see Remark 2 for more details).

Exact methods for solving (1) are typically executed by first writing solutions for low-dimensional problems and then selecting the best solution over these. To see this, for any binary vector $\mathbf{s} = (s_1, \dots, s_p)^\top \in \{0, 1\}^p$, let $X_{[\mathbf{s}]}$ be the matrix of size $n \times |\mathbf{s}|$ created by keeping only columns j of X for which $s_j = 1$, where $j = 1, \dots, p$. Then, for any k , in the exact best subset selection, an optimal \mathbf{s} can be found by solving the problem,

$$\begin{aligned} & \underset{\mathbf{s} \in \{0, 1\}^p}{\text{minimize}} && \frac{1}{n} \|\mathbf{y} - X_{[\mathbf{s}]} \widehat{\boldsymbol{\beta}}_{[\mathbf{s}]}\|_2^2 \\ & \text{subject to} && |\mathbf{s}| = k, \end{aligned} \quad (2)$$

where $\widehat{\boldsymbol{\beta}}_{[\mathbf{s}]}$ is a low-dimensional least squares estimate of elements of $\boldsymbol{\beta}$ with indices corresponding

to non-zero elements of \mathbf{s} , given by

$$\widehat{\boldsymbol{\beta}}_{[\mathbf{s}]} = (X_{[\mathbf{s}]}^\top X_{[\mathbf{s}]})^\dagger X_{[\mathbf{s}]}^\top \mathbf{y}, \quad (3)$$

where A^\dagger denotes the pseudo-inverse of a matrix A . Both (1) and (2) are essentially solving the same problem, because $\widehat{\boldsymbol{\beta}}_{[\mathbf{s}]}$ is the least squares solution when constrained so that $I(\beta_j \neq 0) = s_j$ for all $j = 1, \dots, p$.

It is well-known that solving the exact optimization problem above is in general non-deterministic polynomial-time hard (Natarajan, 1995). For instance, a popular exact method called *leaps-and-bounds* (Furnival and Wilson, 2000) is currently practically useful only for values of p smaller than 30 (Tarr et al, 2018). To overcome this difficulty, the relatively recent method by Bertsimas et al (2016) elegantly reformulates the best subset selection problem (1) as a mixed integer optimization and demonstrates that the problem can be solved for p much larger than 30 using modern mixed integer optimization solvers such as in the commercial software Gurobi (Gurobi Optimization, limited liability company, 2022) (which is not freely available except for an initial short period). As the name suggests, the formulation of mixed integer optimization has both continuous and discrete constraints. Although, the mixed integer optimization approach is faster than the exact methods for large p , its implementation via Gurobi remains slow from a practical point of view (Hazimeh and Mazumder, 2020a).

Due to computational constraints of mixed integer optimization, other popular existing methods for best subset selection are still very common in practice, these include forward stepwise selection, the *least absolute shrinkage and selection operator* (generally known as the Lasso), and their variants. Forward stepwise selection follows a greedy approach, starting with an empty model (or intercept-only model), and iteratively adding the variable that is most suitable for inclusion (Efroymson, 1966; Hocking and Leslie, 1967). On the other hand, the Lasso (Tibshirani, 1996) solves a convex relaxation of the highly non-convex best subset selection problem by replacing the discrete \mathcal{L}_0 -norm $\|\boldsymbol{\beta}\|_0$ in (1) with the \mathcal{L}_1 -norm $\|\boldsymbol{\beta}\|_1$. This clever relaxation makes the Lasso fast, significantly faster than mixed-integer optimization solvers. However, it is important to note that

Lasso solutions typically do not yield the best subset solution (Hazimeh and Mazumder, 2020a; Zhu et al, 2020) and in essence solve a different problem than exhaustive best subset selection approaches. In summary, there exists a trade-off between speed and accuracy when selecting an existing best subset selection method.

With the aim to develop a method that performs best subset selection as fast as the existing fast methods without compromising the accuracy, in this paper, we design COMBSS, a novel *continuous optimization method towards best subset selection*.

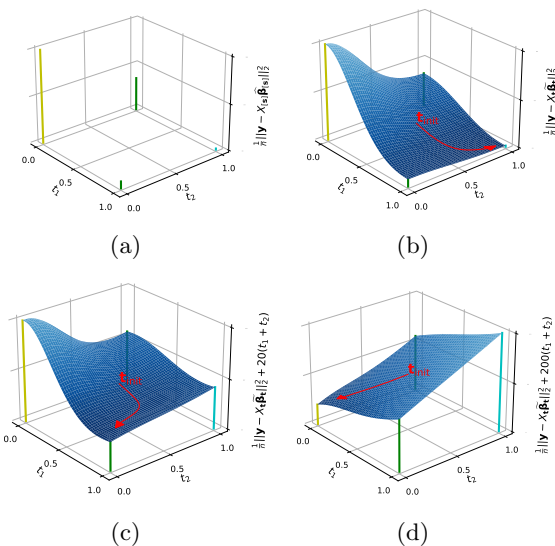


Fig. 1: Illustration of the workings of COMBSS for an example data with $p = 2$. Plot (a) shows the objective function of the exact method (2) for $\mathbf{s} \in \{0, 1\}^2$. Observe that the best subsets correspond to $k = 0$, $k = 1$, and $k = 2$ are $(1, 1)^\top$, $(1, 0)^\top$, and $(0, 0)^\top$, respectively. Plots (b) - (d) show the objective function of our optimization method (4) for different values of the parameter λ . In each of these three plots, the curve (in red) shows the execution of a basic gradient descent algorithm that, starting at the initial point $\mathbf{t}_{init} = (0.5, 0.5)^\top$, converges towards the best subsets of sizes 0, 1, and 2, respectively.

Our continuous optimization method can be described as follows. Instead of the binary vector

space $\{0, 1\}^p$ as in the exact methods, we consider the whole hyper-cube $[0, 1]^p$ and for each $\mathbf{t} \in [0, 1]^p$, we consider a new estimate $\tilde{\boldsymbol{\beta}}_{\mathbf{t}}$ (defined later in Section 2) so that we have the following well-defined continuous extension of the exact problem (2):

$$\underset{\mathbf{t} \in [0, 1]^p}{\text{minimize}} \quad \frac{1}{n} \|\mathbf{y} - X_{\mathbf{t}} \tilde{\boldsymbol{\beta}}_{\mathbf{t}}\|_2^2 + \lambda \sum_{j=1}^p t_j, \quad (4)$$

where $X_{\mathbf{t}}$ is obtained from X by multiplying the j th column of X by t_j for each $j = 1, \dots, p$, and the tuning parameter λ controls the sparsity of the solution obtained, analogous to selecting the best k in the exact optimization. Our construction of $\tilde{\boldsymbol{\beta}}_{\mathbf{t}}$ guarantees that $\|\mathbf{y} - X_{\mathbf{s}} \tilde{\boldsymbol{\beta}}_{\mathbf{s}}\|_2 = \|\mathbf{y} - X_{[\mathbf{s}]} \hat{\boldsymbol{\beta}}_{[\mathbf{s}]}\|_2$ at the corner points \mathbf{s} of the hypercube $[0, 1]^p$ where $\hat{\boldsymbol{\beta}}_{[\mathbf{s}]}$ exists, and the new objective function $\|\mathbf{y} - X_{\mathbf{t}} \tilde{\boldsymbol{\beta}}_{\mathbf{t}}\|_2^2$ is smooth over the hypercube.

While COMBSS aims to find sets of models that are candidates for the best subset of variables, an important property is that it has no discrete constraints, unlike the exact optimization problem (2) or the mixed integer optimization formulation. As a consequence, our method can take advantage of standard continuous optimization methods, such as gradient descent methods, by starting at an interior point on the hypercube $[0, 1]^p$ and iteratively moving towards a corner that minimizes the objective function. See Fig. 1 for an illustration of our method. In the implementation, we move the box constrained problem (4) to an equivalent unconstrained problem so that the gradient descent method can run without experiencing boundary issues.

The rest of the paper is organized as follows: In Section 2, we describe the mathematical framework of the proposed method COMBSS. In Section 3, we first establish the continuity of the objective functions involved in COMBSS, and then we derive expressions for their gradients. These gradients are exploited for conducting continuous optimization. Complete details of COMBSS algorithm are presented in Section 4. In Section 5, we discuss roles of the tuning parameters that control the surface shape of the objective functions and the sparsity of the solutions obtained. Section 6 provides steps for efficient implementation of COMBSS using some popular linear algebra techniques. Simulation results

comparing COMBSS with existing popular methods are presented in Section 7. We conclude the paper with some brief remarks in Section 8. Proofs all our theoretical results are provided in Appendix A.

2 Continuous Extension of Best Subset Selection Problem

To see our continuous extension of the exact best subset selection optimization problem (2), for $\mathbf{t} = (t_1, \dots, t_p)^\top \in [0, 1]^p$, define $T_{\mathbf{t}} = \text{Diag}(\mathbf{t})$, the diagonal matrix with the diagonal elements being t_1, \dots, t_p , and let $X_{\mathbf{t}} = XT_{\mathbf{t}}$. For a fixed constant $\delta > 0$, define

$$L_{\mathbf{t}} = L_{\mathbf{t}}(\delta) = \frac{1}{n} \left[X_{\mathbf{t}}^\top X_{\mathbf{t}} + \delta (I - T_{\mathbf{t}}^2) \right], \quad (5)$$

where we suppress δ for ease of reading. Intuitively, $L_{\mathbf{t}}$ can be seen as a ‘convex combination’ of the matrices $(X^\top X)/n$ and $\delta I/n$, because $X_{\mathbf{t}}^\top X_{\mathbf{t}} = T_{\mathbf{t}} X^\top X T_{\mathbf{t}}$ and thus

$$L_{\mathbf{t}} = T_{\mathbf{t}} \frac{X^\top X}{n} T_{\mathbf{t}} + (I - T_{\mathbf{t}}) \frac{\delta I}{n} (I - T_{\mathbf{t}}). \quad (6)$$

Using this notation, now define

$$\tilde{\boldsymbol{\beta}}_{\mathbf{t}} = L_{\mathbf{t}}^\dagger \left(\frac{X_{\mathbf{t}}^\top \mathbf{y}}{n} \right), \quad \mathbf{t} \in [0, 1]^p. \quad (7)$$

We need $L_{\mathbf{t}}^\dagger$ in (7) so that $\tilde{\boldsymbol{\beta}}_{\mathbf{t}}$ is defined for all $\mathbf{t} \in [0, 1]^p$. However, from the way we conduct optimization, we need to compute $\tilde{\boldsymbol{\beta}}_{\mathbf{t}}$ only for $\mathbf{t} \in [0, 1]^p$. We later show in Theorem 1 that for all $\mathbf{t} \in [0, 1]^p$, $L_{\mathbf{t}}$ is invertible and thus in the implementation of our method, $\tilde{\boldsymbol{\beta}}_{\mathbf{t}}$ always takes the form $\tilde{\boldsymbol{\beta}}_{\mathbf{t}} = L_{\mathbf{t}}^{-1} X_{\mathbf{t}}^\top \mathbf{y}/n$, eliminating the need to compute any computationally expensive pseudo-inverse.

With the support of these observations, an immediate well-defined generalization of the best subset selection problem (1) is

$$\begin{aligned} & \underset{\mathbf{t} \in [0, 1]^p}{\text{minimize}} && \frac{1}{n} \|\mathbf{y} - X_{\mathbf{t}} \tilde{\boldsymbol{\beta}}_{\mathbf{t}}\|_2^2 \\ & \text{subject to} && \sum_{j=1}^p t_j = k. \end{aligned} \quad (8)$$

Instead of solving the constrained problem (8), by defining a Lagrangian function

$$f_\lambda(\mathbf{t}) = \frac{1}{n} \|\mathbf{y} - X_{\mathbf{t}} \tilde{\boldsymbol{\beta}}_{\mathbf{t}}\|_2^2 + \lambda \sum_{j=1}^p t_j, \quad (9)$$

for a tunable parameter $\lambda > 0$, we aim to solve

$$\underset{\mathbf{t} \in [0, 1]^p}{\text{minimize}} f_\lambda(\mathbf{t}). \quad (10)$$

By defining $g_\lambda(\mathbf{w}) = f_\lambda(\mathbf{t}(\mathbf{w}))$, we reformulate the box constrained problem (10) into an equivalent unconstrained problem,

$$\underset{\mathbf{w} \in \mathbb{R}^p}{\text{minimize}} g_\lambda(\mathbf{w}), \quad (11)$$

where the mapping $\mathbf{t} = \mathbf{t}(\mathbf{w})$ is

$$t_j(w_j) = 1 - \exp(-w_j^2), \quad j = 1, \dots, p. \quad (12)$$

The unconstrained optimization problem (11) is equivalent to the box constrained problem (10), because $1 - \exp(-u^2) < 1 - \exp(-v^2)$ if and only if $u^2 < v^2$, for any $u, v \in \mathbb{R}$.

Remark 1 The non-zero parameter δ is important in the expression of the proposed estimator $\tilde{\boldsymbol{\beta}}_{\mathbf{t}}$, as in (7), not only to make $L_{\mathbf{t}}$ invertible for $\mathbf{t} \in [0, 1]^p$, but also to make the surface of $f_\lambda(\mathbf{t})$ to have smooth transitions from one corner to another over the hypercube. For example consider a situation where $X^\top X$ is invertible. Then, for any interior point $\mathbf{t} \in (0, 1)^p$, since $T_{\mathbf{t}}^{-1}$ exists, the optimal solution to $\min_{\boldsymbol{\beta}} \|\mathbf{y} - X_{\mathbf{t}} \boldsymbol{\beta}\|_2^2/n$ after some simplification is $T_{\mathbf{t}}^{-1} (X^\top X) X_{\mathbf{t}}^\top \mathbf{y}$. As a result, the corresponding minimum loss is $\|\mathbf{y} - X (X^\top X) X_{\mathbf{t}}^\top \mathbf{y}\|_2^2/n$, which is a constant for all \mathbf{t} over the interior of the hypercube. Hence, the surface of the loss function would have jumps at the borders while being flat over the interior of the hypercube. Clearly, such a loss function is not useful for conducting continuous optimization.

Remark 2 The proposed method and the corresponding theoretical results presented in this paper easily extend to linear models with intercept term. More generally, if we want to keep some features in the model, say features $j = 1, 2$, and 4, then we enforce $t_j = 1$ for $j = 1, 2, 4$, and conduct subset selection only over the remaining features by taking $\mathbf{t} = (1, 1, t_3, 1, t_5, \dots, t_p)^\top$ and optimize t_3, t_5, \dots, t_p .

3 Continuity and Gradients of the Objective Function

In this section, we first prove that the objective function $g_\lambda(\mathbf{w})$ of the unconstrained optimization problem (11) is continuous on \mathbb{R}^p and then we derive its gradients. En-route, we also establish the relationship between $\widehat{\beta}_{[s]}$ and $\widetilde{\beta}_{\mathbf{t}}$ which are respectively defined by (3) and (7). This relationship is useful in understanding the relationship between our method and the exact optimization (2).

Theorem 1 shows that for all $\mathbf{t} \in [0, 1]^p$, the matrix $L_{\mathbf{t}}$, which is defined in (5), is symmetric positive-definite and hence invertible.

Theorem 1 *For any $\mathbf{t} \in [0, 1]^p$, $L_{\mathbf{t}}$ is symmetric positive-definite and $\widetilde{\beta}_{\mathbf{t}} = L_{\mathbf{t}}^{-1} X_{\mathbf{t}}^\top \mathbf{y} / n$.*

Recall that $\widehat{\beta}_{[s]}$ exists only on some corner points \mathbf{s} of the hyper-cube $[0, 1]^p$, where the matrix $X_{[s]}^\top X_{[s]}$ is invertible. Theorem 2 establishes a relationship between $\widetilde{\beta}_{\mathbf{s}}$ and $\widehat{\beta}_{[s]}$ at these corner points $\mathbf{s} \in \{0, 1\}^p$. Towards this, for any point $\mathbf{s} \in \{0, 1\}^p$ and a vector $\mathbf{u} \in \mathbb{R}^p$, we write $(\mathbf{u})_+$ (respectively, $(\mathbf{u})_0$) to denote the *sliced* vector of dimension $|\mathbf{s}|$ (respectively, $p - |\mathbf{s}|$) created from \mathbf{u} by removing all its elements with the indices j with $s_j = 0$ (respectively, with $s_j > 0$). For instance, if $\mathbf{u} = (2, 3, 4, 5)^\top$ and $\mathbf{s} = (1, 0, 1, 0)^\top$, then $(\mathbf{u})_+ = (2, 4)$ and $(\mathbf{u})_0 = (3, 5)$.

Theorem 2 *For any $\mathbf{s} \in \{0, 1\}^p$, $(\widetilde{\beta}_{\mathbf{s}})_+ = \widehat{\beta}_{[s]}$ and $(\widetilde{\beta}_{[s]})_0 = \mathbf{0}$. Furthermore, we have $X_{[s]} \widehat{\beta}_{[s]} = X_{\mathbf{s}} \widetilde{\beta}_{\mathbf{s}}$.*

As an immediate consequence of Theorem 2, we have $\|\mathbf{y} - X_{[s]} \widehat{\beta}_{[s]}\|_2^2 = \|\mathbf{y} - X_{\mathbf{s}} \widetilde{\beta}_{\mathbf{s}}\|_2^2$. Therefore, the objective function of the exact optimization problem (2) is identical to the objective function of its extended optimization problem (8) (with $\lambda = 0$) at the corner points $\mathbf{s} \in \{0, 1\}^p$.

Our next result, Theorem 3, shows that $f_\lambda(\mathbf{t})$ is a continuous function on $[0, 1]^p$.

Theorem 3 *The function $f_\lambda(\mathbf{t})$ defined in (9) is continuous over $[0, 1]^p$ in the sense that for any sequence $\mathbf{t}^{(1)}, \mathbf{t}^{(2)}, \dots \in [0, 1]^p$ converging to $\mathbf{t} \in [0, 1]^p$, the*

limit $\lim_{l \rightarrow \infty} f_\lambda(\mathbf{t}^{(l)})$ exists and

$$f_\lambda(\mathbf{t}) = \lim_{l \rightarrow \infty} f_\lambda(\mathbf{t}^{(l)}).$$

Corollary 1 establishes the continuity of g_λ on \mathbb{R}^p . This is a simple consequence of Theorem 3, because from the definition, $g_\lambda(\mathbf{w}) = f_\lambda(\mathbf{t}(\mathbf{w}))$ with $\mathbf{t}(\mathbf{w}) = \mathbf{1} - \exp(-\mathbf{w} \odot \mathbf{w})$. Here and afterwards, in an expression with vectors, $\mathbf{1}$ denotes a vector of all ones of appropriate dimension, \odot denotes the element-wise product of two vectors, and the exponential function, $\exp(\cdot)$, is also applied element-wise.

Corollary 1 The objective function $g_\lambda(\mathbf{w})$ is continuous at every point $\mathbf{w} \in \mathbb{R}^p$.

As mentioned earlier, our continuous optimization method uses a gradient descent method to solve the problem (11). Towards that we need to obtain the gradients of $g_\lambda(\mathbf{w})$. Theorem 4 provides an expression of the gradient $\nabla g_\lambda(\mathbf{w})$.

Theorem 4 *For every $\mathbf{w} \in \mathbb{R}^p$, with $\mathbf{t} = \mathbf{t}(\mathbf{w})$ is defined by (12),*

$$\nabla f_\lambda(\mathbf{t}) = \zeta_{\mathbf{t}} + \lambda \mathbf{1}, \quad \mathbf{t} \in (0, 1)^p,$$

and

$$\nabla g_\lambda(\mathbf{w}) = (\zeta_{\mathbf{t}} + \lambda \mathbf{1}) \odot (2\mathbf{w} \odot \exp(-\mathbf{w} \odot \mathbf{w})), \quad \mathbf{w} \in \mathbb{R}^p,$$

where

$$\zeta_{\mathbf{t}} = 2 \left(\widetilde{\beta}_{\mathbf{t}} \odot (\mathbf{a}_{\mathbf{t}} - \mathbf{d}_{\mathbf{t}}) \right) - 2 (\mathbf{b}_{\mathbf{t}} \odot \mathbf{c}_{\mathbf{t}}), \quad (13)$$

with

$$\mathbf{a}_{\mathbf{t}} = \frac{1}{n} [X^\top X (\mathbf{t} \odot \widetilde{\beta}_{\mathbf{t}}) - X^\top \mathbf{y}],$$

$$\mathbf{b}_{\mathbf{t}} = \mathbf{a}_{\mathbf{t}} - n^{-1} \delta (\mathbf{t} \odot \widetilde{\beta}_{\mathbf{t}}),$$

$$\mathbf{c}_{\mathbf{t}} = L_{\mathbf{t}}^{-1} (\mathbf{t} \odot \mathbf{a}_{\mathbf{t}}), \quad \text{and}$$

$$\mathbf{d}_{\mathbf{t}} = \frac{1}{n} [X^\top X - \delta I] (\mathbf{t} \odot \mathbf{c}_{\mathbf{t}}).$$

Figure 2 illustrates the typical convergence behavior of \mathbf{t} for an example dataset during the execution of a basic gradient descent algorithm for minimizing $g_\lambda(\mathbf{w})$ using the gradient ∇g_λ given in Theorem 4. Here, \mathbf{w} is mapped to \mathbf{t} via (12) at each iteration.

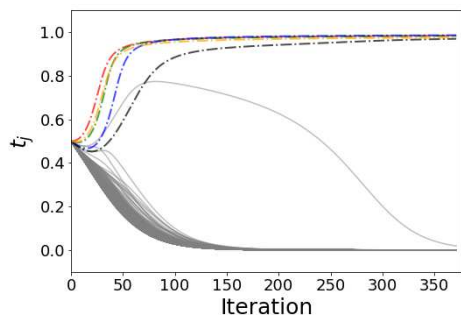


Fig. 2: Convergence of \mathbf{t} for a dataset with $p = 1000$ and $n = 100$ during the execution of basic gradient descent. Solid lines correspond to $\beta_j = 0$ and remaining 5 curves (with line style $-\cdot-$) correspond to $\beta_j \neq 0$. The dataset is generated using Beta-type 1 model presented in Section 7 of preprint of this work (Moka et al, 2022) with $k_0 = 5$, $n = 100$, $p = 1000$, $\rho = 0.8$, and signal-to-noise ratio of 5. The parameters $\lambda = 0.1$ and $\delta = n$; see Section 5 for more discussion on how to choose λ and δ .

4 Subset Selection Algorithms

Our algorithm COMBSS as stated in Algorithm 1, takes the data $[X, \mathbf{y}]$, tuning parameters δ, λ , and an initial point $\mathbf{w}^{(0)}$ as input, and returns either a single model or multiple models of different sizes as output. It is executed in three steps.

In Step 1, $\text{GradientDescent}(\mathbf{w}^{(0)}, \nabla g_\lambda)$ calls a gradient descent method, such as the well known *adam* optimizer, for minimizing the objective function $g_\lambda(\mathbf{w})$, which takes $\mathbf{w}^{(0)}$ as the initial point and uses the gradient function ∇g_λ for updating the vector \mathbf{w} in each iteration; see, for example, Kochenderfer and Wheeler (2019) for a review of popular gradient based optimization methods. It terminates when a predefined termination condition is satisfied and returns the sequence $\mathbf{w}_{\text{path}} = (\mathbf{w}^{(0)}, \mathbf{w}^{(1)}, \dots)$ of all the points \mathbf{w} visited during its execution, where $\mathbf{w}^{(l)}$ denotes the point obtained in the l th iteration. Usually, a robust termination condition is to terminate when the change in \mathbf{w} (or, equivalently, in $\mathbf{t}(\mathbf{w})$) is significantly small over a set of consecutive iterations.

Selecting the initial point $\mathbf{w}^{(0)}$ requires few considerations. From Theorem 4, for any $j = 1, \dots, p$, we have $t_j(w_j) = 0$ if and only if $w_j = 0$ and $\partial g_\lambda(\mathbf{w})/\partial w_j = 0$ if $w_j = 0$. Hence,

Algorithm 1 COMBSS ($X, \mathbf{y}, \delta, \lambda, \mathbf{w}^{(0)}$)

- 1: $\mathbf{w}_{\text{path}} \leftarrow \text{GradientDescent}(\mathbf{w}^{(0)}, \nabla g_\lambda)$
 - 2: Obtain \mathbf{t}_{path} from \mathbf{w}_{path} using the map (12)
 - 3: $\mathcal{M} \leftarrow \text{SubsetMap}(\mathbf{t}_{\text{path}})$
 - 4: **return** \mathcal{M}
-

if we start the gradient descent algorithm with $w_j^{(0)} = 0$ for some j , both w_j and t_j can continue to take 0 forever. As a result, we might not learn the optimal value for w_j (or, equivalently for t_j). Thus, it is important to select all the elements of $\mathbf{w}^{(0)}$ away from 0.

Consider the second argument, ∇g_λ , in the gradient descent method. From Theorem 4, observe that computing the gradient $\nabla g_\lambda(\mathbf{w})$ involves finding the values of the expression of the form $L_t^{-1}\mathbf{u}$ twice, first for computing $\tilde{\beta}_t$ (using (7)) and then for computing the vector \mathbf{c}_t (defined in Theorem 4). Since L_t is of dimension $p \times p$, computing the matrix inversion L_t^{-1} can be computationally demanding particularly in high-dimensional cases ($n < p$), where p can be very large; see, for example, Golub and Van Loan (1996). Since L_t is invertible, observe that $L_t^{-1}\mathbf{u}$ is the unique solution of the linear equation $L_t\mathbf{z} = \mathbf{u}$. In Section 6, we first use the well-known Woodbury matrix identity to convert this p -dimensional linear equation problem to an n -dimensional linear equation problem, which is then solved using the conjugate gradient method, a popular linear equation solver. Moreover, again from Theorem 4, notice that $\nabla g_\lambda(\mathbf{w})$ depends on both the tuning parameters δ and λ . Specifically, δ is required for computing L_t and λ is used in the penalty term $\lambda \sum_{j=1}^p t_j$ of the objective function. In Section 5 we provide more details on the roles of these two parameters and instructions on how to choose them.

In Step 2, we obtain the sequence $\mathbf{t}_{\text{path}} = (\mathbf{t}^{(0)}, \mathbf{t}^{(1)}, \dots)$ from \mathbf{w}_{path} by using the map (12), that is, $\mathbf{t}^{(l)} = \mathbf{t}(\mathbf{w}^{(l)}) = \mathbf{1} - \exp(-\mathbf{w}^{(l)} \odot \mathbf{w}^{(l)})$ for each l .

Finally, in Step 3, $\text{SubsetMap}(\mathbf{t}_{\text{path}})$ takes the sequence \mathbf{t}_{path} as input to find a set of models \mathcal{M} correspond to the input parameter λ . In the following subsections, we describe two versions of SubsetMap .

4.1 Subset Map Version 1

One simple implementation of `SubsetMap` is stated as Algorithm 2 which we call `SubsetMapV1` (where V1 stands for *version 1*) and it requires only the final point in the sequence \mathbf{t}_{path} and returns only one model using a predefined threshold parameter $\tau \in [0, 1)$.

Algorithm 2 `SubsetMapV1` ($\mathbf{t}_{\text{path}}, \tau$)

- 1: Take \mathbf{t} to be the final point of \mathbf{t}_{path}
 - 2: **for** $j = 1$ to $j = p$ **do**
 - 3: $s_j \leftarrow \mathbb{I}(t_j > \tau)$
 - 4: **end for**
 - 5: **return** $\mathbf{s} = (s_1, \dots, s_p)^\top$
-

Due to the tolerance allowed by the termination condition of the gradient descent, some w_j in the final point of \mathbf{w}_{path} can be almost zero but not exactly zero, even though they are meant to converge to zero. As a result, the corresponding t_j also take values close to zero but not exactly zero because of the mapping from \mathbf{w} to \mathbf{t} . Therefore, the threshold τ helps in mapping the insignificantly small t_j to 0 and all other t_j to 1. In practice, we call `ModelSelection` ($X, \mathbf{y}, \delta, \lambda, \mathbf{w}^{(0)}$) for each λ over a grid of values. When `SubsetMapV1` is used, larger the value of λ , higher the sparsity in the resulting model \mathbf{s} . Thus, we can control the sparsity of the output model using λ . Since we only care about the last point in \mathbf{t}_{path} in this version, an intuitive option for initialization is to take $\mathbf{w}^{(0)}$ to be such that $\mathbf{t}(\mathbf{w}^{(0)}) = (1/2, \dots, 1/2)^\top$, the mid-point on the hypercube $[0, 1]^p$, as it is at an equal distance from all the corner points, of which one is the (unknown) target solution of the best subset selection problem.

In the preprint Moka et al (2022), we demonstrated the efficacy of COMBSS using `SubsetMapV1` in predicting the true model of the data. In almost all the cases, we observe superior performance of COMBSS in comparison to existing popular methods.

4.2 Subset Map Version 2

Ideally, there is a value of λ for each k such that the output model \mathbf{s} obtained by `SubsetMapV1` has exactly k non-zero elements. However, when the ultimate goal is to find a best suitable model \mathbf{s}

for a given $k \leq q$ such that $|\mathbf{s}| = k$, for some $q \ll \min(n, p)$, since λ is selected over a grid, we might not obtain any model for some values of k . Furthermore, for a given size k , if there are two models with almost the same mean square error, then the optimization may have difficulty in distinguishing them. Addressing this difficulty may involve fine tuning of hyper-parameters of the optimization algorithm.

To overcome these challenges without any hyper-parameter tuning and reduce the reliance on the parameter λ , we consider the other points in \mathbf{t}_{path} . In particular, we propose a more optimal implementation of `SubsetMap`, which we call `SubsetMapV2` and is stated as Algorithm 3. The key idea of this version is that as the gradient descent progresses over the surface of $f_\lambda(\mathbf{t})$, it can point towards some corners of the hypercube $[0, 1]^p$ before finally moving towards the final corner. Considering all these corners, we can refine the results. Specifically, this version provides for each λ a model for every $k \leq q$. In this implementation, λ is seen as a parameter that allows us to explore the surface of $f_\lambda(\mathbf{t})$ rather than as a sparsity parameter.

Algorithm 3 `SubsetMapV2` (\mathbf{t}_{path})

- 1: $\mathcal{M}_k \leftarrow \{\}$ for each $k \leq q$
 - 2: **for** each $\mathbf{t} = (t_1, \dots, t_p)^\top$ in \mathbf{t}_{path} **do**
 - 3: Let $t_{j_1}, t_{j_2}, \dots, t_{j_q}$ be the q largest elements of \mathbf{t} in the descending order
 - 4: **for** $k = 1$ to q **do**
 - 5: Take $\mathbf{s}_k \in \{0, 1\}^p$ with non-zero elements only at j_1, \dots, j_k
 - 6: $\mathcal{M}_k \leftarrow \mathcal{M}_k \cup \{\mathbf{s}_k\}$
 - 7: **end for**
 - 8: **end for**
 - 9: **for** $k = 1$ to $k = q$ **do**
 - 10: $\mathbf{s}_k^* \leftarrow \operatorname{argmin}_{\mathbf{s} \in \mathcal{M}_k} \frac{1}{n} \|\mathbf{y} - X_{[s]} \hat{\boldsymbol{\beta}}_{[s]}\|_2^2$
 - 11: **end for**
 - 12: **return** $\mathcal{M} = \{\mathbf{s}_1^*, \dots, \mathbf{s}_q^*\}$
-

For the execution of `SubsetMapV2`, we start at Step 1 with an empty set of models \mathcal{M}_k for each $k \leq q$. In Step 2, for each \mathbf{t} in \mathbf{t}_{path} , we consider the sequence of indices j_1, \dots, j_q such that $t_{j_1} \geq t_{j_2} \geq \dots \geq t_{j_q}$. Then, for each $k \leq q$, we take \mathbf{s}_k to be a binary vector with 1's only at j_1, \dots, j_k

and add \mathbf{s}_k to the set \mathcal{M}_k . With this construction, it is clear that \mathcal{M}_k consists of models of size k , of which we pick a best candidate \mathbf{s}_k^* as show at Step 3. Finally, the algorithm returns the set consists of $\mathbf{s}_1^*, \dots, \mathbf{s}_q^*$ correspond to the given λ . When the main COMBSS is called for a grid of m values of λ with `SubsetMapV2`, then for each $k \leq q$ we obtain at most m models and among them the model with the minimum mean squared error is selected as the final best model for k . Since this version of COMBSS explores the surface, we can refine results further by starting from different initial points $\mathbf{w}^{(0)}$.

Remark 3 It is not hard to observe that for each λ , if the model obtained by Algorithm 2 is of a size $k \leq q$, then this model is present in \mathcal{M}_k of Algorithm 3, and hence, COMBSS with `SubsetMapV2` always provides the same or a better solution than COMBSS with `SubsetMapV1`.

5 Roles of Tuning Parameters

In this section, we provide insights on how the tuning parameters δ and λ influence the objective function $f_\lambda(\mathbf{t})$ (or, equivalently $g_\lambda(\mathbf{w})$) and hence the convergence of the algorithm.

5.1 Controlling the Shape of $f_\lambda(\mathbf{t})$ through δ

The normalized cost $\|\mathbf{y} - X_t \tilde{\boldsymbol{\beta}}_t\|_2^2/n$ provides an estimator of the error variance. For any fixed \mathbf{t} , we expect this variance (and hence the objective function $f_\lambda(\mathbf{t})$) to be almost the same for all relatively large values of n , particularly, in situations where the errors ϵ_i are independent and identically distributed. This is the case at all the corner points $\mathbf{s} \in \{0, 1\}^p$, because at these corner points, from Theorem 2, $X_{\mathbf{s}} \tilde{\boldsymbol{\beta}}_{\mathbf{s}} = X_{[s]} \hat{\boldsymbol{\beta}}_{[s]}$, which is independent of δ . We would like to have a similar behavior at all the interior points $\mathbf{t} \in (0, 1)^p$ as well, so that for each \mathbf{t} , the function $f_\lambda(\mathbf{t})$ is roughly the same for all large values of n . Such consistent behavior is helpful in guaranteeing that the convergence paths of the gradient descent method are approximately the same for large values of n .

Figure 3 shows surface plots of $f_\lambda(\mathbf{t})$ for different values of n and δ for an example dataset

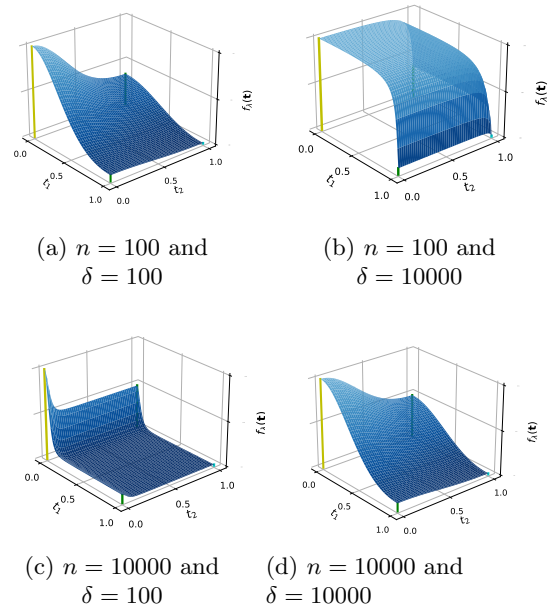


Fig. 3: Illustration of how δ effects the objective function $f_\lambda(\mathbf{t})$ (with $\lambda = 0$). A dataset consists of 10000 samples generated from the illustrative linear model used in Figure 1. For (a) and (b), 100 samples from the same dataset are used.

obtained from a linear model with $p = 2$. Surface plots (a) and (d) correspond to $\delta = n$, and as we can see, the shape of the surface of $f_\lambda(\mathbf{t})$ over $[0, 1]^p$ is very similar in both these plots.

To make this observation more explicit, we now show that the function $f_\lambda(\mathbf{t})$, at any \mathbf{t} , takes almost the same value for all large n if we keep $\delta = cn$, for a fixed constant $c > 0$, under the assumption that the data samples are independent and identically distributed (this assumption simplifies the following discussion; however, the conclusion holds more generally).

Observe that

$$\frac{1}{n} \|\mathbf{y} - X_t \tilde{\boldsymbol{\beta}}_t\|_2^2 = \frac{\mathbf{y}^\top \mathbf{y}}{n} - 2\boldsymbol{\gamma}_t^\top \frac{X^\top \mathbf{y}}{n} + \boldsymbol{\gamma}_t^\top \frac{X^\top X}{n} \boldsymbol{\gamma}_t,$$

where $\boldsymbol{\gamma}_t = n^{-1} T_t L_t^{-1} T_t X^\top \mathbf{y}$. Since, under the independent and identically distributed assumption, $\mathbf{y}^\top \mathbf{y}/n$, $X^\top \mathbf{y}/n$, and $X^\top X/n$ converge element-wise as n increases and T_t is independent of n , we would like to choose δ such that L_t^{-1} also

converges as n increases. Now recall from (6) that

$$L_{\mathbf{t}} = T_{\mathbf{t}} \left(\frac{X^{\top} X}{n} \right) T_{\mathbf{t}} + \frac{\delta}{n} (I - T_{\mathbf{t}}^2).$$

It is then evident that the choice $\delta = cn$ for a fixed constant c , independent of n , makes $L_{\mathbf{t}}$ converging as n increases. Specifically, the choice $c = 1$ justifies the behavior observed in Figure 3.

5.2 Sparsity Controlling through λ

Intuitively, the larger the value of λ the sparser the solution offered by COMBSS using SubsetMapV1, when all other parameters are fixed. We now strengthen this understanding mathematically. From Theorem 4,

$$\nabla f_{\lambda}(\mathbf{t}) = \boldsymbol{\zeta}_{\mathbf{t}} + \lambda \mathbf{1}, \quad \mathbf{t} \in (0, 1)^p,$$

and

$$\nabla g_{\lambda}(\mathbf{w}) = (\boldsymbol{\zeta}_{\mathbf{t}} + \lambda \mathbf{1}) \odot (2\mathbf{w} \odot \exp(-\mathbf{w} \odot \mathbf{w})),$$

for $\mathbf{w} \in \mathbb{R}^p$, where $\boldsymbol{\zeta}_{\mathbf{t}}$, given by (13), is independent of λ . Note the following property of $\boldsymbol{\zeta}_{\mathbf{t}}$.

Proposition 5 *For any $j = 1, \dots, p$, if all t_i for $i \neq j$ are fixed,*

$$\lim_{t_j \downarrow 0} \boldsymbol{\zeta}_{\mathbf{t}}(j) = 0.$$

This result implies that for any $j = 1, \dots, p$, we have $\lim_{t_j \downarrow 0} \partial f_{\lambda}(\mathbf{t}) / \partial t_j = \lambda$, where $\lim_{t_j \downarrow 0}$ denotes the existence of the limit for any sequence of t_j that converges to 0 from the right. Since $\boldsymbol{\zeta}_{\mathbf{t}}$ is independent of λ , the above limit implies that there is a window $(0, a_j)$ such that the slope $\partial f_{\lambda}(\mathbf{t}) / \partial t_j > 0$ for $t_j \in (0, a_j)$ and also the window size increases (i.e., a_j increases) as λ increases. As a result, for the function $g_{\lambda}(\mathbf{w})$, there exists a constant $a'_j > 0$ such that

$$\frac{\partial g_{\lambda}(\mathbf{w})}{\partial w_j} \begin{cases} < 0, & \text{for } -a'_j < w_j < 0 \\ > 0, & \text{for } 0 < w_j < a'_j. \end{cases}$$

In other words, for positive λ , there is a ‘valley’ on the surface of $g_{\lambda}(\mathbf{w})$ along the line $w_j = 0$ and the valley becomes wider as λ increases. In summary, the larger the values of λ the more w_j (or,

equivalently t_j) have tendency to move towards 0 by the optimization algorithm and then a sparse model is selected (i.e, small number k of variables chosen). At the extreme value $\lambda_{\max} = \|y\|_2^2 / n$, all t_j are forced towards 0 and thus the null model will be selected.

6 Efficient Implementation of COMBSS

In this section, we focus on efficient implementation of COMBSS using the *conjugate gradient* method, the *Woodbury matrix identity*, and the *Banachiewicz Inversion Formula*.

6.1 Low- vs High-dimension

Recall the expression of $L_{\mathbf{t}}$ from (5):

$$L_{\mathbf{t}} = \frac{1}{n} \left[X_{\mathbf{t}}^{\top} X_{\mathbf{t}} + \delta (I - T_{\mathbf{t}}^2) \right].$$

We have noticed earlier from Theorem 4 that for computing $\nabla g_{\lambda}(\mathbf{w})$, twice we evaluate matrix-vector products of the form $L_{\mathbf{t}}^{-1} \mathbf{u}$, which is the unique solution of the linear equation $L_{\mathbf{t}} \mathbf{z} = \mathbf{u}$. Solving linear equations efficiently is one of the important and well-studied problems in the field of linear algebra. Among many elegant approaches for solving linear equations, the conjugate gradient method is well-suited for our problem as $L_{\mathbf{t}}$ is symmetric positive-definite; see, for example, Golub and Van Loan (1996).

The running time of the conjugate gradient method for solving the linear equation $A\mathbf{z} = \mathbf{u}$ depends on the dimension of A . For our algorithm, since $L_{\mathbf{t}}$ is of dimension $p \times p$, the conjugate gradient method can return a good approximation of $L_{\mathbf{t}}^{-1} \mathbf{u}$ within $O(p^2)$ time by fixing the maximum number of iterations taken by the conjugate gradient method. This is true for both low-dimensional models (where $p < n$) and high-dimensional models (where $n < p$).

We now specifically focus on high-dimensional models and transform the problem of solving the p -dimensional linear equation $L_{\mathbf{t}} \mathbf{z} = \mathbf{u}$ to the problem of solving an n -dimensional linear equation problem. This approach is based on a well-known result in linear algebra called the Woodbury matrix identity. Since we are calling

the gradient descent method for solving a n -dimensional problem, instead of p -dimensional, we can achieve a much lower overall computational complexity for the high-dimensional models. The following result is a consequence of the Woodbury matrix identity, which is stated as Lemma 2 in Appendix A.

Theorem 6 For $\mathbf{t} \in [0, 1]^p$, let $S_{\mathbf{t}}$ be a p -dimensional diagonal matrix with the j th diagonal element being $n/\delta(1 - t_j^2)$ and $\tilde{L}_{\mathbf{t}} = I + X_{\mathbf{t}}S_{\mathbf{t}}X_{\mathbf{t}}^{\top}/n$. Then,

$$L_{\mathbf{t}}^{-1}\mathbf{u} = (S_{\mathbf{t}}\mathbf{u}) - \frac{1}{n}S_{\mathbf{t}}X_{\mathbf{t}}^{\top}\tilde{L}_{\mathbf{t}}^{-1}(X_{\mathbf{t}}S_{\mathbf{t}}\mathbf{u}).$$

The above expression suggests that instead of solving the p -dimensional problem $L_{\mathbf{t}}\mathbf{z} = \mathbf{u}$ directly, we can first solve the n -dimensional problem $\tilde{L}_{\mathbf{t}}\mathbf{z} = (X_{\mathbf{t}}S_{\mathbf{t}}\mathbf{u})$ and substitute the result in the above expression to get the value of $L_{\mathbf{t}}^{-1}\mathbf{u}$.

6.2 A Dimension Reduction Approach

During the execution of the gradient descent algorithm, Step 1 of Algorithm 1, some of w_j (and hence the corresponding t_j) can reach zero. Particularly, for basic gradient descent and similar methods, once w_j reaches zero it remains zero until the algorithm terminates, because the update of \mathbf{w} in the l th iteration of the basic gradient descent depends only on the gradient $g_{\lambda}(\mathbf{w}^{(l)})$, whose j th element

$$\frac{\partial g_{\lambda}(\mathbf{w}^{(l)})}{\partial w_j} = 0 \quad \text{if } w_j^{(l)} = 0. \quad (14)$$

Because (14) holds, we need to only focus on $\partial g_{\lambda}(\mathbf{w})/\partial w_j$ associated with $w_j \neq 0$ in order to reduce the cost of computing the gradient $\nabla g_{\lambda}(\mathbf{w})$. To simplify the notation, let $\mathcal{P} = \{1, \dots, p\}$ and for any $\mathbf{t} \in [0, 1]^p$, let $\mathcal{Z}_{\mathbf{t}}$ be the set of indices of the zero elements of \mathbf{t} , that is,

$$\mathcal{Z}_{\mathbf{t}} = \{j : t_j = 0, j = \mathcal{P}\}. \quad (15)$$

Similar to the notation used in Theorem 2, for a vector $\mathbf{u} \in \mathbb{R}^p$, we write $(\mathbf{u})_{+}$ (respectively, $(\mathbf{u})_0$) to denote the vector of dimension $p - |\mathcal{Z}_{\mathbf{t}}|$ (respectively, $|\mathcal{Z}_{\mathbf{t}}|$) constructed from \mathbf{u} by removing all its elements with the indices in $\mathcal{Z}_{\mathbf{t}}$ (respectively,

in $\mathcal{P} \setminus \mathcal{Z}_{\mathbf{t}}$). Similarly, for a matrix A of dimension $p \times p$, we write $(A)_{+}$ (respectively, $(A)_0$) to denote the new matrix constructed from A by removing its rows and columns with the indices in $\mathcal{Z}_{\mathbf{t}}$ (respectively, in $\mathcal{P} \setminus \mathcal{Z}_{\mathbf{t}}$). Then we have the following result.

Theorem 7 Suppose $\mathbf{t} \in [0, 1]^p$. Then,

$$(L_{\mathbf{t}})_{+} = \frac{1}{n} \left[(T_{\mathbf{t}})_{+} \left(X^{\top} X \right)_{+} (T_{\mathbf{t}})_{+} + \delta (I - (T_{\mathbf{t}})_{+}) \right].$$

Furthermore, we have

$$\begin{aligned} (L_{\mathbf{t}}^{-1})_0 &= \frac{n}{\delta} I, \\ (L_{\mathbf{t}}^{-1})_{+} &= ((L_{\mathbf{t}})_{+})^{-1}, \end{aligned} \quad (16)$$

$$\begin{aligned} (\tilde{\beta}_{\mathbf{t}})_0 &= \mathbf{0}, \\ (\tilde{\beta}_{\mathbf{t}})_{+} &= ((L_{\mathbf{t}})_{+})^{-1} \left((\mathbf{t})_{+} \odot \left(\frac{X^{\top} \mathbf{y}}{n} \right)_{+} \right), \end{aligned} \quad (17)$$

$$\begin{aligned} (\mathbf{c}_{\mathbf{t}})_0 &= \mathbf{0}, \\ (\mathbf{c}_{\mathbf{t}})_{+} &= ((L_{\mathbf{t}})_{+})^{-1} \left((\mathbf{t})_{+} \odot (\mathbf{a}_{\mathbf{t}})_{+} \right). \end{aligned} \quad (18)$$

In Theorem 7, (16) shows that for every $j \in \mathcal{Z}_{\mathbf{t}}$, all the off-diagonal elements of the j th row as well as the j th column of $L_{\mathbf{t}}^{-1}$ are zero while its j th diagonal element is n/δ , and all other elements of $L_{\mathbf{t}}^{-1}$ (which constitute the sub-matrix $(L_{\mathbf{t}}^{-1})_{+}$) depend only on $(L_{\mathbf{t}})_{+}$, which can be computed using only the columns of the design matrix X with indices in $\mathcal{P} \setminus \mathcal{Z}_{\mathbf{t}}$. As a consequence, (17) and (18) imply that computing $\tilde{\beta}_{\mathbf{t}}$ and $\mathbf{c}_{\mathbf{t}}$ is equal to solving p_{+} -dimensional linear equations of the form $(L_{\mathbf{t}}^{-1})_{+} \mathbf{z} = \mathbf{v}$, where $p_{+} = p - |\mathcal{Z}_{\mathbf{t}}|$. Since $p_{+} \leq p$, solving such a p_{+} -dimensional linear equation using the conjugate gradient can be faster than solving the original p -dimensional linear equation of the form $L_{\mathbf{t}}\mathbf{z} = \mathbf{u}$.

In summary, for a vector $\mathbf{t} \in [0, 1]^p$ with some elements being 0, the values of $f_{\lambda}(\mathbf{t})$ and $\nabla f_{\lambda}(\mathbf{t})$ do not depend on the columns j of X where $t_j = 0$. Therefore, we can reduce the computational complexity by removing all the columns j of the design matrix X where $t_j = 0$.

6.3 Making Our Algorithm Fast

In Section 6.2, we noted that when some elements t_j of \mathbf{t} are zero, it is faster to compute the objective functions $f_\lambda(\mathbf{t})$ and $g_\lambda(\mathbf{t})$ and their gradients $\nabla f_\lambda(\mathbf{t})$ and $\nabla g_\lambda(\mathbf{t})$ by ignoring the columns j of the design matrix X . In Section 5.2, using Proposition 5, we further noted that for any $\lambda > 0$ there is a ‘valley’ on the surface of $g_\lambda(\mathbf{w})$ along $w_j = 0$ for all $j = 1, \dots, p$, and thus for any j , when w_j (or, equivalently, t_j) is sufficiently small during the execution of the gradient descent method, it will eventually become zero. Using these observations, in the implementation of our method, to reduce the computational cost of estimating the gradients, it is wise to map w_j (and t_j) to 0 when w_j is almost zero. We incorporate this truncation idea into our algorithm as follows.

We first fix a small constant η , say at 0.001. As we run the gradient descent algorithm, when t_j becomes smaller than η for some $j \in \mathcal{P}$, we take t_j and w_j to be zero and we stop updating them; that is, t_j and w_j will continue to be zero until the gradient descent algorithm terminates. In each iteration of the gradient descent algorithm, the design matrix is updated by removing all the columns corresponding to zero t_j 's. If the algorithm starts at \mathbf{w} with all non-zero elements, the effective dimension p_+ , which denotes the number of columns in the updated design matrix, monotonically decreases starting from p . In an iteration, if $p_+ > n$, we can use Theorem 6 to reduce the complexity of computing the gradients. However, when p_+ falls below n , we directly use conjugate gradient for computing the gradients without invoking Theorem 6.

Using a dataset, Fig. 4 illustrates the substantial improvement in the speed of our algorithm when the above mentioned improvement ideas are incorporated in its implementation.

7 Simulation Experiments

Our method is available through Python and R codes provided at <https://github.com/saratmoka/COMBSS-Python-VIGNETTE> and <https://github.com/benoit-liquet/COMBSS-R-VIGNETTE>, respectively. The code includes examples where p is as large as of order 10,000. This code further allows to replicate our

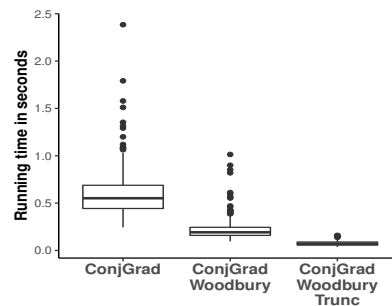
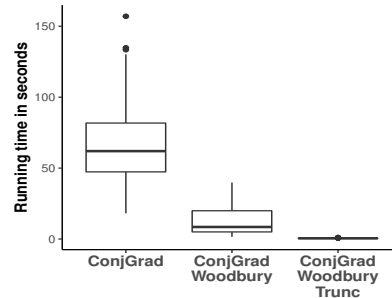
(a) $p = 1000$ (b) $p = 5000$

Fig. 4: Running time of our algorithm at $\lambda = 0.1$ for an example dataset using the adam optimizer, a popular gradient based method. These boxplots are based on 300 replications. Here we compare running times for COMBSS with SubsetMapV1 using only conjugate gradient (ConjGrad), conjugate gradient with Woodbury matrix identity (ConjGrad-Woodbury), and conjugate gradient with both Woodbury matrix identity and truncation improvement (ConjGrad-Woodbury-Trunc). For the truncation, $\eta = 0.001$. The dataset is generated using the Beta-type 1 model presented in Section 7 of (Moka et al, 2022) with $k_0 = 5$, $n = 100$, $\rho = 0.8$, and the signal-to-noise ratio of 5.

empirical research in this article and the simulation studies presented in Moka et al (2022) where efficacy of our method is compared to three popular existing methods, namely, forward selection (FS), Lasso, and mixed integer optimization (MIO). In Moka et al (2022), we focused on demonstrating (using SubsetMapV1) the efficacy in predicting the true model of the data.

Here, our focus is on demonstrating the efficacy of our method in retrieving best subsets of given sizes, meaning our ability to solve (1) using `SubsetMapV2`. We also compare our approach to forward selection, Lasso, mixed integer optimization and `L0Learn` (Hazimeh and Mazumder, 2020b).

7.1 Simulation design

The data is generated from the linear model:

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \text{where } \boldsymbol{\epsilon} \sim N_n(0, \sigma^2 \mathbb{I}_n), \quad (19)$$

where each row of the predictor matrix X is generated from a multivariate normal distribution with zero mean and covariance matrix Σ with diagonal elements $\Sigma_{j,j} = 1$ and off-diagonal elements $\Sigma_{i,j} = \rho^{|i-j|}$, $i \neq j$, for some correlation parameter $\rho \in (-1, 1)$. In order to investigate a challenging situation, we use $\rho = 0.8$ to mimic strong correlation between predictors. For each simulation, we fix the signal-to-noise ratio (SNR) and compute the variance σ^2 of the noise $\boldsymbol{\epsilon}$ using

$$\sigma^2 = \frac{\boldsymbol{\beta}^\top \Sigma \boldsymbol{\beta}}{\text{SNR}}.$$

We consider the following two simulation settings:

- **Beta-type 2:** The first $k_0 = 10$ components of $\boldsymbol{\beta}$ are equal to 1 and all other components of $\boldsymbol{\beta}$ are equal to 0.
- **Beta-type 3:** The first $k_0 = 10$ components of $\boldsymbol{\beta}$ are equal to $\beta_i = 0.5^{i-1}$, for $i = 1, \dots, k_0$ and all other components of $\boldsymbol{\beta}$ are equal to 0.

Both Beta-type 2 and Beta-type 3 assumes strong correlation between the active predictors. Beta-type 3 differs from Beta-type 2 by presenting a signal decaying exponentially to 0. Beta-type 2 has been used in Moka et al (2022) and so to avoid confusion we call the decay signal as Beta-type 3.

For both these types, we investigate the performance of our method in low- and high-dimensional settings. For the low-dimensional setting, we take $n = 100$ and $p = 20$ for $\text{SNR} \in \{0.5, 1, 2, \dots, 8\}$, while for the high-dimensional setting, $n = 100$ and $p = 1000$ for $\text{SNR} \in \{2, 3, \dots, 8\}$.

In the low-dimensional setting, FS and MIO were tuned over $k = 0, \dots, 20$. In this simulation

we ran MIO through the R package `bestsubset` offered in Hastie et al (2018) while we ran `L0Learn` through the R package `L0Learn` offered in Hazimeh et al (2023). For the high dimensional setting, we do not include MIO due to time computational constraints posed by MIO. In particular, the MIO based on the Gurobi optimizer is quite time consuming for high dimensional cases (see Hastie et al (2020)).

In low- and high-dimensional settings, the Lasso was tuned for 50 values of λ ranging from $\lambda_{\max} = \|X^T \mathbf{y}\|_\infty$ to a small fraction of λ_{\max} on a log scale, as per the default in `bestsubset` package. In both the low- and high-dimensional settings, COMBSS with `SubsetMapV2` was called four times starting at four different initial points $\mathbf{t}^{(0)}$: $(0.5, \dots, 0.5)^\top$, $(0.99, \dots, 0.99)^\top$, $(0.75, \dots, 0.75)^\top$, and $(0.3, \dots, 0.3)^\top$. For each call we used 24 values of λ on a grid as follows. Starting from $\lambda_{\max} = \|\mathbf{y}\|_2^2/n$, half of λ values were generated by $\{\lambda_l = \lambda_{\max}/2^{(l)}, l = 1, \dots, 12\}$. From this sequence, the remaining λ values were created by $\{(\lambda_{l+1} + \lambda_l)/2 : l = 1, \dots, 12\}$.

7.2 Low-dimensional case

In low dimensional case, we use the exhaustive method to find the exact solution of the best subset for any subset size ranging from 1 to p . Then, we assess our method in retrieving the exact best subset for each subset size. Figure 5, shows the frequency (over 50 replications) of retrieving the exact best subset (provided by exhaustive search) for any subset size from $k = 1, \dots, p$. For each SNR level, MIO as expected retrieves perfectly the *true* best subset of any model size. Then COMBSS gives the best results to retrieve the best subset compared to FS, Lasso and `L0Learn`. Similar behaviours are reported for Beta-type 3 in Figure 6.

7.3 High-dimensional case

Note that the exact best subset is unknown for the high dimensional case since it is computationally impractical to conduct an exhaustive search over all the subsets of sizes 1 to $p = 1000$. Hence, to assess the performance of our method in retrieving a *competing* best subset, we compare the best subset obtained from COMBSS for two different subset sizes: 5 and 10. For this comparison, we use

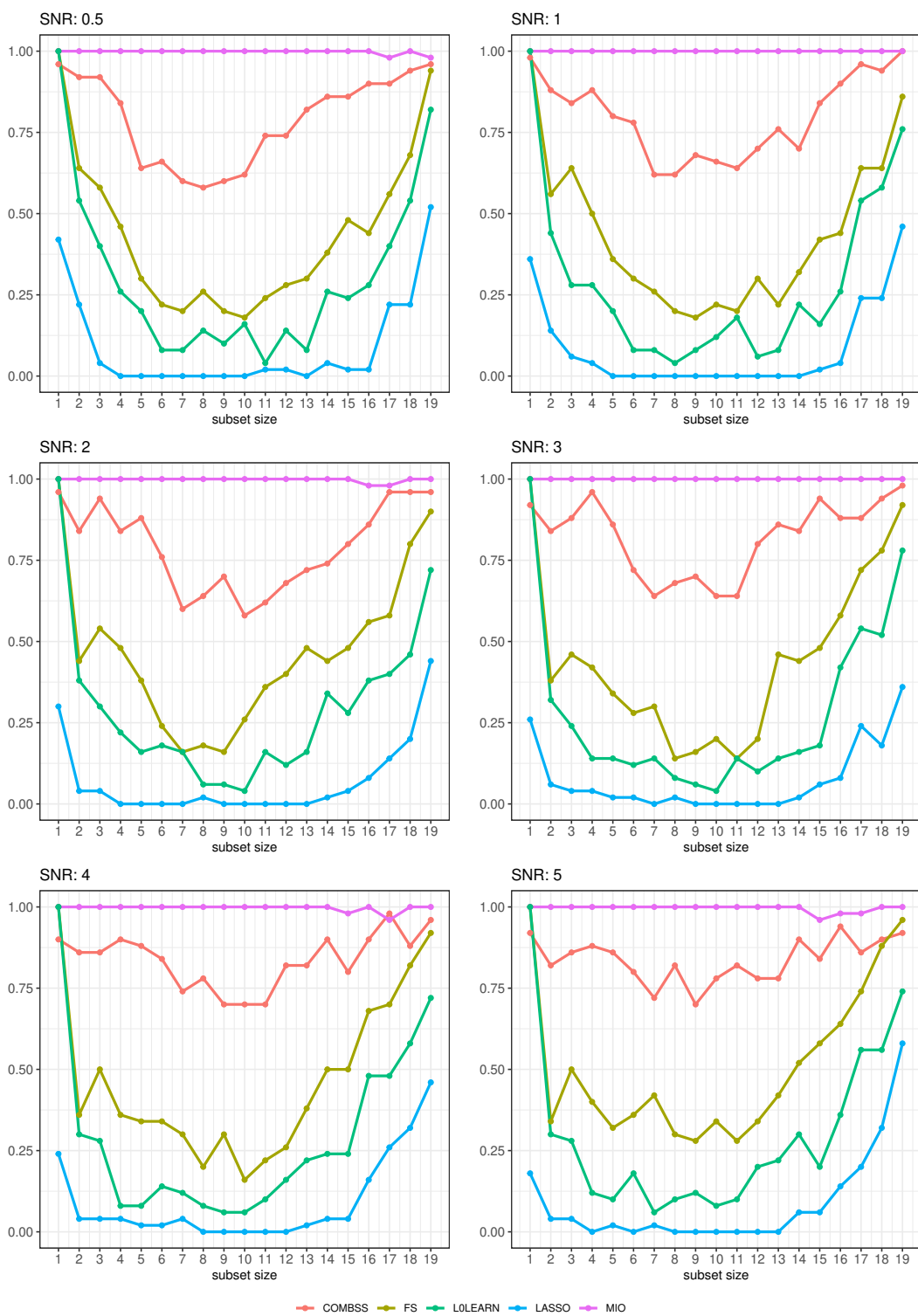


Fig. 5: Frequency (over 50 replications) of retrieving the exact best subset for any subset size from $k = 1, \dots, p$ for beta-type 2 case.



Fig. 6: Frequency (over 50 replications) of retrieving the exact best subset for any subset size from $k = 1, \dots, p$ for beta-type 3 case.

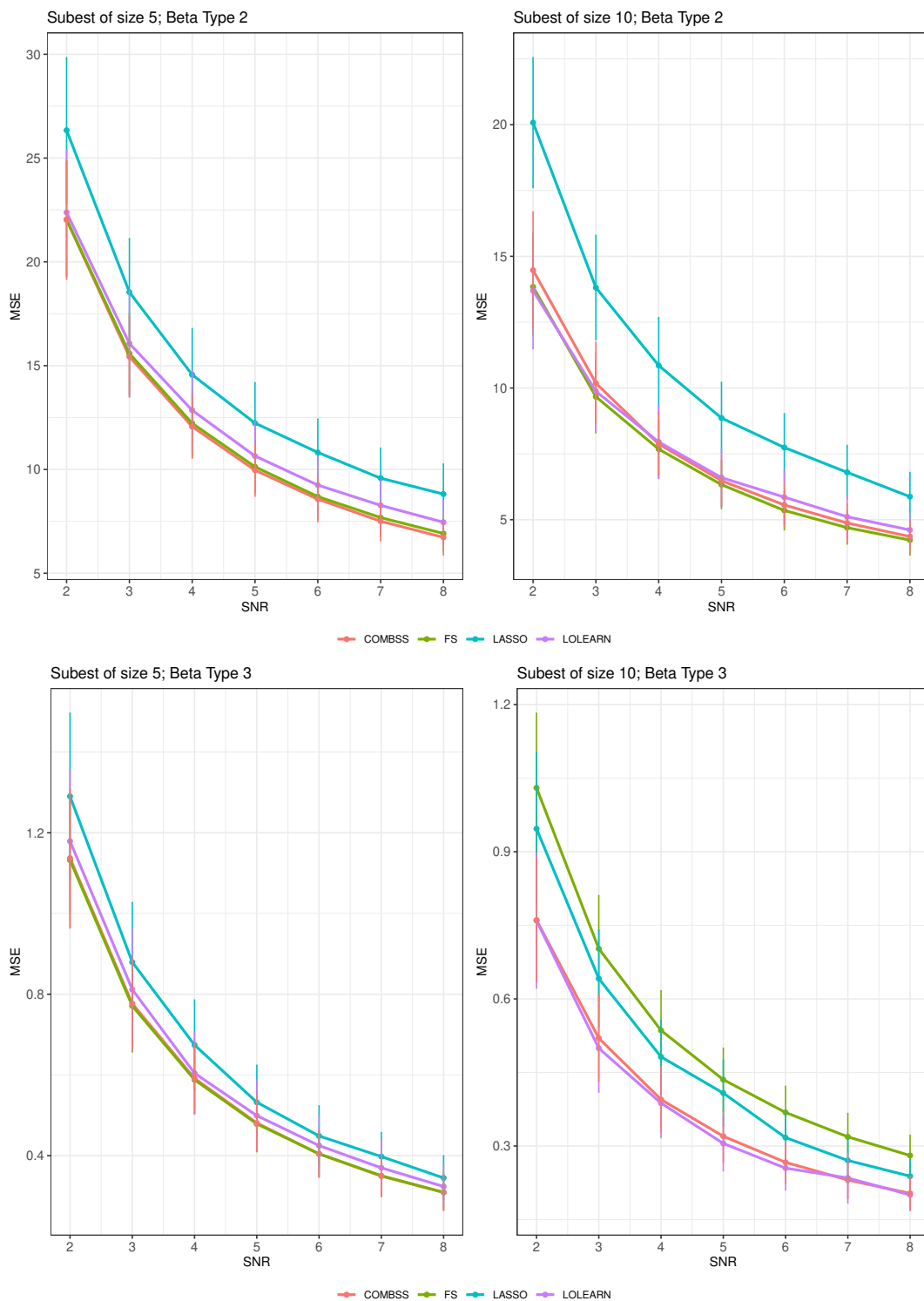


Fig. 7: Ability of COMBSS for providing a *competing* best subset for subset size 5 and 10 comparing to FS, Lasso and LOLearn. Top plots are for beta-type 2 and bottom plots correspond to beta-type 3.

the mean square of error (MSE) of the dataset to evaluate which method is providing a better subset for size 5 and 10. Figure 7 presents these results over 50 replications for SNR values from 2 to 8. Overall, COMBSS is consistently same or better than other methods for providing a *competing* best subset. On the other hand none of the alternative methods is consistent across all the cases.

8 Conclusion and Discussion

In this paper, we have introduced COMBSS, a novel continuous optimization method towards best subset selection in linear regression. The key goal of COMBSS is to extend the highly difficult discrete constrained best subset selection problem to an unconstrained continuous optimization problem. In particular, COMBSS involves extending the objective function of the best subset selection, which is defined at the corners of the hypercube $[0, 1]^p$, to a differentiable function defined on the whole hypercube. For this extended function, starting from an interior point, a gradient descent method is executed to find a corner of the hypercube where the objective function is minimum.

In our algorithm, the primary operations involved are the matrix-vector product, the vector-vector element-wise product, and the scalar-vector product. Additionally, we note that the primary operation for executing a conjugate gradient method for solving a linear equation $Az = \mathbf{u}$ is the matrix-vector product $A\mathbf{u}$. All these operations are known to execute faster on graphics processing unit based computers using parallel programming, which could substantially increase the speed of our method further.

In the preprint Moka et al (2022), we have conducted several simulation experiments in both low-dimensional and high-dimensional setups to illustrate the good performance of COMBSS with SubsetMapV1 for predicting the true model of the data in comparison to the existing popular methods: Forward Stepwise, Lasso, and Mixed Integer Optimization. In this paper, our simulation experiments highlight the ability of COMBSS with SubsetMapV2 for retrieving “exact” best subset for any subset size in comparison to the existing methods. Both of these empirical results highlight the potential of COMBSS for feature extractions.

Our simulation experiments show that the sequence of vectors $\mathbf{t}^{(1)}, \mathbf{t}^{(2)}, \dots$ obtained during the execution of gradient descent always converges. A future work focusing on the theoretical study for establishing such convergence would be useful.

A future direction for finding the best model of a given fixed size k is to explore different options for the penalty term of the objective function $f_\lambda(\mathbf{t})$. Ideally, if we select a sufficiently large penalty for $\sum_{j=1}^p t_j > k$ and 0 otherwise, we can drive the optimization algorithm towards a model of size k that lies along the hyperplane given by $\sum_{j=1}^p t_j = k$. Because such a discrete penalty is not differentiable, we could use smooth alternatives. For instance, the penalty could be taken to be $\lambda(k - \sum_{j=1}^p t_j)^2$ when $\sum_{j=1}^p t_j > k$ and 0 otherwise, for a tuning parameter $\lambda > 0$.

We expect, similarly to the significant body of work that focuses on the Lasso and on MIO, respectively, that there are many avenues that can be explored and investigated for building on the presented COMBSS framework. Particularly, to tackle best subset selection when problems are ultra high dimensional. In this paper, we have opened a novel framework for feature selection and this framework can be extended to other models beyond the linear regression model. For instance, recently Mathur et al (2023) extended the COMBSS framework for solving column subset selection and Nyström approximation problems. Furthermore, our ongoing research focuses on the extensions of COMBSS to non-linear regression problems including logistic regression.

Acknowledgements. Samuel Muller was supported by the Australian Research Council Discovery Project Grant #210100521.

Appendix A Proofs

Proof of Theorem 1 Since both $X_{\mathbf{t}}^\top X_{\mathbf{t}}$ and $T_{\mathbf{t}}$ are symmetric, the symmetry of $L_{\mathbf{t}}$ is obvious. We now show that $L_{\mathbf{t}}$ is positive-definite for $\mathbf{t} \in [0, 1]^p$ by establishing

$$\mathbf{u}^\top L_{\mathbf{t}} \mathbf{u} > 0, \quad \text{for all } \mathbf{u} \in \mathbb{R}^p \setminus \{\mathbf{0}\}. \quad (\text{A1})$$

The matrix $X_{\mathbf{t}}^\top X_{\mathbf{t}}$ is a positive semi-definite, because

$$\mathbf{u}^\top X_{\mathbf{t}}^\top X_{\mathbf{t}} \mathbf{u} = \|X_{\mathbf{t}} \mathbf{u}\|_2^2 \geq 0.$$

In addition, for all $\mathbf{t} \in [0, 1]^p$, the matrix $\delta (I - T_{\mathbf{t}}^2)$ is also a positive-definite because $\delta > 0$, and

$$\begin{aligned} \mathbf{u}^\top (I - T_{\mathbf{t}}^2) \mathbf{u} &= \|\mathbf{u}\|_2^2 - \|T_{\mathbf{t}} \mathbf{u}\|_2^2 \\ &= \sum_{j=1}^p u_j^2 (1 - t_j^2), \end{aligned} \quad (\text{A2})$$

which is strictly positive if $\mathbf{t} \in [0, 1]^p$ and $\mathbf{u} \in \mathbb{R}^p \setminus \{\mathbf{0}\}$. Since positive-definite matrices are invertible, we have $L_{\mathbf{t}}^\dagger = L_{\mathbf{t}}^{-1}$, and thus, $\tilde{\beta}_{\mathbf{t}} = L_{\mathbf{t}}^{-1} X_{\mathbf{t}}^\top \mathbf{y}/n$. \square

Theorem 8 is a collection of results from the literature that we need in our proofs. Results (i) and (ii) of Theorem 8 are well-known in the literature as Banachiewicz inversion lemma (see, e.g., Tian and Takane (2005)), and (iii) is its generalization to Moore–Penrose inverse (See Corollary 3.5 (c) in Castro-González et al (2015)).

Theorem 8 *Let M be a square block matrix of the form*

$$M = \begin{bmatrix} A & C \\ B & D \end{bmatrix}$$

with A being a square matrix. Let the Schur complement $S = D - BA^\dagger C$. Suppose that D is non-singular. Then following holds.

- (i) *If A is non-singular, then M is non-singular if and only if S is non-singular.*
- (ii) *If both A and S are non-singular, then*

$$\begin{aligned} M^{-1} &= \begin{bmatrix} I & -A^{-1}CS^{-1} \\ 0 & S^{-1} \end{bmatrix} \begin{bmatrix} A^{-1} & 0 \\ -BA^{-1} & I \end{bmatrix} \\ &= \begin{bmatrix} I & -A^{-1}C \\ 0 & I \end{bmatrix} \begin{bmatrix} I & 0 \\ -S^{-1}B & S^{-1} \end{bmatrix} \begin{bmatrix} A^{-1} & 0 \\ 0 & I \end{bmatrix}. \end{aligned} \quad (\text{A3})$$

- (iii) *If A is singular, S is non-singular, $BA^\dagger A = B$ and $AA^\dagger C = C$, then*

$$\begin{aligned} M^\dagger &= \begin{bmatrix} I & -A^\dagger CS^{-1} \\ 0 & S^{-1} \end{bmatrix} \begin{bmatrix} A^\dagger & 0 \\ -BA^\dagger & I \end{bmatrix} \\ &= \begin{bmatrix} I & -A^\dagger C \\ 0 & I \end{bmatrix} \begin{bmatrix} I & 0 \\ -S^{-1}B & S^{-1} \end{bmatrix} \begin{bmatrix} A^\dagger & 0 \\ 0 & I \end{bmatrix}. \end{aligned} \quad (\text{A4})$$

Proof of Theorem 2 The inverse of a matrix after a permutation of rows (respectively, columns) is identical to the matrix obtained by applying the same permutation on columns (respectively, rows) on the

inverse of the matrix. Therefore, without loss of generality, we assume that all the zero-elements of $\mathbf{s} \in \{0, 1\}^p$ appear at the end, in the form:

$$\mathbf{s} = (s_1, \dots, s_m, 0, \dots, 0),$$

where m indicates the number of non-zeros in \mathbf{s} . Recall that $X_{[\mathbf{s}]}$ is the matrix of size $n \times |\mathbf{s}|$ created by keeping only columns j of X for which $s_j = 1$. Thus, $L_{\mathbf{s}}$ is given by,

$$\begin{aligned} L_{\mathbf{s}} &= \frac{1}{n} \left[\begin{pmatrix} X_{[\mathbf{s}]}^\top X_{[\mathbf{s}]} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} + \delta \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & I \end{pmatrix} \right] \\ &= \frac{1}{n} \begin{pmatrix} X_{[\mathbf{s}]}^\top X_{[\mathbf{s}]} & \mathbf{0} \\ \mathbf{0} & \delta I \end{pmatrix}. \end{aligned} \quad (\text{A5})$$

From Theorem 8 (i), it is evident that $L_{\mathbf{s}}$ is invertible if and only if $X_{[\mathbf{s}]}^\top X_{[\mathbf{s}]}$ is invertible.

First assume that $X_{[\mathbf{s}]}^\top X_{[\mathbf{s}]}$ is invertible. Then, from Theorem 8 (ii),

$$\begin{aligned} L_{\mathbf{s}}^{-1} &= n \begin{pmatrix} I & \mathbf{0} \\ \mathbf{0} & \frac{1}{\delta} I \end{pmatrix} \begin{pmatrix} (X_{[\mathbf{s}]}^\top X_{[\mathbf{s}]})^{-1} & \mathbf{0} \\ \mathbf{0} & I \end{pmatrix} \\ &= n \begin{pmatrix} (X_{[\mathbf{s}]}^\top X_{[\mathbf{s}]})^{-1} & \mathbf{0} \\ \mathbf{0} & \frac{1}{\delta} I \end{pmatrix}. \end{aligned} \quad (\text{A6})$$

Now recall the notations $(\tilde{\beta}_{\mathbf{s}})_+$ and $(\tilde{\beta}_{\mathbf{s}})_0$ introduced before stating Theorem 2. Then, we use (A6) to obtain

$$\begin{aligned} \begin{pmatrix} (\tilde{\beta}_{\mathbf{s}})_+ \\ (\tilde{\beta}_{\mathbf{s}})_0 \end{pmatrix} &= L_{\mathbf{s}}^{-1} T_{\mathbf{s}} \frac{X_{[\mathbf{s}]}^\top \mathbf{y}}{n} \\ &= \begin{pmatrix} (X_{[\mathbf{s}]}^\top X_{[\mathbf{s}]})^{-1} X_{[\mathbf{s}]}^\top \mathbf{y} \\ \mathbf{0} \end{pmatrix}. \end{aligned}$$

This further guarantees that $X_{[\mathbf{s}]} \hat{\beta}_{[\mathbf{s}]} = X_{\mathbf{s}} \tilde{\beta}_{\mathbf{s}}$.

When $X_{[\mathbf{s}]}^\top X_{[\mathbf{s}]}$ is singular, by replacing the inverse with its pseudo-inverse in the above discussion, and using Theorem 8 (iii) instead of Theorem 8 (ii), we can establish the same conclusions. This is because, the corresponding Schur complement for $L_{\mathbf{s}}$ is $S = nI/\delta$, which is symmetric and positive definite. \square

Proof of Theorem 3 Consider a sequence $\mathbf{t}_1, \mathbf{t}_2, \dots \in [0, 1]^p$ that converges a point $\mathbf{t} \in [0, 1]^p$. We know that the converges easily holds when $\mathbf{t} \in [0, 1]^p$ from the continuity of matrix inversion which states that for any sequence of invertible matrices Z_1, Z_2, \dots that converging to an invertible matrix Z , the sequence of their inverses $Z_1^{-1}, Z_2^{-1}, \dots$ converges to Z^{-1} .

Now using Theorem 8, we prove the convergence when some or all of the elements of \mathbf{t} are equal to 1. Suppose $m\mathbf{t}$ has exactly m elements equal to 1. Using the arguments from the proof of 2, without of loss of

generality assume that all 1s in \mathbf{t} appear together in the first m positions, that is,

$$\mathbf{t} = (\underbrace{1, \dots, 1}_m, \underbrace{t_{m+1}, \dots, t_p}_{p-m \text{ times}}).$$

In that case, by writing

$$T_{\ell,1} = \text{Diag}(t_{\ell,1}, \dots, t_{\ell,m}), \quad \text{and} \\ T_{\ell,2} = \text{Diag}(t_{\ell,m+1}, \dots, t_{\ell,p}),$$

we observe that as $\ell \rightarrow \infty$

$$T_{\ell,1} \rightarrow I, \quad \text{and} \quad T_{\ell,2} \rightarrow T_2 = \text{Diag}(t_{m+1}, \dots, t_p).$$

Further, take

$$F_\ell = \begin{bmatrix} T_{\ell,1} & 0 \\ 0 & I \end{bmatrix},$$

and

$$X = [X_1, X_2],$$

with X_1 denoting the first m columns of X . Similarly, we can write

$$X_{\mathbf{t}_\ell} = [X_{\mathbf{t}_\ell,1}, X_{\mathbf{t}_\ell,2}].$$

We now observe that

$$\begin{aligned} X_{\mathbf{t}_\ell}^\top X_{\mathbf{t}_\ell} &= \begin{bmatrix} X_{\mathbf{t}_\ell,1}^\top X_{\mathbf{t}_\ell,1} & X_{\mathbf{t}_\ell,1}^\top X_{\mathbf{t}_\ell,2} \\ X_{\mathbf{t}_\ell,2}^\top X_{\mathbf{t}_\ell,1} & X_{\mathbf{t}_\ell,2}^\top X_{\mathbf{t}_\ell,2} \end{bmatrix} \\ &= \begin{bmatrix} T_{\ell,1} X_1^\top X_1 T_{\ell,1} & T_{\ell,1} X_1^\top X_2 \\ X_2^\top X_1 T_{\ell,1} & X_2^\top X_2 \end{bmatrix} \\ &= F_\ell \begin{bmatrix} X_1^\top X_1 & X_1^\top X_2 \\ X_2^\top X_1 & X_2^\top X_2 \end{bmatrix} F_\ell. \end{aligned}$$

As a result,

$$\begin{aligned} L_{\mathbf{t}_\ell} &= X_{\mathbf{t}_\ell}^\top X_{\mathbf{t}_\ell} + \delta(I - T_{\mathbf{t}_\ell}^2) \\ &= F_\ell \left(\begin{bmatrix} X_1^\top X_1 & X_1^\top X_2 \\ X_2^\top X_1 & X_2^\top X_2 \end{bmatrix} \right. \\ &\quad \left. + \delta \begin{bmatrix} T_{\mathbf{t}_\ell,1}^{-2} - I & 0 \\ 0 & I - T_{\mathbf{t}_\ell,2}^2 \end{bmatrix} \right) F_\ell. \end{aligned}$$

Now define,

$$\begin{aligned} A_\ell &= X_1^\top X_1 + \delta(T_{\mathbf{t}_\ell,1}^{-2} - I), \\ B_\ell &= X_2^\top X_1, \\ C_\ell &= X_1^\top X_2, \\ D_\ell &= X_2^\top X_2 + \delta(I - T_{\mathbf{t}_\ell,2}^2), \end{aligned}$$

and

$$M_\ell = \begin{bmatrix} A_\ell & C_\ell \\ B_\ell & D_\ell \end{bmatrix}.$$

Since $\mathbf{t}_\ell \in [0, 1]^p$, $L_{\mathbf{t}_\ell}$ is non-singular (see Theorem 1), and hence we have

$$L_{\mathbf{t}_\ell}^{-1} = F_\ell^{-1} M_\ell^{-1} F_\ell^{-1}.$$

Note that the corresponding Schur complement $S_\ell = D_\ell - B_\ell A_\ell^{-1} C_\ell$ is non-singular from Theorem 8 (i). Furthermore, since

$$X_{\mathbf{t}_\ell} = [X_1 \ X_{\mathbf{t}_\ell,2}] F_\ell,$$

$$L_{\mathbf{t}_\ell}^{-1} X_{\mathbf{t}_\ell}^\top = F_\ell^{-1} M_\ell^{-1} \begin{bmatrix} X_1^\top \\ X_{\mathbf{t}_\ell,2}^\top \end{bmatrix},$$

and hence,

$$\begin{aligned} \lim_{\ell \rightarrow \infty} L_{\mathbf{t}_\ell}^{-1} X_{\mathbf{t}_\ell}^\top &= \lim_{\ell \rightarrow \infty} F_\ell \lim_{\ell \rightarrow \infty} \left(M_\ell^{-1} \begin{bmatrix} X_1^\top \\ X_{\mathbf{t}_\ell,2}^\top \end{bmatrix} \right) \\ &= \lim_{\ell \rightarrow \infty} \left(M_\ell^{-1} \begin{bmatrix} X_1^\top \\ X_{\mathbf{t}_\ell,2}^\top \end{bmatrix} \right). \end{aligned}$$

Using (A3),

$$M_\ell^{-1} = \begin{bmatrix} I & -A_\ell^{-1} C_\ell \\ 0 & I \end{bmatrix} \begin{bmatrix} I & 0 \\ -S_\ell^{-1} B_\ell & S_\ell^{-1} \end{bmatrix} \begin{bmatrix} A_\ell^{-1} & 0 \\ 0 & I \end{bmatrix},$$

and hence,

$$M_\ell^{-1} \begin{bmatrix} X_1^\top \\ X_{\mathbf{t}_\ell,2}^\top \end{bmatrix}$$

is equal to

$$\begin{bmatrix} I & -A_\ell^{-1} X_1^\top X_{\mathbf{t}_\ell,2} \\ 0 & I \end{bmatrix} \begin{bmatrix} I & 0 \\ -S_\ell^{-1} B_\ell & S_\ell^{-1} \end{bmatrix} \begin{bmatrix} A_\ell^{-1} X_1^\top & 0 \\ 0 & X_{\mathbf{t}_\ell,2}^\top \end{bmatrix}. \quad (\text{A7})$$

Now by defining

$$\begin{aligned} A &= X_1^\top X_1, \\ B &= X_2^\top X_1, \\ C &= X_1^\top X_2, \quad \text{and} \\ D &= X_2^\top X_2 + \delta(I - T_{\mathbf{t}_\ell,2}^2), \end{aligned}$$

we have

$$L_{\mathbf{t}_\ell} = \begin{bmatrix} A & C \\ B & D \end{bmatrix}.$$

Since $T_{\mathbf{t}_\ell,2} < I$, we can see that D is symmetric positive definite and hence non-singular (this can be established just like the proof of Theorem 1). Furthermore, the corresponding Schur complement $S = D - BA^\dagger C$ is symmetric positive definite, and hence non-singular. The symmetry of S is easy to see from the definition because A and D are symmetric and $B = C^\top$. To see that S is positive definite, for any $\mathbf{x} \in \mathbb{R}^{p-m} \setminus \{\mathbf{0}\}$, let $\mathbf{z} = X_{\mathbf{t}_\ell,2} \mathbf{x}$ and thus

$$\begin{aligned} \mathbf{x}^\top S \mathbf{x} &= \mathbf{z}^\top \mathbf{z} + \delta \mathbf{x}^\top (I - T_{\mathbf{t}_\ell,2}^2) \mathbf{x} - \mathbf{z}^\top X_1 (X_1^\top X_1)^\dagger X_1^\top \mathbf{z} \\ &> \mathbf{z}^\top \mathbf{z} - \mathbf{z}^\top X_1 (X_1^\top X_1)^\dagger X_1^\top \mathbf{z} \end{aligned}$$

$$= z^\top \left(I - X_1 (X_1^\top X_1)^\dagger X_1^\top \right) z.$$

Since $(I - X_1 (X_1^\top X_1)^\dagger X_1^\top)$ is a projection matrix and hence positive definite, S is also positive definite.

In addition, using the singular value decomposition (SVD) $X_1 = U_1 \Delta_1 V_1^\top$, we have

$$\begin{aligned} BA^\dagger A &= X_{t,2}^\top X_1 (X_1^\top X_1)^\dagger (X_1^\top X_1) \\ &= X_{t,2}^\top U_1 \Delta_1 (\Delta_1^\top \Delta_1)^\dagger (\Delta_1^\top \Delta_1) V_1 \\ &= X_{t,2}^\top U_1 \Delta_1 V_1 \\ &= X_{t,2}^\top X_1 = B. \end{aligned}$$

Similarly, we can show that $AA^\dagger C = C$. Thus, using (A4), $L_t^\dagger X_t^\top$ is equal to

$$\begin{bmatrix} I & -A^\dagger X_1^\top X_{t,2} \\ 0 & I \end{bmatrix} \begin{bmatrix} I & 0 \\ -S^{-1}B & S^{-1} \end{bmatrix} \begin{bmatrix} A^\dagger X_1^\top & 0 \\ 0 & X_{t,2} \end{bmatrix}. \quad (\text{A8})$$

Since $\lim_{\ell \rightarrow \infty} X_{t,\ell,2} = X_{t,2}$ and $\lim_{\ell \rightarrow \infty} B_\ell = B$, from (A7) and (A8), to show that

$$\lim_{\ell \rightarrow \infty} L_{t,\ell}^{-1} X_{t,\ell}^\top = L_t^\dagger X_t^\top,$$

it is enough to show that

$$\lim_{\ell \rightarrow \infty} S_\ell^{-1} = S^{-1}, \quad (\text{A9})$$

$$\lim_{\ell \rightarrow \infty} A_\ell^{-1} X_1^\top = A^\dagger X_1^\top. \quad (\text{A10})$$

Since S and each of S_ℓ are non-singular, (A9) holds from the continuity of matrix inversion. Now observe that

$$\begin{aligned} A^\dagger X_1^\top &= \left(X_1^\top X_1 \right)^\dagger X_1^\top \\ &= V_1 \left(\Delta_1^\top \Delta_1 \right)^\dagger \Delta_1^\top U_1^\top \\ &= V_1 \Delta_1^\dagger U_1^\top \\ &= X_1^\dagger, \end{aligned}$$

To establish (A10), we need to show that $\bar{X} = \lim_{\ell \rightarrow \infty} A_\ell^{-1} X_1^\top$ is equal to X_1^\dagger . Towards this, define

$$\begin{aligned} \eta_\ell &= \max_{i=1,\dots,m} (1/t_{\ell,i}^2 - 1), \\ \epsilon_\ell &= \min_{i=1,\dots,m} (1/t_{\ell,i}^2 - 1). \end{aligned}$$

Then, we observe that both η_ℓ and ϵ_ℓ are strictly positive and going to zero as $\ell \rightarrow \infty$. Thus,

$$X_1^\top X_1 + \delta \epsilon_\ell I \leq A_\ell \leq X_1^\top X_1 + \delta \eta_\ell I,$$

where for any two symmetric positive semi-definite matrices Z and Z' , we write $Z \geq Z'$ if $Z - Z'$ is also positive semi-definite. Let

$$\underline{A}_\ell = X_1^\top X_1 + \delta \epsilon_\ell I \quad \text{and} \quad \bar{A}_\ell = X_1^\top X_1 + \delta \eta_\ell I.$$

Thus, $\underline{A}_\ell^{-1} \geq A_\ell^{-1} \geq \bar{A}_\ell^{-1}$, or, alternatively,

$$A_\ell^{-1} - \bar{A}_\ell^{-1} \leq \underline{A}_\ell^{-1} - \bar{A}_\ell^{-1}.$$

Now for any matrix norm, denoting as $\|\cdot\|$, using the triangular inequality,

$$\begin{aligned} \|A^{-1} X_1^\top - X_1^\dagger\| &= \|(A_\ell^{-1} - \bar{A}_\ell^{-1}) X_1^\top + (\bar{A}_\ell^{-1} X_1^\top - X_1^\dagger)\| \\ &\leq \|(A_\ell^{-1} - \bar{A}_\ell^{-1}) X_1^\top\| + \|(\bar{A}_\ell^{-1} X_1^\top - X_1^\dagger)\| \\ &\leq \|(\underline{A}_\ell^{-1} - \bar{A}_\ell^{-1}) X_1^\top\| + \|(\bar{A}_\ell^{-1} X_1^\top - X_1^\dagger)\|. \end{aligned} \quad (\text{A11})$$

Using the SVD of $X_1 = U_1 \Delta_1 V_1^\top$, we get the SVD of $(\underline{A}_\ell^{-1} - \bar{A}_\ell^{-1}) X_1^\top$ as

$$V_1 \left((\Delta_1^\top \Delta_1 + \delta \epsilon_\ell I)^{-1} \Delta_1^\top - (\Delta_1^\top \Delta_1 + \delta \eta_\ell I)^{-1} \Delta_1^\top \right) U_1^\top.$$

That is, suppose σ_i is the i th singular value of X_1 , then the i th singular value of $(\underline{A}_\ell^{-1} - \bar{A}_\ell^{-1}) X_1^\top$ is 0 if $\sigma_i = 0$, otherwise, it is

$$\frac{\sigma_i}{\sigma_i^2 + \delta \epsilon_\ell} - \frac{\sigma_i}{\sigma_i^2 + \delta \eta_\ell} = \frac{\sigma_i \delta (\eta_\ell - \epsilon_\ell)}{(\sigma_i^2 + \delta \epsilon_\ell)(\sigma_i^2 + \delta \eta_\ell)},$$

which goes to zero and thus the first term in (A11) goes to zero. The second term in (A11) also converges to zero because of the limit definition of pseudo-inverse that states that for any matrix Z

$$Z^\dagger = \lim_{\epsilon > 0} \left(Z^\top Z + \epsilon I \right)^{-1} Z^\top.$$

This completes the proof. \square

For proving Theorem 4, we use Lemma 1, which obtains the partial derivatives of $\tilde{\beta}_t$ with respect to the elements of t .

Lemma 1 For any $t \in (0, 1)^p$, the partial derivative $\frac{\partial \tilde{\beta}_t}{\partial t_j}$ for each $j = 1, \dots, p$ is equal to

$$L_t^{-1} \left[E_j - E_j Z T_t L_t^{-1} T_t - T_t Z E_j L_t^{-1} T_t \right] \left(\frac{X_t^\top y}{n} \right),$$

where $Z = n^{-1} (X_t^\top X_t - \delta I)$ and E_j is a square matrix of dimension $p \times p$ with 1 at the (j, j) th position and 0 everywhere else.

Proof of Lemma 1 Existence of $\tilde{\beta}_t$ for every $t \in (0, 1)^p$ and $\delta > 0$ follows from Theorem 1 which states that L_t is positive-definite and hence guarantees the invertibility of L_t . Since $\tilde{\beta}_t = L_t^{-1} T_t X_t^\top y/n$, using matrix calculus, for any $j = 1, \dots, p$,

$$\frac{\partial \tilde{\beta}_t}{\partial t_j} = \frac{\partial (L_t^{-1} T_t)}{\partial t_j} \left(\frac{X_t^\top y}{n} \right)$$

$$\begin{aligned}
&= \left[\frac{\partial L_t^{-1}}{\partial t_j} T_t + L_t^{-1} \frac{\partial T_t}{\partial t_j} \right] \left(\frac{X^\top \mathbf{y}}{n} \right) \\
&= \left[L_t^{-1} \frac{\partial T_t}{\partial t_j} - L_t^{-1} \frac{\partial L_t}{\partial t_j} L_t^{-1} T_t \right] \left(\frac{X^\top \mathbf{y}}{n} \right),
\end{aligned}$$

where we used differentiation of an invertible matrix which implies

$$\frac{\partial L_t^{-1}}{\partial t_j} = -L_t^{-1} \frac{\partial L_t}{\partial t_j} L_t^{-1}.$$

Since $\partial T_t / \partial t_j = E_j$, and the fact that $L_t = T_t Z T_t + \delta I / n$, we get

$$\begin{aligned}
\frac{\partial L_t}{\partial t_j} &= \frac{\partial T_t}{\partial t_j} Z T_t + T_t Z \frac{\partial T_t}{\partial t_j} \\
&= E_j Z T_t + T_t Z E_j.
\end{aligned}$$

Therefore, $L_t^{-1} \frac{\partial T_t}{\partial t_j} - L_t^{-1} \frac{\partial L_t}{\partial t_j} L_t^{-1} T_t$ is equal to

$$\begin{aligned}
&L_t^{-1} E_j - L_t^{-1} E_j Z T_t L_t^{-1} T_t - L_t^{-1} T_t Z E_j L_t^{-1} T_t \\
&= L_t^{-1} \left[E_j - E_j Z T_t L_t^{-1} T_t - T_t Z E_j L_t^{-1} T_t \right].
\end{aligned}$$

This completes the proof Lemma 1. \square

Proof of Theorem 4 To obtain the gradient $\nabla f_\lambda(\mathbf{t})$ for $\mathbf{t} \in (0, 1)^p$, let $\boldsymbol{\gamma}_t = T_t \tilde{\boldsymbol{\beta}}_t = \mathbf{t} \odot \tilde{\boldsymbol{\beta}}_t$. Then,

$$\begin{aligned}
\|\mathbf{y} - X_t \tilde{\boldsymbol{\beta}}_t\|_2^2 &= \|\mathbf{y} - X \boldsymbol{\gamma}_t\|_2^2 \\
&= \mathbf{y}^\top \mathbf{y} - 2 \boldsymbol{\gamma}_t^\top (X^\top \mathbf{y}) + \boldsymbol{\gamma}_t^\top (X^\top X) \boldsymbol{\gamma}_t.
\end{aligned} \tag{A12}$$

Consequently,

$$\begin{aligned}
\frac{\partial f_\lambda(\mathbf{t})}{\partial t_j} &= \frac{1}{n} \frac{\partial}{\partial t_j} \left[\|\mathbf{y} - X_t \tilde{\boldsymbol{\beta}}_t\|_2^2 \right] + \lambda \\
&= -\frac{2}{n} \left(\frac{\partial \boldsymbol{\gamma}_t}{\partial t_j} \right)^\top (X^\top \mathbf{y}) \\
&\quad + \frac{2}{n} \left(\frac{\partial \boldsymbol{\gamma}_t}{\partial t_j} \right)^\top (X^\top X) \boldsymbol{\gamma}_t + \lambda \\
&= \frac{2}{n} \left(\frac{\partial \boldsymbol{\gamma}_t}{\partial t_j} \right)^\top \left[(X^\top X) \boldsymbol{\gamma}_t - (X^\top \mathbf{y}) \right] + \lambda \\
&= 2 \left(\frac{\partial \boldsymbol{\gamma}_t}{\partial t_j} \right)^\top \mathbf{a}_t + \lambda,
\end{aligned} \tag{A13}$$

where $\mathbf{a}_t = n^{-1} [X^\top X \boldsymbol{\gamma}_t - X^\top \mathbf{y}]$. From the definitions of $\tilde{\boldsymbol{\beta}}_t$ and $\boldsymbol{\gamma}_t$,

$$\begin{aligned}
\frac{\partial \boldsymbol{\gamma}_t}{\partial t_j} &= \frac{\partial T_t \tilde{\boldsymbol{\beta}}_t}{\partial t_j} \\
&= \frac{\partial T_t}{\partial t_j} \tilde{\boldsymbol{\beta}}_t + T_t \frac{\partial \tilde{\boldsymbol{\beta}}_t}{\partial t_j}, \\
&= E_j \tilde{\boldsymbol{\beta}}_t + T_t L_t^{-1} \left[E_j - E_j Z T_t L_t^{-1} T_t \right]
\end{aligned}$$

$$-T_t Z E_j L_t^{-1} T_t \left(\frac{X^\top \mathbf{y}}{n} \right),$$

which is obtained using Lemma 1 and the fact that $\partial T_t / \partial t_j = E_j$ and $Z = n^{-1} (X^\top X - \delta I)$. This in-turn yields that $\frac{\partial \boldsymbol{\gamma}_t}{\partial t_j}$ is equal to

$$\begin{aligned}
&E_j \tilde{\boldsymbol{\beta}}_t + T_t L_t^{-1} \left[E_j \left(\frac{X^\top \mathbf{y}}{n} \right) - E_j Z \boldsymbol{\gamma}_t - T_t Z E_j \tilde{\boldsymbol{\beta}}_t \right] \\
&= E_j \tilde{\boldsymbol{\beta}}_t - T_t L_t^{-1} E_j \mathbf{b}_t - T_t L_t^{-1} T_t Z E_j \tilde{\boldsymbol{\beta}}_t,
\end{aligned} \tag{A14}$$

where we recall that

$$\mathbf{b}_t = Z \boldsymbol{\gamma}_t - \left(\frac{X^\top \mathbf{y}}{n} \right) = \mathbf{a}_t - \frac{\delta}{n} \boldsymbol{\gamma}_t.$$

For a further simplification, recall that $\mathbf{c}_t = L_t^{-1} (\mathbf{t} \odot \mathbf{a}_t)$ and $\mathbf{d}_t = Z (\mathbf{t} \odot \mathbf{c}_t)$. Then, from (A14), the matrix $\partial \boldsymbol{\gamma}_t / \partial \mathbf{t}$ of dimension $p \times p$, with j th column being $\partial \boldsymbol{\gamma}_t / \partial t_j$, can be expressed as

$$\begin{aligned}
\frac{\partial \boldsymbol{\gamma}_t}{\partial \mathbf{t}} &= \text{Diag}(\tilde{\boldsymbol{\beta}}_t) - T_t L_t^{-1} \text{Diag}(\mathbf{b}_t) \\
&\quad - T_t L_t^{-1} T_t Z \text{Diag}(\tilde{\boldsymbol{\beta}}_t).
\end{aligned} \tag{A15}$$

From (A13), with $\mathbf{1}$ representing a vector of all ones, $\nabla f_\lambda(\mathbf{t})$ can be expressed as

$$\begin{aligned}
\nabla f_\lambda(\mathbf{t}) &= 2 \text{Diag}(\tilde{\boldsymbol{\beta}}_t) \mathbf{a}_t - 2 \text{Diag}(\mathbf{b}_t) L_t^{-1} T_t \mathbf{a}_t \\
&\quad - 2 \text{Diag}(\tilde{\boldsymbol{\beta}}_t) Z T_t L_t^{-1} T_t \mathbf{a}_t + \lambda \mathbf{1} \\
&= 2 (\tilde{\boldsymbol{\beta}}_t \odot \mathbf{a}_t) - 2 \text{Diag}(\mathbf{b}_t) \mathbf{c}_t \\
&\quad - 2 \text{Diag}(\tilde{\boldsymbol{\beta}}_t) Z T_t \mathbf{c}_t + \lambda \mathbf{1} \\
&= 2 (\tilde{\boldsymbol{\beta}}_t \odot \mathbf{a}_t) - 2 (\mathbf{b}_t \odot \mathbf{c}_t) - 2 (\tilde{\boldsymbol{\beta}}_t \odot \mathbf{d}_t) + \lambda \mathbf{1} \\
&= 2 (\tilde{\boldsymbol{\beta}}_t \odot (\mathbf{a}_t - \mathbf{d}_t)) - 2 (\mathbf{b}_t \odot \mathbf{c}_t) + \lambda \mathbf{1} \\
&= \boldsymbol{\zeta}_t + \lambda \mathbf{1},
\end{aligned}$$

where

$$\boldsymbol{\zeta}_t = 2 (\tilde{\boldsymbol{\beta}}_t \odot (\mathbf{a}_t - \mathbf{d}_t)) - 2 (\mathbf{b}_t \odot \mathbf{c}_t).$$

Finally, recall that $g_\lambda(\mathbf{w}) = f_\lambda(\mathbf{t}(\mathbf{w}))$, $\mathbf{w} \in \mathbb{R}^p$, where the map $\mathbf{t}(\mathbf{w}) = \mathbf{1} - \exp(-\mathbf{w} \odot \mathbf{w})$ and

$$f_\lambda(\mathbf{t}) = \frac{1}{n} \|\mathbf{y} - X_t \tilde{\boldsymbol{\beta}}_t\|_2^2 + \lambda \sum_{j=1}^p t_j.$$

Then, from the chain rule of differentiation, for each $j = 1, \dots, p$,

$$\frac{\partial g_\lambda(\mathbf{w})}{\partial w_j} = \frac{\partial f_\lambda(\mathbf{t})}{\partial t_j} (2w_j \exp(-w_j^2)).$$

Alternatively, in short,

$$\nabla g_\lambda(\mathbf{w}) = \nabla f_\lambda(\mathbf{t}) \odot (2\mathbf{w} \odot \exp(-\mathbf{w} \odot \mathbf{w})). \tag{A16}$$

\square

Proof of Proposition 5 From Theorem 4, we obtain $\nabla f_\lambda(\mathbf{t})$ as follows,

$$\nabla f_\lambda(\mathbf{t}) = \zeta_{\mathbf{t}} + \lambda \mathbf{1},$$

where $\zeta_{\mathbf{t}} \in \mathbb{R}^p$, recalling from Theorem 4, is given by

$$\zeta_{\mathbf{t}} = 2 \left(\tilde{\beta}_{\mathbf{t}} \odot (\mathbf{a}_{\mathbf{t}} - \mathbf{d}_{\mathbf{t}}) \right) - 2 (\mathbf{b}_{\mathbf{t}} \odot \mathbf{c}_{\mathbf{t}}), \quad (\text{A17})$$

with $\mathbf{a}_{\mathbf{t}} = n^{-1}[X^\top X(\mathbf{t} \odot \tilde{\beta}_{\mathbf{t}}) - X^\top \mathbf{y}]$, $\mathbf{b}_{\mathbf{t}} = \mathbf{a}_{\mathbf{t}} - n^{-1}\delta(\mathbf{t} \odot \tilde{\beta}_{\mathbf{t}})$, $\mathbf{c}_{\mathbf{t}} = L_{\mathbf{t}}^{-1}(\mathbf{t} \odot \mathbf{a}_{\mathbf{t}})$, and $\mathbf{d}_{\mathbf{t}} = n^{-1}[X^\top X - \delta I](\mathbf{t} \odot \mathbf{c}_{\mathbf{t}})$.

From the definition of $\tilde{\beta}_{\mathbf{t}}$, we can show that for any j , if we fix t_i for all $i \neq j$, then the j th component of $\tilde{\beta}_{\mathbf{t}}$ goes to zero as $t_j \downarrow 0$, that is, $\lim_{t_j \downarrow 0} \tilde{\beta}_{\mathbf{t}}(j) = 0$. Similarly, $\lim_{t_j \downarrow 0} \mathbf{c}_{\mathbf{t}}(j) = 0$. Therefore, from the expression of $\zeta_{\mathbf{t}}$ in (A17), we have $\lim_{t_j \downarrow 0} \zeta_{\mathbf{t}}(j) = 0$. \square

For proving Theorem 6, we use Lemma 2 which is well-known as the Woodbury matrix identity or Duncan Inversion Formula; we refer to Woodbury (1950) for a proof of Lemma 2.

Lemma 2 For any conformable matrices A, B_1 , and B_2 , and C , the matrix $(A + B_1 C B_2)^{-1}$ is equal to

$$A^{-1} - A^{-1} B_1 (C^{-1} + B_2 A^{-1} B_1)^{-1} B_2 A^{-1}.$$

Proof of Theorem 6 Recall the expression of $L_{\mathbf{t}}$:

$$L_{\mathbf{t}} = \frac{1}{n} \left[X_{\mathbf{t}}^\top X_{\mathbf{t}} + \delta (I - T_{\mathbf{t}}^2) \right].$$

From the definition, for $\mathbf{t} \in [0, 1]^p$,

$$S_{\mathbf{t}} = \frac{n}{\delta} (I - T_{\mathbf{t}}^2)^{-1},$$

which exists. Further, if we take

$$A = \frac{\delta}{n} (I - T_{\mathbf{t}}^2), \quad B_1 = B_2 = \frac{1}{\sqrt{n}} X_{\mathbf{t}}^\top, \quad \text{and} \quad C = I,$$

in Lemma 2, then,

$$L_{\mathbf{t}}^{-1} = S_{\mathbf{t}} - \frac{1}{n} S_{\mathbf{t}} X_{\mathbf{t}}^\top \left(I + \frac{1}{n} X_{\mathbf{t}} S_{\mathbf{t}} X_{\mathbf{t}}^\top \right)^{-1} X_{\mathbf{t}} S_{\mathbf{t}}.$$

Since $S_{\mathbf{t}}$ is a diagonal matrix,

$$\tilde{L}_{\mathbf{t}} = I + \frac{1}{n} X_{\mathbf{t}} S_{\mathbf{t}} X_{\mathbf{t}}^\top$$

is a symmetric positive-definite matrix of dimension $n \times n$. Thus, $\tilde{L}_{\mathbf{t}}^{-1}$ exists and

$$L_{\mathbf{t}}^{-1} = S_{\mathbf{t}} - \frac{1}{n} S_{\mathbf{t}} X_{\mathbf{t}}^\top \tilde{L}_{\mathbf{t}}^{-1} X_{\mathbf{t}} S_{\mathbf{t}}. \quad \square$$

Proof of Theorem 7 For the same reasons mentioned in the proof of Theorem 2, without loss of generality we can assume that all the zero elements of \mathbf{t} appear at the end of \mathbf{t} ; that is, \mathbf{t} is of the form:

$$\mathbf{t} = (t_1, \dots, t_m, 0, \dots, 0),$$

where m indicates the number of non-zero elements in \mathbf{t} . Then, $L_{\mathbf{t}}$ is given by

$$\begin{aligned} L_{\mathbf{t}} &= \frac{1}{n} \left[\begin{pmatrix} (X_{\mathbf{t}}^\top X_{\mathbf{t}})_+ & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} + \delta \begin{pmatrix} I - (T_{\mathbf{t}}^2)_+ & \mathbf{0} \\ \mathbf{0} & I \end{pmatrix} \right] \\ &= \begin{pmatrix} (L_{\mathbf{t}})_+ & \mathbf{0} \\ \mathbf{0} & \frac{\delta}{n} I \end{pmatrix}. \end{aligned} \quad (\text{A18})$$

Since $\mathbf{t} \in [0, 1]^p$, from Theorem 1, $L_{\mathbf{t}}$ is a positive-definite matrix. Every principle submatrix of a positive-definite matrix is also positive-definite. Thus, $(L_{\mathbf{t}})_+$ is positive-definite and hence invertible. Using Theorem 8 (ii),

$$\begin{aligned} L_{\mathbf{t}}^{-1} &= \begin{pmatrix} ((L_{\mathbf{t}})_+)^{-1} & \mathbf{0} \\ \mathbf{0} & I \end{pmatrix} \begin{pmatrix} I & \mathbf{0} \\ \mathbf{0} & \frac{\delta}{n} I \end{pmatrix} \\ &= \begin{pmatrix} ((L_{\mathbf{t}})_+)^{-1} & \mathbf{0} \\ \mathbf{0} & \frac{\delta}{n} I \end{pmatrix}. \end{aligned} \quad (\text{A19})$$

Now observe that

$$\begin{pmatrix} X_{\mathbf{t}}^\top \mathbf{y} \\ n \end{pmatrix} = \begin{pmatrix} \left(\left(\frac{X_{\mathbf{t}}^\top \mathbf{y}}{n} \right)_+ \right) \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} (\mathbf{t})_+ \odot \left(\frac{X_{\mathbf{t}}^\top \mathbf{y}}{n} \right)_+ \\ \mathbf{0} \end{pmatrix}.$$

Since $\tilde{\beta}_{\mathbf{t}} = L_{\mathbf{t}}^{-1} X_{\mathbf{t}}^\top \mathbf{y} / n$, using (A19), we establish the desired expressions for $(\tilde{\beta}_{\mathbf{t}})_0$ and $(\tilde{\beta}_{\mathbf{t}})_+$. Similar arguments yield the desired expressions for $(c_{\mathbf{t}})_0$ and $(c_{\mathbf{t}})_+$; hence for conciseness that proof is omitted here. \square

References

- Bertsimas D, King A, Mazumder R (2016) Best subset selection via a modern optimization lens. *The Annals of Statistics* 44(2):813 – 852
- Castro-González N, Martínez-Serrano M, Robles J (2015) Expressions for the moore–penrose inverse of block matrices involving the schur complement. *Linear Algebra and its Applications* 471:353–368
- Froymson MA (1966) Stepwise regression—a backward and forward look. Presented at the Eastern Regional Meetings of the Institute of Mathematical Statistics, Florham Park, New Jersey

- Fan J, Li R (2006) Statistical challenges with high dimensionality: feature selection in knowledge discovery. In: International Congress of Mathematicians. Vol. III. Eur. Math. Soc., Zürich, p 595–622
- Fan J, Lv J (2010) A selective overview of variable selection in high dimensional feature space. *Statist Sinica* 20(1):101–148
- Furnival GM, Wilson RW (2000) Regressions by leaps and bounds. *Technometrics* 42:69–79
- Golub GH, Van Loan CF (1996) Matrix computations, 3rd edn. Johns Hopkins Studies in the Mathematical Sciences, Johns Hopkins University Press, Baltimore, MD
- Gurobi Optimization, limited liability company (2022) Gurobi Optimizer Reference Manual. URL <https://www.gurobi.com>
- Hastie T, Tibshirani R, Tibshirani R (2018) bestsubset: Tools for best subset selection in regression. R package version 1.0.10
- Hastie T, Tibshirani R, Tibshirani R (2020) Best subset, forward stepwise or lasso? analysis and recommendations based on extensive comparisons. *Statistical Science* 35(4):579–592
- Hazimeh H, Mazumder R (2020a) Fast best subset selection: coordinate descent and local combinatorial optimization algorithms. *Oper Res* 68(5):1517–1537
- Hazimeh H, Mazumder R (2020b) Fast best subset selection: Coordinate descent and local combinatorial optimization algorithms. *Operations Research* 68(5):1517–1537
- Hazimeh H, Mazumder R, Nonet T (2023) L0Learn: Fast Algorithms for Best Subset Selection. R package version 2.1.0
- Hocking RR, Leslie RN (1967) Selection of the best subset in regression analysis. *Technometrics* 9:531–540
- Kochenderfer MJ, Wheeler TA (2019) Algorithms for optimization. MIT Press, Cambridge, MA
- Mathur A, Moka S, Botev Z (2023) Column subset selection and Nyström approximation via continuous optimization. [2304.09678](https://arxiv.org/abs/2304.09678)
- Miller A (2019) Subset selection in regression, Monographs on Statistics and Applied Probability, vol 95. Chapman & Hall/CRC, Boca Raton, FL
- Moka S, Lique B, Zhu H, et al (2022) COMBSS: Best subset selection via continuous optimization. URL <https://arxiv.org/abs/2205.02617>
- Müller S, Welsh AH (2010) On model selection curves. *Intl Statist Reviews* 78(2):240–256
- Natarajan BK (1995) Sparse approximate solutions to linear systems. *SIAM J Comput* 24(2):227–234
- Tarr G, Muller S, Welsh AH (2018) mplot: An r package for graphical model stability and variable selection procedures. *J Statist Software* 83(9):1–28
- Tian Y, Takane Y (2005) Schur complements and Banachiewicz-Schur forms. *Electron J Linear Algebra* 13:405–418
- Tibshirani R (1996) Regression shrinkage and selection via the lasso. *J Roy Statist Soc Ser B* 58(1):267–288
- Woodbury MA (1950) Inverting modified matrices. Princeton University, Princeton, N. J., statistical Research Group, Memo. Rep. no. 42,
- Zhu J, Wen C, Zhu J, et al (2020) A polynomial algorithm for best-subset selection problem. *Proc Natl Acad Sci USA* 117(52):33,117–33,123