

A Global Two-stage Algorithm for Non-convex Penalized High-dimensional Linear Regression Problems

Peili Li*

Min Liu[†]Zhou Yu[‡]

Abstract

By the asymptotic oracle property, non-convex penalties represented by minimax concave penalty (MCP) and smoothly clipped absolute deviation (SCAD) have attracted much attentions in high-dimensional data analysis, and have been widely used in signal processing, image restoration, matrix estimation, etc. However, in view of their non-convex and non-smooth characteristics, they are computationally challenging. Almost all existing algorithms converge locally, and the proper selection of initial values is crucial. Therefore, in actual operation, they often combine a warm-starting technique to meet the rigid requirement that the initial value must be sufficiently close to the optimal solution of the corresponding problem. In this paper, based on the DC (difference of convex functions) property of MCP and SCAD penalties, we aim to design a global two-stage algorithm for the high-dimensional least squares linear regression problems. A key idea for making the proposed algorithm to be efficient is to use the primal dual active set with continuation (PDASC) method, which is equivalent to the semi-smooth Newton (SSN) method, to solve the corresponding sub-problems. Theoretically, we not only prove the global convergence of the proposed algorithm, but also verify that the generated iterative sequence converges to a d -stationary point. In terms of computational performance, the abundant research of simulation and real data show that the algorithm in this paper is superior to the latest SSN method and the classic coordinate descent (CD) algorithm for solving non-convex penalized high-dimensional linear regression problems.

Keywords: High-dimensional linear regression, global convergence, two-stage algorithm, primal dual active set with continuation algorithm, difference of convex functions.

*School of Statistics, KLATASDS-MOE, East China Normal University, Shanghai 200062, PR China (Email: plli@sfs.ecnu.edu.cn).

[†]School of Mathematics and Statistics, Wuhan University, Wuhan 430072, PR China (Email: mliuf@whu.edu.cn).

[‡]School of Statistics, KLATASDS-MOE, East China Normal University, Shanghai 200062, PR China (Email: zyu@stat.ecnu.edu.cn).

1 Introduction

In this paper, we mainly consider the following high-dimensional linear regression model:

$$y = X\beta^* + \varepsilon, \quad (1.1)$$

where $y \in \mathbb{R}^n$ is the response vector, $X \in \mathbb{R}^{n \times p}$ is the design matrix, $\varepsilon \in \mathbb{R}^n$ is the noise vector, and β^* is the underlying regression coefficient. In the high-dimensional settings, the number of predictors p is usually larger or much larger than the number of observations n . At this time, we usually assume that β^* is sparse, that is, only a small part of its elements are non-zero. If this idea is expressed in the parameter estimation models, it is natural to add the constraint $\|\beta\|_0 \leq s$, where $\|\beta\|_0$ denotes the number of non-zero elements in β , and $s > 0$ is a tuning parameter which controls the sparsity level. However, the non-convexity and discontinuity of the ℓ_0 pseudo-norm make it NP-hard to solve the corresponding problems Natarajan (1995). Especially in the high-dimensional settings, it is very challenging to design a feasible algorithm that can achieve accurate solutions. Therefore, various surrogates of the ℓ_0 pseudo-norm have been proposed in the existing literature and have been widely studied in statistics, optimization, computational mathematics, machine learning and other fields.

The first type of surrogate functions is mainly the well-known ℓ_1 norm Chen, Donoho, and Saunders (Chen et al.); Fan et al. (2014); Tibshirani (1996), and its corresponding Lagrangian form of least squares linear regression model is the following convex but non-smooth minimization problem:

$$\min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \|X\beta - y\|^2 + \lambda \|\beta\|_1 \right\}, \quad (1.2)$$

where $\|\beta\|_1 = \sum_{i=1}^p |\beta_i|$ denotes the ℓ_1 norm of the vector β , $\lambda > 0$ is a regularization parameter. In view of the good characteristics of the above model, it has received extensive attention in different application fields. Theoretically, under certain conditions on the design matrix X and the sparsity level of the underlying regression coefficient β^* , the min-

minimizers of (1.2) have attractive statistical properties Candes and Tao (2005); Meinshausen and Buhlmann (2006); Zhao and Yu (2006). Numerically, the convexity of (1.2) has led to many fast and effective algorithms, such as least angle regression (LARS) Efron et al. (2004), alternating direction method of multipliers (ADMM) Boyd et al. (2011), coordinate descent (CD) method Wu and Lange (2008) and semi-smooth Newton (SSN) method (or equivalent primal dual active set (PDAS) algorithm) Hintermüller et al. (2002); Li et al. (2018) etc. It is worth emphasizing that the PDAS algorithm in Fan et al. (2014) not only has the local superlinear convergence which can be obtained by reformulating it in the SSN framework, but also has the locally one step convergence under certain conditions. In addition, the continuation technique on the regularization parameter globalizes the convergence of the algorithm. In this paper, we will apply it to solve internal sub-problems, and one can see Section 3.2 for details.

Although the convexity of ℓ_1 penalty makes the corresponding problem computationally attractive, there still exists bias in its estimator. Therefore, scholars proposed the second type of surrogate functions for ℓ_0 pseudo-norm, which mainly contains some non-convex penalties, such as the minimax concave penalty (MCP) Zhang (2010a), the smoothly clipped absolute deviation (SCAD) penalty Fan and Li (2001), capped ℓ_1 Zhang (2010b) and bridge Frank and Friedman (1993); Fu (1998) etc. Numerous studies have shown that, compared with a convex relaxation with the ℓ_1 norm, a proper non-convex penalty method can achieve a sparse estimation with fewer measurements, and is more robust against noises Chartrand (2007); Chen and Gu (2014). Therefore, non-convex penalties have been widely used in various sparse learning problems Breheny and Huang (2011); Chartrand (2007); Chen and Gu (2014); Gong et al. (2013); Huang et al. (2021); Li et al. (2017); Mazumder et al. (2011).

In this paper, we mainly focus on the least squares regression model with MCP or

SCAD penalty, i.e.,

$$\min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \|X\beta - y\|^2 + \sum_{i=1}^p \rho(\beta_i; \lambda, \tau) \right\}, \quad (1.3)$$

where $\rho(\cdot; \lambda, \tau)$ is the MCP or SCAD penalty, which are respectively defined by

$$\rho_{mcp}(t; \lambda, \tau) := \lambda \int_0^{|t|} \max\left(0, 1 - \frac{s}{\lambda\tau}\right) ds = \begin{cases} \frac{\lambda^2\tau}{2}, & |t| > \lambda\tau, \\ \lambda\left(|t| - \frac{t^2}{2\lambda\tau}\right), & |t| \leq \lambda\tau, \quad \tau > 1, \end{cases} \quad (1.4)$$

$$\rho_{scad}(t; \lambda, \tau) := \lambda \int_0^{|t|} \min\left(1, \frac{\max(0, \lambda\tau - s)}{\lambda(\tau - 1)}\right) ds = \begin{cases} \frac{\lambda^2(\tau+1)}{2}, & |t| > \lambda\tau, \\ \frac{\lambda\tau|t| - \frac{1}{2}(t^2 + \lambda^2)}{\tau-1}, & \lambda < |t| \leq \lambda\tau, \quad \tau > 2, \\ \lambda|t|, & |t| \leq \lambda, \end{cases} \quad (1.5)$$

Here τ is a given parameter which controls the concavity of the corresponding penalty. When proposing MCP and SCAD penalties, their authors established that the regression models with MCP and SCAD penalties have the so-called oracle property, that is, in an asymptotic sense, they perform as well as if the analyst had known in advance which coefficients were zero and which were nonzero.

However, non-convex and non-smooth characteristics of the objective function make the numerical calculation of model (1.3) very challenging. There are several typical algorithms in the existing literature, and here is a simple summary in chronological order. Firstly, the authors of Fan and Li (2001); Hunter and Li (2005) proposed a local quadratic approximation (LQA) algorithm and its slightly perturbed version. They suggested iteratively, locally approximating the penalty function by a quadratic function, and then using a modified Newton-Raphson algorithm to solve the corresponding problem. However, the behavior of deleting small coefficients or choosing the size of perturbation will cause numerical instability. To overcome this difficulty, Zou and Li Zou and Li (2008) proposed a new unified algorithm based on the local linear approximation (LLA), and calculated the resulting LASSO problem by LARS algorithm. However, LLA used the path-tracing LARS algorithm to update the regression coefficients, so it is inherently inefficient to some extent.

Then, the coordinate descent (CD) type algorithms were designed for the least squares regression models penalized by MCP and SCAD Breheny and Huang (2011); Mazumder et al. (2011). The numerical results showed that the performance of this algorithm is better than that of LLA. However, the CD-type algorithm requires many iterations in the pursuit of high accuracy, because its convergence rate is sub-linear or linear locally Li and Pong (2018). In addition, it has been proved that each non-convex surrogate function of ℓ_0 pseudo-norm can be expressed as the difference of two convex functions Ahn et al. (2017); Le Thi et al. (2015). Therefore, based on the DC (difference of convex functions) property of the non-convex functions, Li et al. Li et al. (2017) proposed a DC proximal Newton (DCPN) method for the general nonlinear problems with non-convex penalty. They firstly used multistage convex relaxation to transform the original optimization into sequences of LASSO regularized nonlinear regressions. Then, in each stage, they used the second order Taylor expansion to approximate the nonlinear loss functions, and adopted the Proximity Newton method in Lee et al. (2014) to solve the convex sub-problem. Under the conditions of locally restricted strong convexity and Hessian smoothness, they proved their algorithm is locally quadratic convergent within each stage of convex relaxation. Recently, Shi et al. Shi et al. (2018) and Huang et al. Huang et al. (2021) respectively proposed SSN and PDAS algorithms for the model (1.3), and their convergence rates are all locally super-linear.

After in-depth study of the relevant literature, we can find that above-mentioned algorithms are all locally convergent, so they generally combine various warm-starting techniques in actual operations. This inspires us to design an effective calculation method with global convergence to weaken the rigid requirement that the initial value must be sufficiently close to the optimal solution. Here we will design a global two-stage algorithm based on the DC expression of MCP and SCAD penalties. From Ahn et al. (2017); Le Thi et al. (2015); Tang et al. (2020), we know that MCP and SCAD penalties can be reformulated

as:

$$\sum_{i=1}^p \rho_{mcp}(\beta_i; \lambda, \tau) = \lambda \|\beta\|_1 - q_{mcp}(\beta), \quad \tau > 1, \quad (1.6)$$

$$\sum_{i=1}^p \rho_{scad}(\beta_i; \lambda, \tau) = \lambda \|\beta\|_1 - q_{scad}(\beta), \quad \tau > 2, \quad (1.7)$$

where $q_{mcp}(\beta) = \sum_{i=1}^p q_{mcp}(\beta_i; \lambda, \tau)$, $q_{scad}(\beta) = \sum_{i=1}^p q_{scad}(\beta_i; \lambda, \tau)$, and

$$q_{mcp}(t; \lambda, \tau) = \begin{cases} \lambda|t| - \frac{\lambda^2\tau}{2}, & |t| > \lambda\tau \\ \frac{t^2}{2\tau}, & |t| \leq \lambda\tau \end{cases}, \quad q_{scad}(t; \lambda, \tau) = \begin{cases} \lambda|t| - \frac{\lambda^2(\tau+1)}{2}, & |t| > \lambda\tau \\ \frac{(|t|-\lambda)^2}{2(\tau-1)}, & \lambda < |t| \leq \lambda\tau \\ 0, & |t| \leq \lambda \end{cases}.$$

The functions $q_{mcp}(\beta)$ and $q_{scad}(\beta)$ are continuously differentiable with

$$\frac{\partial q_{mcp}(\beta)}{\partial \beta_i} = \begin{cases} \lambda \operatorname{sign}(\beta_i), & |\beta_i| > \lambda\tau \\ \frac{\beta_i}{\tau}, & |\beta_i| \leq \lambda\tau \end{cases}, \quad \frac{\partial q_{scad}(\beta)}{\partial \beta_i} = \begin{cases} \lambda \operatorname{sign}(\beta_i), & |\beta_i| > \lambda\tau \\ \frac{\operatorname{sign}(\beta_i)(|\beta_i|-\lambda)}{\tau-1}, & \lambda < |\beta_i| \leq \lambda\tau \\ 0, & |\beta_i| \leq \lambda \end{cases}.$$

Therefore, the original model (1.3) can be rewritten as follows,

$$\min_{\beta \in \mathbb{R}^p} \left\{ f(\beta) := \frac{1}{2} \|X\beta - y\|^2 + \lambda \|\beta\|_1 - q(\beta) \right\}, \quad (1.8)$$

where $q : \mathbb{R}^p \rightarrow \mathbb{R}$ is q_{mcp} or q_{scad} , which is a convex smooth function. Together with the motivation from global and super-linear proximal majorization-minimization (PMM) algorithm in Tang et al. (2020), which is proposed for nonconvex square-root-loss regression problems, we are thus inspired to adopt the PMM framework for solving the least squares model (1.8). A key idea for making the proposed algorithm to be efficient is to use the PDASC algorithm for solving the corresponding sub-problems. Specifically, in the first stage, by directly removing the second term $-q(\beta)$ and adding a proximal term $\frac{\sigma}{2} \|\beta\|^2$, we will use the PDASC method in Fan et al. (2014) to solve the obtained convex sub-problem, which can get an initial point of the second stage. Then in the second stage, we linearize the second term $-q(\beta)$ with respect to the current iteration β^k and add an appropriate proximal term $\frac{\sigma}{2} \|\beta - \beta^k\|^2$, then directly use the PDASC method to iteratively solve the resulting problem.

The remainder of this paper is organized as follows. In Section 2, we present some preliminaries for our subsequent developments. In Section 3, we describe the two-stage algorithm detailly. We establish the algorithm's convergence in Section 4. In Section 5, we report numerical experiments to show the efficiency of the algorithm, and do performance comparisons with the latest SSN method and the classic CD algorithm. Finally, we conclude our paper in Section 6.

2 Preliminaries

We denote the set of all proper lower semicontinuous convex functions on \mathbb{R}^p as $\mathcal{L}(\mathbb{R}^p)$. For a given $f \in \mathcal{L}(\mathbb{R}^p)$, The proximal mapping of f is defined as

$$Prox_f(x) := \arg \min_{y \in \mathbb{R}^p} \left\{ f(y) + \frac{1}{2} \|y - x\|^2 \right\}, \quad \forall x \in \mathbb{R}^p.$$

Then, from Micchelli et al. (2011), we have

$$z \in \partial f(y) \Leftrightarrow y = Prox_f(y + z). \quad (2.1)$$

The proximal operator of $\|\cdot\|_1$ is given by the pointwise soft-thresholding operator Donoho and Johnstone (1995):

$$Prox_{\lambda \|\cdot\|_1}(x) = S_\lambda(x), \quad (2.2)$$

where

$$y = S_\lambda(x) \Leftrightarrow y_i = \max\{|x_i| - \lambda, 0\} \text{sign}(x_i). \quad (2.3)$$

The subdifferential of any $f \in \mathcal{L}(\mathbb{R}^p)$ is a set-value mapping defined by

$$\partial f(x) := \{z \in \mathbb{R}^p : f(y) \geq f(x) + \langle z, y - x \rangle, \forall y \in \mathbb{R}^p\}.$$

The subdifferential of $f = \|\cdot\|_1$ is the pointwise set-value sign function $\text{Sign}(x)$ Donoho and Johnstone (1995), i.e.,

$$z \in \text{Sign}(x) \Leftrightarrow z_i \begin{cases} = 1, & x_i > 0 \\ \in [-1, 1], & x_i = 0 \\ = -1, & x_i < 0 \end{cases} . \quad (2.4)$$

The classical Fermat's rule for proper lower semicontinuous convex functions Rockafellar (2015) asserts

$$\mathbf{0} \in \partial f(x^*) \Leftrightarrow x^* \text{ is a global minimizer of } f, \quad (2.5)$$

where $\mathbf{0}$ denotes a column vector whose elements are all 0. If the function f is locally Lipschitz continuous near x^* and directionally differentiable at x^* , then $0 \in \partial f(x^*)$ is equivalent to the directional-stationarity (d-stationarity) of x^* , that is

$$f'(x^*; h) := \lim_{\delta \rightarrow 0} \frac{f(x^* + \delta h) - f(x^*)}{\delta} \geq 0, \forall h \in \mathbb{R}^p.$$

In this paper, we will prove that the iterative sequence of the proposed algorithm converges to a d-stationarity point of problem (1.8).

3 Algorithm

In this section, we will propose a two-stage proximal majorization-minimization (PMM) algorithm for model (1.3), and the internal sub-problem with ℓ_1 penalty will be approximately solved by the primal dual active set with continuation (PDASC) method in Fan et al. (2014).

3.1 PMM algorithm

The PMM algorithm contains two stages, where the first stage provides a good initial point for the second stage. Another key idea to make PMM algorithm effective is to use the PDASC algorithm for solving the corresponding subproblems. Specifically, in the first stage, we get a nonsmooth convex subproblem with ℓ_1 penalty by directly removing the concave term $-q(\beta)$ and adding a proximal term $\frac{\sigma}{2}\|\beta\|^2$. Then we use PDASC method to approximately solve the obtained subproblem so that the corresponding KKT residual satisfies a prescribed termination criterion. Next, the solution obtained in the first stage is used as the initial value of the second stage. In the second stage, we linearize the concave

term $-q(\beta)$ with respect to the current iteration β^k and add an appropriate proximal term $\frac{\sigma}{2}\|\beta - \beta^k\|^2$. Then we also use PDASC to solve the corresponding convex sub-problem so that the error vector satisfies a preset accuracy condition. In this way, the second stage is looped and the penalty parameter σ is updated iteratively until the iteration sequence satisfies the termination condition given in advance.

Given $\sigma > 0$, $\tilde{\beta} \in \mathbb{R}^p$ and $\tilde{v} \in \mathbb{R}^p$, we consider the following minimization problem in each iteration:

$$\min_{\beta \in \mathbb{R}^p} J(\beta; \sigma, \tilde{\beta}, \tilde{v}) := \frac{1}{2}\|X\beta - y\|^2 + \lambda\|\beta\|_1 - q(\tilde{\beta}) - \langle \tilde{v}, \beta - \tilde{\beta} \rangle + \frac{\sigma}{2}\|\beta - \tilde{\beta}\|^2. \quad (3.1)$$

Obviously, the above model is a convex problem with ℓ_1 penalty, which can be effectively solved by PDASC method. Next, we summarize the iterative framework of PMM algorithm in Algorithm 3.1.

PMM algorithm

Step 1. Take $\sigma^1 > 0$, $\sigma^{2,0} > 0$. Compute

$$\beta^0 = \arg \min_{\beta \in \mathbb{R}^p} \{J(\beta; \sigma^1, \mathbf{0}, \mathbf{0})\} \quad (3.2)$$

by PDASC method such that the corresponding KKT residual satisfies a prescribed termination criterion. For $k = 0, 1, 2, \dots$, do the following operations iteratively.

Step 2. Compute

$$\beta^{k+1} = \arg \min_{\beta \in \mathbb{R}^p} \{J(\beta; \sigma^{2,k}, \beta^k, \nabla q(\beta^k)) + \langle \delta^k, \beta - \beta^k \rangle\}$$

by PDASC method such that the error vector δ^k satisfies

$$\|\delta^k\| \leq \frac{\sigma^{2,k}}{4} \|\beta^{k+1} - \beta^k\|. \quad (3.3)$$

Step 3. Check the prescribed stopping condition, if stop, denote the last iteration by $\hat{\beta}$. Else, update $\sigma^{2,k+1} = \gamma\sigma^{2,k}$ with $\gamma \in (0, 1)$ and set $k := k + 1$.

Remark 3.1. *It should be pointed out that, we do not need to calculate the dual problem of the corresponding subproblem. Because the sub-problem here is essentially a convex problem with ℓ_1 penalty, which can be directly and effectively solved by the PDASC method. This part is different from Tang et al. (2020).*

Remark 3.2. *Through the verification of many experiments and the communication with the authors in Tang et al. (2020), we found that if we use PDASC to solve $\min_{\beta \in \mathbb{R}^p} \{J(\beta; \sigma^{2,k}, \beta^k, \nabla q(\beta^k))\}$ in the second stage so that the corresponding KKT residual satisfies a prescribed accuracy, such as $1e - 6$, then the condition (3.3) is automatically contented. Therefore, in our subsequent numerical experiments, the termination conditions of all sub-problems are set as the corresponding KKT residuals are sufficiently small. And the inequality (3.3) is mainly used for theoretical analysis.*

3.2 The PDASC method for sub-problems

From Fan et al. (2014), we can see that the design idea of PDASC method is inspired by the first order optimality system of (3.1), which can be seen in the following Lemma 3.1.

Lemma 3.1. *$\beta^* \in \mathbb{R}^p$ is a global minimizer of (3.1) if and only if there exists a $d^* \in \mathbb{R}^p$ such that the following KKT system holds:*

$$(X^\top X + \sigma I)\beta^* + d^* = X^\top y + \tilde{v} + \sigma \tilde{\beta}, \quad (3.4)$$

$$\beta^* = S_\lambda(\beta^* + d^*). \quad (3.5)$$

Proof. By (2.5), we can have

$$\beta^* \in \mathbb{R}^p \text{ is a minimizer of (3.1)} \Leftrightarrow \mathbf{0} \in \partial J(\beta^*; \sigma, \tilde{\beta}, \tilde{v}).$$

Obviously, $\partial J(\beta^*; \sigma, \tilde{\beta}, \tilde{v}) = (X^\top X + \sigma I)\beta^* - X^\top y - \tilde{v} - \sigma \tilde{\beta} + \lambda \partial \|\cdot\|_1(\beta^*)$. Therefore, there exists $d^* \in \lambda \partial \|\cdot\|_1(\beta^*)$ such that

$$(X^\top X + \sigma I)\beta^* + d^* = X^\top y + \tilde{v} + \sigma \tilde{\beta}. \quad (3.6)$$

In addition, (2.1) and (2.2) imply

$$d^* \in \lambda \partial \|\cdot\|_1(\beta^*) \Leftrightarrow \beta^* = \text{Prox}_{\lambda \|\cdot\|_1}(\beta^* + d^*) = S_\lambda(\beta^* + d^*).$$

□

Based on Lemma 3.1, we can directly apply the PDASC method to solve problem (3.1), which is exhibited in Algorithm 3.2.

PDASC method with $(\sigma, \tilde{\beta}, \tilde{v}) \in \mathbb{R}_{++} \times \mathbb{R}^p \times \mathbb{R}^p$

Step 0. Given $\lambda_0 \geq \|X^\top y\|_\infty$, the active set $\mathcal{A}(\lambda_0) = \emptyset$, $\beta(\lambda_0) = \mathbf{0}$, $d(\lambda_0) = X^\top y$, $\mu \in (0, 1)$, $K_{max} \in \mathbb{N}$. For $j = 0, 1, \dots$, do the following operations iteratively.

Step 1. Let $\lambda_j = \mu \lambda_{j-1}$, $\mathcal{A}_0 = \mathcal{A}(\lambda_{j-1})$, $(\beta^0, d^0) = (\beta(\lambda_{j-1}), d(\lambda_{j-1}))$. For $k = 1, 2, \dots, K_{max}$, do the following operations iteratively.

Step 1.1. Compute the active and inactive sets \mathcal{A}_k and \mathcal{I}_k :

$$\begin{aligned} \mathcal{A}_k^+ &= \{i \in [p] : \beta_i^{k-1} + d_i^{k-1} > \lambda\}, \\ \mathcal{A}_k^- &= \{i \in [p] : \beta_i^{k-1} + d_i^{k-1} < -\lambda\}, \\ \mathcal{A}_k &= \mathcal{A}_k^+ \cup \mathcal{A}_k^-, \quad \mathcal{I}_k = \mathcal{A}_k^c. \end{aligned} \tag{3.7}$$

Step 1.2. Check stopping criterion $\mathcal{A}_k = \mathcal{A}_{k-1}$.

Step 1.3. Update the primal and dual variables β^k and d^k respectively by

$$\begin{aligned} \beta_{\mathcal{I}_k}^k &= \mathbf{0}_{\mathcal{I}_k}, \quad d_{\mathcal{A}_k}^k = \lambda[\mathbf{1}_{\mathcal{A}_k^+}; -\mathbf{1}_{\mathcal{A}_k^-}], \\ \beta_{\mathcal{A}_k}^k &= (X_{\mathcal{A}_k}^\top X_{\mathcal{A}_k} + \sigma I_{\mathcal{A}_k})^{-1} (X_{\mathcal{A}_k}^\top y + \tilde{v}_{\mathcal{A}_k} + \sigma \tilde{\beta}_{\mathcal{A}_k} - d_{\mathcal{A}_k}^k), \\ d_{\mathcal{I}_k}^k &= X_{\mathcal{I}_k}^\top y + \tilde{v}_{\mathcal{I}_k} + \sigma \tilde{\beta}_{\mathcal{I}_k} - X_{\mathcal{I}_k}^\top X_{\mathcal{A}_k} \beta_{\mathcal{A}_k}^k. \end{aligned} \tag{3.8}$$

Step 2. Set $\tilde{k} = \min(K_{max}, k)$, $\mathcal{A}(\lambda_j) = \{i \in [p] : \beta_i^{\tilde{k}} + d_i^{\tilde{k}} > \lambda\} \cup \{i \in [p] : \beta_i^{\tilde{k}} + d_i^{\tilde{k}} < -\lambda\}$ and $(\beta(\lambda_j), d(\lambda_j)) = (\beta^{\tilde{k}}, d^{\tilde{k}})$.

Step 3. Check stop condition, if stop, employ the high-dimensional Bayesian information criterion (HBIC) to choose the optimal regularization parameter $\hat{\lambda}$ and denote the corresponding $\beta(\hat{\lambda})$ by $\hat{\beta}$. Else, $j := j + 1$.

Remark 3.3. For the step 10 in Algorithm 3.2, the high-dimensional Bayesian information criterion (HBIC) Wang et al. (2013) chooses the optimal $\hat{\lambda}$ by

$$\hat{\lambda} = \arg \min_{\lambda \in [\lambda_{min}, \lambda_{max}]} \left\{ HBIC(\lambda) := \log\left(\frac{1}{n} \|X\beta(\lambda) - y\|^2\right) + \frac{\log(\log(n))\log(p)}{n} \|\beta(\lambda)\|_0 \right\},$$

where λ_{min} and λ_{max} will be specified in numerical tests.

4 Convergence analysis

We firstly describe the convergence result of the algorithm in our first stage. Since $J(\beta; \sigma^1, \mathbf{0}, \mathbf{0})$ is bounded below, we can get the following result from (Hofmann and Hohage, 2011, Proposition 4.19) and (Tang et al., 2020, Theorem 4.2).

Theorem 4.1. Let $\bar{J}(\sigma^1) := \min_{\beta \in \mathbb{R}^p} \{J(\beta; \sigma^1, \mathbf{0}, \mathbf{0})\}$. Then we have

$$\lim_{\sigma^1 \rightarrow 0} \bar{J}(\sigma^1) = \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \|X\beta - y\|^2 + \lambda \|\beta\|_1 \right\}.$$

Proof. For any $\sigma^1 > 0$ and $\beta \in \mathbb{R}^p$, we have

$$\bar{J}(\sigma^1) \leq \frac{1}{2} \|X\beta - y\|^2 + \lambda \|\beta\|_1 + \frac{\sigma^1}{2} \|\beta\|^2.$$

Therefore, $\lim_{\sigma^1 \rightarrow 0} \bar{J}(\sigma^1) \leq \frac{1}{2} \|X\beta - y\|^2 + \lambda \|\beta\|_1$. Combining with the arbitrariness of β , we can get

$$\lim_{\sigma^1 \rightarrow 0} \bar{J}(\sigma^1) \leq \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \|X\beta - y\|^2 + \lambda \|\beta\|_1 \right\}.$$

In addition, since $\frac{\sigma^1}{2} \|\beta\|^2 \geq 0$ for any $\beta \in \mathbb{R}^p$, so

$$\bar{J}(\sigma^1) \geq \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \|X\beta - y\|^2 + \lambda \|\beta\|_1 \right\},$$

and then

$$\lim_{\sigma^1 \rightarrow 0} \bar{J}(\sigma^1) \geq \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \|X\beta - y\|^2 + \lambda \|\beta\|_1 \right\}.$$

Hence, we can get the desired result. \square

Then, we will analyze the convergence of PMM algorithm. Denote

$$J_k(\beta) := J(\beta; \sigma^{2,k}, \beta^k, \nabla q(\beta^k)).$$

At the k -th iteration of stage II, we have that

$$\beta^{k+1} = \arg \min_{\beta \in \mathbb{R}^p} \{ J_k(\beta) + \langle \delta^k, \beta - \beta^k \rangle \} \quad (4.9)$$

such that condition (3.3) is satisfied. The following lemma shows the descent property of the function J_k .

Lemma 4.1. *Let β^{k+1} be an approximate solution of the subproblem in the k -th iteration such that (3.3) holds. Then we have*

$$J_k(\beta^k) \geq J_k(\beta^{k+1}) - \frac{\sigma^{2,k}}{4} \|\beta^{k+1} - \beta^k\|^2.$$

Proof. From the convexity of function J_k , we have $J_k(\beta^k) - J_k(\beta^{k+1}) \geq \langle \partial J_k(\beta^{k+1}), \beta^k - \beta^{k+1} \rangle$. In addition, we can get $-\delta^k \in \partial J_k(\beta^{k+1})$ from (4.9). Therefore, we obtain

$$J_k(\beta^k) - J_k(\beta^{k+1}) \geq \langle \delta^k, \beta^{k+1} - \beta^k \rangle \geq -\|\delta^k\| \cdot \|\beta^{k+1} - \beta^k\|.$$

Combining with condition (3.3), it is easy to get the desired result

$$J_k(\beta^k) \geq J_k(\beta^{k+1}) - \frac{\sigma^{2,k}}{4} \|\beta^{k+1} - \beta^k\|^2.$$

\square

Next we recall the equivalent expression of a d-stationary point of (1.8) in the following lemma, which is similar to that in Cui et al. (2018); Pang et al. (2017); Tang et al. (2020).

Lemma 4.2. *The vector $\bar{\beta} \in \mathbb{R}^p$ is a d-stationary point of (1.8) if and only if there exist $\sigma > 0$ such that*

$$\bar{\beta} \in \arg \min_{\beta \in \mathbb{R}^p} \{J(\beta; \sigma, \bar{\beta}, \nabla q(\bar{\beta}))\}.$$

Proof. The proof is similar to (Tang et al., 2020, Lemma 4.2), so it is omitted here. \square

Now we present the main result of this section on the subsequential convergence of $\{\beta^k\}$ to a d-stationary point of (1.8).

Theorem 4.2. *Assume $\{\sigma^{2,k}\}$ is a convergent sequence. Let $\{\beta^k\}$ be the sequence generated by the PMM algorithm. The following two statements hold.*

1. *The function sequence $\{f(\beta^k)\}$ is convergent, and $\lim_{k \rightarrow \infty} \|\beta^{k+1} - \beta^k\| = 0$.*
2. *Every accumulation point of the sequence $\{\beta^k\}$, if exists, is a d-stationary point of (1.8).*

Proof. 1. We can easily get $f(\beta^k) = J_k(\beta^k)$. Then from Lemma 4.1, we have

$$\begin{aligned} f(\beta^k) &= J_k(\beta^k) \geq J_k(\beta^{k+1}) - \frac{\sigma^{2,k}}{4} \|\beta^{k+1} - \beta^k\|^2 \\ &= \frac{1}{2} \|X\beta^{k+1} - y\|^2 + \lambda \|\beta^{k+1}\|_1 - q(\beta^k) - \langle \nabla q(\beta^k), \beta^{k+1} - \beta^k \rangle + \frac{\sigma^{2,k}}{4} \|\beta^{k+1} - \beta^k\|^2 \\ &= f(\beta^{k+1}) + \frac{\sigma^{2,k}}{4} \|\beta^{k+1} - \beta^k\|^2 + q(\beta^{k+1}) - q(\beta^k) - \langle \nabla q(\beta^k), \beta^{k+1} - \beta^k \rangle \\ &\geq f(\beta^{k+1}) + \frac{\sigma^{2,k}}{4} \|\beta^{k+1} - \beta^k\|^2. \end{aligned}$$

The last inequality is derived from the convexity of q . Therefore the sequence $\{f(\beta^k)\}$ is non-increasing. Since $f(\beta)$ is bounded below, the sequence $\{f(\beta^k)\}$ converges, and then the sequence $\{\|\beta^{k+1} - \beta^k\|\}$ converges to zero.

2. Let β^∞ be the limit of a convergent subsequence $\{\beta^k\}_{k \in K_0}$. We can easily prove that $\{\beta^{k+1}\}_{k \in K_0}$ also converges to β^∞ . From the definition of β^{k+1} , we can get

$$J_k(\beta) + \langle \delta^k, \beta - \beta^k \rangle \geq J_k(\beta^{k+1}) + \langle \delta^k, \beta^{k+1} - \beta^k \rangle, \quad \forall \beta \in \mathbb{R}^p.$$

Further,

$$J_k(\beta) \geq J_k(\beta^{k+1}) + \langle \delta^k, \beta^{k+1} - \beta \rangle \geq J_k(\beta^{k+1}) - \|\delta^k\| \cdot \|\beta^{k+1} - \beta\|, \quad \forall \beta \in \mathbb{R}^p.$$

Letting $k(\in K_0) \rightarrow \infty$, we obtain that $J_\infty(\beta) \geq J_\infty(\beta^\infty), \forall \beta \in \mathbb{R}^p$. Equivalently,

$$J(\beta; \sigma^{2,\infty}, \beta^\infty, \nabla q(\beta^\infty)) \geq J(\beta^\infty; \sigma^{2,\infty}, \beta^\infty, \nabla q(\beta^\infty)), \forall \beta \in \mathbb{R}^p,$$

where $\sigma^{2,\infty} = \lim_{k \rightarrow \infty} \sigma^{2,k} \geq 0$. Then, we can conclude

$$\beta^\infty \in \arg \min_{\beta \in \mathbb{R}^p} \{J(\beta; \sigma^{2,\infty}, \beta^\infty, \nabla q(\beta^\infty))\}.$$

From Lemma 4.2, we can easily obtain the desired result. □

5 Numerical Experiments

In this section, we will use multiple sets of simulated and real examples to illustrate the performance of the proposed PMM algorithm for non-convex penalized high-dimension linear regression problems. The specific layout is that we first use some examples to illustrate the behavior of PMM algorithm, and then highlight the effectiveness and comparability through numerical comparison with the latest SSN method in Shi et al. (2018) and the classic CD algorithm in Breheny and Huang (2011). All the experiments are performed with Microsoft Windows 10 and MATLAB R2019a, and run on a PC with an Intel Core i7-9700 CPU at 3.00 GHz and 16 GB of memory.

5.1 Experiments setting

In the simulation experiments, we generate the $n \times p$ matrix X whose rows are drawn independently from $\mathcal{N}(0, \Sigma)$ with $\Sigma_{ij} = \kappa^{|i-j|}, 1 \leq i, j \leq p$, where κ is the correlation coefficient of matrix X . In order to generate the target regression coefficient $\beta^* \in \mathbb{R}^p$, we randomly select a subset of $\{1, \dots, p\}$ to form the active set \mathcal{A}^* with $|\mathcal{A}^*| = K < n$. Let $R = m_2/m_1$, where $m_2 = \|\beta_{\mathcal{A}^*}^*\|_{max}$ and $m_1 = \|\beta_{\mathcal{A}^*}^*\|_{min} = 1$. Then the K nonzero coefficients in β^* are uniformly distributed in $[m_1, m_2]$. The response variable is generated by $y = X\beta^* + \varepsilon$ where $\varepsilon \in \mathbb{R}^n$ is the additive Gaussian noise and generated independently from $\mathcal{N}(0, \sigma_1^2 I_n)$.

To select the optimal regularization parameter, we set $\lambda_{max} = \|X^\top b\|_\infty$ and $\lambda_{min} = 10^{-10}\lambda_{max}$. Then an equal-distributed partition on log-scale is employed to divide the interval $[\lambda_{min}, \lambda_{max}]$ into $N = 100$ subintervals. For the parameter τ , unless otherwise specified, we set $\tau = 2.7$ and $\tau = 3.7$ for the MCP and SCAD penalties, respectively. Due to the locally one step convergence of the PDAS method for ℓ_1 regularized least squares problems, we set $K_{max} = 1$. And we use the following two relative KKT residuals R_{kkt}^1 and R_{kkt}^2 to measure the accuracy of the approximate optimal solutions in different stages,

$$R_{kkt}^1 := \frac{\left\| \beta - \text{Prox}_{\lambda\|\cdot\|_1} \left(\beta - X^\top (X\beta - y) \right) \right\|}{1 + \|\beta\| + \|X^\top (X\beta - y)\|}, \quad (5.10)$$

$$R_{kkt}^2 := \frac{\left\| \beta - \text{Prox}_{\lambda\|\cdot\|_{1-q(\cdot)}} \left(\beta - X^\top (X\beta - y) \right) \right\|}{1 + \|\beta\| + \|X^\top (X\beta - y)\|}, \quad (5.11)$$

where the closed form of $\text{Prox}_{\lambda\|\cdot\|_{1-q(\cdot)}}$ can refer to Gong et al. (2013). Then the PDASC method for solving the internal subproblems is terminated if $R_{kkt}^1 < 1e - 6$, and the PMM algorithm will be terminated if $R_{kkt}^2 < 1e - 6$. In addition, we fix some low-impact parameters, such as $\sigma^1 = \sigma^{2,0} = \gamma = 0.1$. The values of other parameters will be given in the context of specific issues.

In addition, for the purpose of highlighting the efficiency and accuracy of PMM algorithm in the subsequent simulation comparison, we compare it with the latest SSN method and the classic CD algorithm from the perspective of the following four indicators based on 100 independent experiments:

- The average CPU time (Time, in seconds);
- The average ℓ_2 relative error: $RE := \frac{\sum_{m=1}^{100} \left(\frac{\|\hat{\beta}^{(m)} - \beta^*\|_2}{\|\beta^*\|_2} \right)}{100}$;
- The average estimated model size: $MS := \frac{\sum_{m=1}^{100} |\hat{\mathcal{A}}^{(m)}|}{100}$;
- The proportion of correct models: $CM := \frac{\sum_{m=1}^{100} \delta_{\{\hat{\mathcal{A}}^{(m)} = \mathcal{A}^*\}}}{100}$,

where $\hat{\beta}$ and $\hat{\mathcal{A}}$ are the estimated regression coefficient and active set, respectively. $|\mathcal{A}|$

indicates the length of set \mathcal{A} , and $\delta\{\hat{\mathcal{A}}^{(m)} = \mathcal{A}^*\} = \begin{cases} 1, & \hat{\mathcal{A}}^{(m)} = \mathcal{A}^* \\ 0, & \hat{\mathcal{A}}^{(m)} \neq \mathcal{A}^* \end{cases}$. Clearly, the smaller Time, the faster calculation speed. And the smaller RE, the closer MS approaches to K , the closer CM approaches to 100%, the higher the solution quality.

5.2 The behavior of PMM algorithm

In this part, we analyze the computational behavior of the PMM algorithm based on 100 independent experiments and consider the problem setting with $n = 300$, $p = 1000$, $K = 10$, $\sigma_1 = 0.1$, $R = 100$. Here we only give the results related to the MCP penalty, since SCAD penalty will produce a similar phenomenon.

Firstly, we utilize a box plot to investigate the performance of variable selection and parameter estimation for the PMM algorithm. To achieve the goal, we generate a coefficient matrix X with $\kappa = 0.2$ and a fixed true regression parameter β^* , whose 10 non-zero elements are $\beta_{30}^* = 6$, $\beta_{198}^* = -11$, $\beta_{269}^* = -10$, $\beta_{395}^* = 25$, $\beta_{442}^* = -8$, $\beta_{495}^* = 100$, $\beta_{637}^* = -9$, $\beta_{766}^* = -10$, $\beta_{777}^* = 5$, $\beta_{865}^* = 1$. In view of the large p , we only describe the estimation effect of non-zero elements in β^* on the left side of Figure 1. Obviously, for each non-zero element, the estimated results fluctuate very little in 100 independent experiments, which fully illustrates the effectiveness and stability of the PMM algorithm. In addition, the private experiment shows that the positions which should be zero are all 0. Therefore, we conclude that the PMM algorithm can simultaneously realize variable selection and parameter estimation.

Next, we examine the calculation speed of PMM algorithm from the perspective of the number of iterations. Based on 100 independent experiments, we show the average number of iterations with different sparsity levels on the right side of Fig. 1. In view of the stop condition $\|\beta(\lambda_j)\|_0 \geq n/\log(p)$ in step 10 of PDASC method, here we consider $K = 5 : 5 : 40$, which means that the sparsity level K varies from 5 to 40 by step 5. In addition, we also take the correlation into consideration and set $\kappa = [0.3; 0.5; 0.7]$. It can

be seen that for the three correlation coefficients, the average number of iterations of the PMM algorithm does not exceed 4. This phenomena fully illustrates that the calculation speed of the PMM algorithm is very fast.

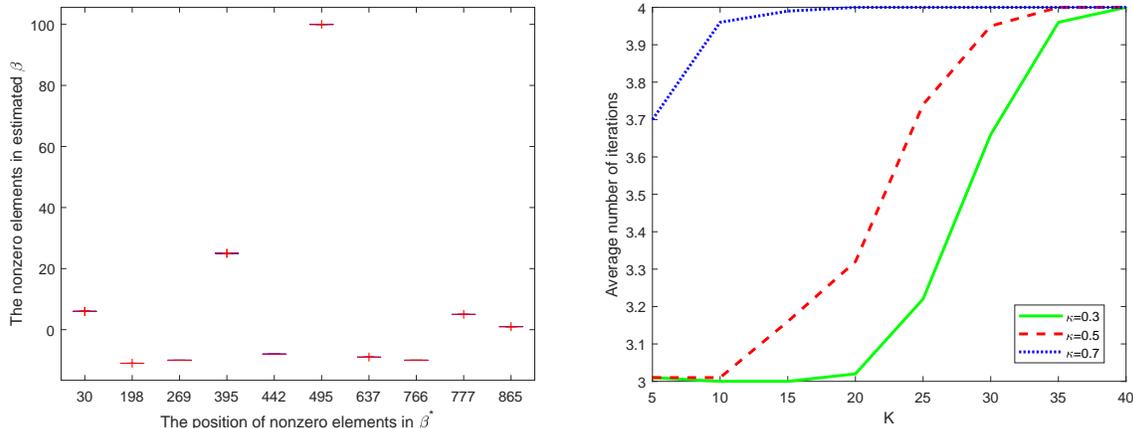


Figure 1: The behavior of PMM algorithm for MCP penalized linear regression problems based on 100 independent experiments

5.3 Comparison with SSN algorithm

In this part, we compare the PMM algorithm with the latest SSN algorithm for solving non-convex penalized high-dimensional linear regression problems based on 100 independent experiments. We set $p = 2000$ with $n = \lfloor \frac{p}{5} \rfloor$ and $K = \lfloor \frac{n}{2\log(p)} \rfloor$, where $\lfloor x \rfloor$ denotes the integer part of x for $x \geq 0$. Here, we set $\sigma_1 = 0.1$ and consider three levels of correlation, i.e., $\kappa = [0.3; 0.5; 0.7]$. It can be observed from the MATLAB package of SSN algorithm that the authors in Shi et al. (2018) lead into a key parameter “Weight” which represents the step size in the programming process. After testing, we find that the effectiveness of SSN algorithm is heavily dependent on this parameter. Here we only consider two values of 0.5 and 0.9. In addition, for the sake of fairness, we use the same continuation method for the regularization parameter in SSN algorithm, unify the maximum number of iterations to 1, and other parameters are consistent with their original papers. Simulation results are

summarized in Table 1.

Table 1: Simulation results of SSN and PMM algorithms

κ	Weight	Penalty	Method	Time	MS	CM	RE
0.3	0.5	MCP	SSN	0.04	26.01	99%	1.00e-4
			PMM	0.07	26.00	100%	1.48e-4
		SCAD	SSN	0.05	26.00	100%	1.00e-4
			PMM	0.07	26.00	100%	1.48e-4
	0.9	MCP	SSN	0.05	26.40	78%	2.00e-4
			PMM	0.08	26.00	100%	1.48e-4
0.5	0.5	MCP	SSN	0.04	25.95	98%	5.00e-4
			PMM	0.08	26.00	100%	1.41e-4
SCAD		SSN	0.05	25.96	99%	4.00e-4	
		PMM	0.07	26.00	100%	1.41e-4	
0.9	MCP	SSN	0.06	26.33	80%	2.00e-4	
		PMM	0.08	26.00	100%	1.41e-4	
	SCAD	SSN	0.06	26.62	65%	1.60e-4	
		PMM	0.07	26.00	100%	1.41e-4	
0.7	0.5	MCP	SSN	0.03	23.94	68%	6.60e-2
			PMM	0.09	26.00	100%	1.46e-4
		SCAD	SSN	0.05	25.40	93%	1.50e-2
			PMM	0.09	26.00	100%	1.46e-4
	0.9	MCP	SSN	0.04	24.96	60%	2.76e-2
			PMM	0.09	26.00	100%	1.46e-4
SCAD	SSN	0.06	26.42	67%	1.61e-4		
	PMM	0.09	26.00	100%	1.46e-4		

From the information in Table 1, we can see that the calculation speed of SSN algorithm is very fast, which thanks to its local super-linear convergence. However, since its performance is heavily dependent on the selection of the step size, the results under the fixed step sizes 0.5 and 0.9 are incomparable with PMM algorithm at present. Therefore,

in view of the fact that the SSN algorithm need to carefully adjust the step size under different problem settings, we will only compare the algorithm in this paper with the classic CD algorithm detailly in the subsequent numerical experiments.

5.4 Comparison with CD algorithm

In this section, we compare our PMM algorithm with the CD algorithm in Breheny and Huang (2011) for solving (1.3) which is summarized in Algorithm 5.4. To be fair, we here use the same continuation method for regularization parameter λ and the same stop condition $R_{kkt}^2 < 1e - 6$ at step 8. In addition, we also set $K_{max} = 1$ to improve the calculation speed of the CD algorithm.

CD algorithm

Step 0. Given λ , $\beta^0 = \mathbf{0}$, $r^0 = y - X\beta^0$, $K_{max} \in \mathbb{N}$. For $k = 0, 1, \dots, K_{max}$, do the following operations iteratively.

Step 1. For $i = 1, 2, \dots, p$, do the following operations iteratively.

Step 1.1. Calculate $z_i^k = X_i^\top r^k + \beta_i^k$, where X_i is the i th column of X and $r^k = y - X\beta^k$ is the current residual value.

Step 1.2. Update $\beta_i^{k+1} = \text{Prox}_{\lambda\|\cdot\|_1 - q(\cdot)}(z_i^k)$.

Step 1.3. Update $r^{k+1} = r^k - (\beta_i^{k+1} - \beta_i^k)X_i$.

Step 2. Check stop condition, if stop, denote the last iteration by $\hat{\beta}$. Else, $k := k + 1$.

5.4.1 Efficiency and accuracy

In this part, we compare the efficiency and accuracy of the PMM algorithm and the CD algorithm based on 100 independent experiments. We set $p = 2000$ and 5000 with $n = \lfloor \frac{p}{5} \rfloor$

and $K = \lfloor \frac{n}{2 \log(p)} \rfloor$. We consider three levels of correlation ($\kappa = [0.3; 0.5; 0.7]$) and two levels of noises ($\sigma_1 = [0.1; 1]$). Simulation results are summarized in Table 2.

From the results of MS, CM and RE in Table 2, it can be concluded that for each combination of (p, κ, σ_1) , the PMM algorithm can always achieve variable selection and parameter estimation very accurately. In addition, the PMM algorithm has better speed performance than CD algorithm for both MCP and SCAD, and PMM is about 2 ~ 9 times faster than CD. In particular, for given penalty and method, the CPU time increases with the increase of p , and decreases with the increase of σ_1 , but does not change much for different κ . In addition, it can be found that larger p can improve the accuracy of both CD and PMM, while larger σ_1 has the opposite effect. Overall, the simulation results in Table 2 illustrate that PMM outperforms CD in terms of CPU time while producing solutions of comparable quality.

5.4.2 Influence of model parameters

We now consider the effects of each of the model parameters $(n, p, K, \kappa, \sigma_1, \tau)$ on the performance of PMM and CD algorithms. Here we only give the results related to the MCP penalty, since SCAD penalty will produce a similar phenomenon. Based on 10 independent replications, we compare the performance of the considered methods in terms of average positive discovery rate (APDR), average false discovery rate (AFDR) and average combined discovery rate (ACDR) Luo and Chen (2014) defined as follows:

$$\text{APDR} = \frac{1}{10} \sum \frac{|\hat{\mathcal{A}} \cap \mathcal{A}^*|}{|\mathcal{A}^*|}, \quad \text{AFDR} = \frac{1}{10} \sum \frac{|\hat{\mathcal{A}} \cap \mathcal{A}^{*c}|}{|\hat{\mathcal{A}}|}, \quad \text{ACDR} = \text{APDR} + (1 - \text{AFDR}),$$

where \mathcal{A}^* denotes the true active set and \mathcal{A}^{*c} denotes the complement of \mathcal{A}^* . Results of APDR, AFDR and ACDR over 10 independent replications are given in Fig. 2-4, respectively. The parameters for solvers are set as follows.

- Influence of the sample size n : We set $p = 1000$, $K = 10$, $\tau = 2.7$, $\kappa = 0.2$, $\sigma_1 = 0.1$, and take $n = 20$ to 200 with a step size 20.

- Influence of the dimension p : We set $n = 200$, $K = 50$, $\tau = 2.7$, $\kappa = 0.2$, $\sigma_1 = 0.1$, and take $p = 500$ to 1000 with a step size 100 .
- Influence of the sparsity level K : We set $n = 200$, $p = 1000$, $\tau = 2.7$, $\kappa = 0.2$, $\sigma_1 = 0.1$, and take $K = 10$ to 50 with a step size 10 .
- Influence of the correlation level κ : We set $n = 200$, $p = 1000$, $K = 40$, $\tau = 2.7$, $\sigma_1 = 0.1$, and take $\kappa = 0.1$ to 0.7 with a step size 0.1 .
- Influence of the noise level σ_1 : We set $n = 200$, $p = 1000$, $K = 10$, $\tau = 2.7$, $\kappa = 0.2$, and take $\sigma_1 \in \{0.1, 0.5, 1.0, 1.5, 2.0, 2.5\}$.
- Influence of the concavity parameter τ : We set $n = 200$, $p = 1000$, $K = 10$, $\kappa = 0.2$, $\sigma_1 = 0.1$, and take $\tau \in \{1.1, 2.7, 5, 10\}$.

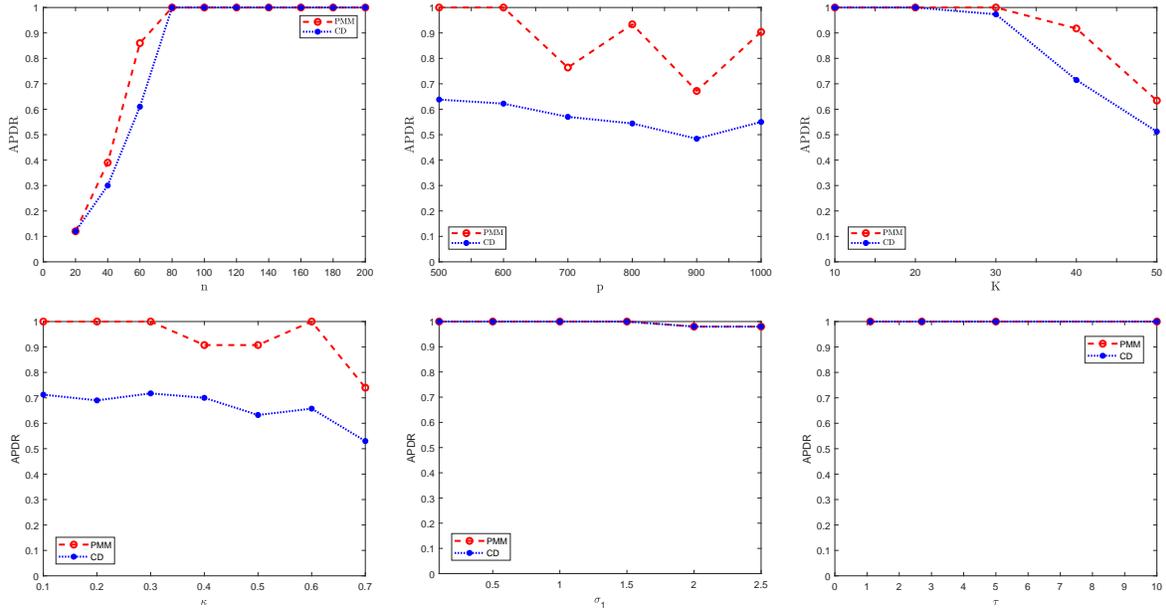


Figure 2: Numerical results of the influence of the model parameters on APDR

From the definitions of APDR, AFDR and ACDR, we can conclude that the closer APDR approaches to 1, the closer AFDR approaches to 0, and the closer ACDR approaches to 2, the higher the solution quality. From Fig 2-4, we can see that for the change interval of

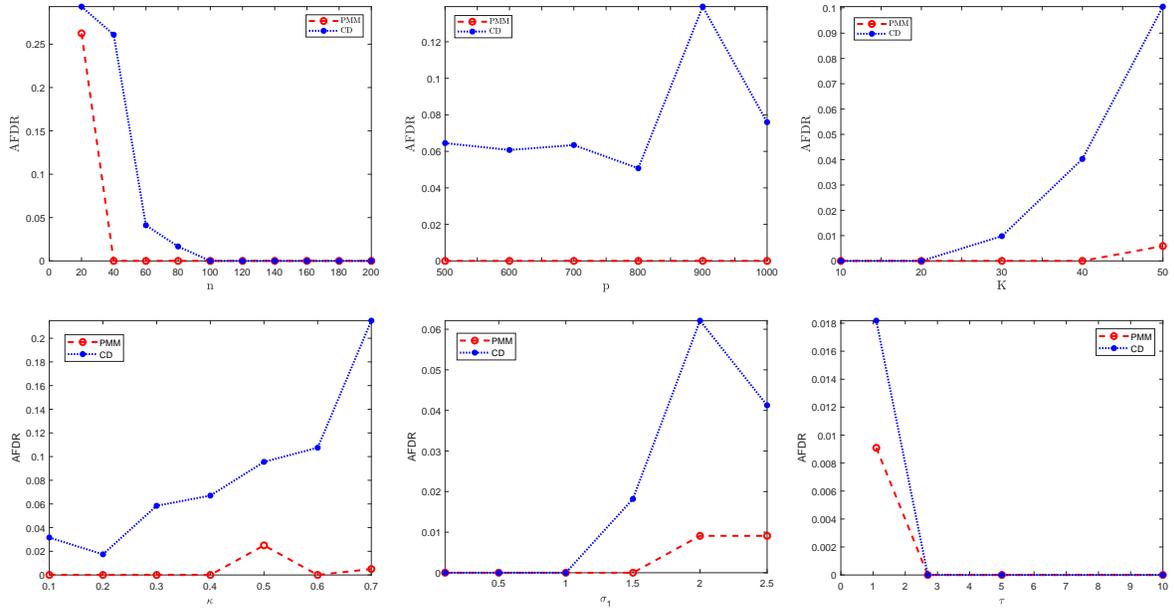


Figure 3: Numerical results of the influence of the model parameters on AFDR

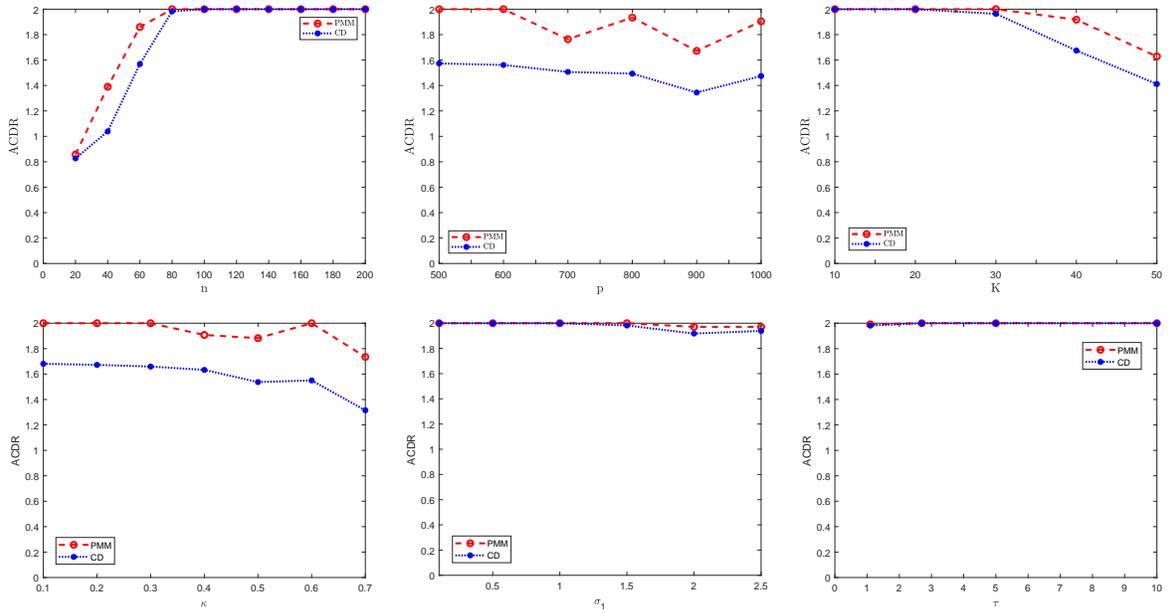


Figure 4: Numerical results of the influence of the model parameters on ACDR

different parameters, PMM can always achieve more expected results. Therefore, compared with CD, PMM is more robust to the considered parameters for solving the MCP penalized least squares problems.

5.4.3 Numerical comparison with real data

In this subsection, we test CD and PMM algorithms with the test instances (X, y) obtained from large-scale regression problems in the LIBSVM data sets, which is available at <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>. These data sets are collected from UCI, StatLib, Delve, 10-K Corpus, GWF01a. For computational efficiency, zero columns in X are removed. As suggested in Huang et al. (2010), in addition to the data sets **log1p.E2006.train** and **E2006.train**, we expand the original features of the remaining data sets by using polynomial basis functions over those features. For example, the last digit in **abalone7** indicates that an order 7 polynomial is used to generate the basis functions. This naming convention is also used in the rest of the expanded data sets. These test instances are quite difficult in terms of the problem dimensions and the largest eigenvalue of XX^\top , which is denoted as $\lambda_{max}(XX^\top)$, one can refer to the first three columns of Table 3. It is worth noting that for these difficult real data, we appropriately reduce the accuracy requirements in the termination conditions. In addition to setting $R_{kkt}^2 < 5e - 2$ for **space_ga9** and $R_{kkt}^2 < 8e - 3$ for **bodyfat7**, we set $R_{kkt}^2 < 5e - 3$ as the termination condition for all other data.

Table 3 reports the detailed numerical results for CD and PMM in solving large-scale regression problems. In the table, “NNZ” denotes the number of nonzeros in the estimated solution, and other symbols are the same as the previous simulation experiment. From the results in Table 3, we can see that PMM can solve all the instances to the desired accuracy despite the huge dimensions and the possibly badly conditioned data sets. More specifically, PMM is able to solve the instance **log1p.E2006.train** with approximately 4.3 million features to accuracy $R_{kkt}^2 = 9.98e - 6$ in 95 seconds. But CD only meets the accuracy requirement for **E2006.train**. In addition, for solving these data sets, CD needs much more time than PMM. For example, for the instance **cpusmall7**, we can see that PMM is at least 144 times faster than CD. The superior numerical performance of PMM indicates that

it is a robust, high-performance solver for MCP/SCAD penalized high-dimensional linear regression problems.

6 Conclusion

Based on the DC property of MCP and SCAD penalties, we developed a global two-stage algorithm for the MCP/SCAD penalized linear regression problems in high-dimensional settings. A key idea for making the proposed algorithm to be efficient is to use the PDASC algorithm to solve the corresponding sub-problems. We established the global convergence of the proposed algorithm and verified the iterative sequence converges to a d-stationary point of the considered problems. Finally, a large number of inspiring numerical experiments have once again verified the effectiveness of the proposed algorithm.

Since each non-convex penalty can be expressed as the difference of two convex functions, the research in this paper can be directly extended to other non-convex penalized high-dimensional linear regression problems. In addition, extending the algorithm in this paper to the regression problems with other loss functions is also a very interesting and promising research direction.

Acknowledgements

The work of Zhou Yu is supported in part by the National Natural Science Foundation of China (Grant No. 11971170).

References

Ahn, M., J.-S. Pang, and J. Xin (2017). Difference-of-convex learning: directional stationarity, optimality, and sparsity. *SIAM Journal on Optimization* 27(3), 1637–1665.

- Boyd, S., N. Parikh, and E. Chu (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers.
- Breheny, P. and J. Huang (2011). Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection. *The Annals of Applied Statistics* 5(1), 232–253.
- Candes, E. J. and T. Tao (2005). Decoding by linear programming. *IEEE Transactions on Information Theory* 51(12), 4203–4215.
- Chartrand, R. (2007). Exact reconstruction of sparse signals via nonconvex minimization. *IEEE Signal Processing Letters* 14(10), 707–710.
- Chen, L. and Y. Gu (2014). The convergence guarantees of a non-convex approach for sparse recovery. *IEEE Transactions on Signal Processing* 62(15), 3754–3767.
- Chen, S. S., D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Review* 43(1), 129–159.
- Cui, Y., J.-S. Pang, and B. Sen (2018). Composite difference-max programs for modern statistical estimation problems. *SIAM Journal on Optimization* 28(4), 3344–3374.
- Donoho, D. L. and I. M. Johnstone (1995). Adapting to unknown smoothness via wavelet shrinkage. *Journal of the American Statistical Association* 90(432), 1200–1224.
- Efron, B., T. Hastie, I. Johnstone, and R. Tibshirani (2004). Least angle regression. *The Annals of Statistics* 32(2), 407–499.
- Fan, J. and R. Li (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* 96(456), 1348–1360.
- Fan, Q., Y. Jiao, and X. Lu (2014). A primal dual active set algorithm with continuation for compressed sensing. *IEEE Transactions on Signal Processing* 62(23), 6276–6285.

- Frank, L. E. and J. H. Friedman (1993). A statistical view of some chemometrics regression tools. *Technometrics* 35(2), 109–135.
- Fu, W. J. (1998). Penalized regressions: the bridge versus the lasso. *Journal of Computational and Graphical Statistics* 7(3), 397–416.
- Gong, P., C. Zhang, Z. Lu, J. Huang, and J. Ye (2013). A general iterative shrinkage and thresholding algorithm for non-convex regularized optimization problems. *Proceedings of the 30th International Conference on Machine Learning* 28(2), 37–45.
- Hintermüller, M., K. Ito, and K. Kunisch (2002). The primal-dual active set strategy as a semismooth Newton method. *SIAM Journal on Optimization* 13(3), 865–888.
- Hofmann, B. and T. Hohage (2011). Generalized Tikhonov regularization: Basic theory and comprehensive results on convergence rates. *Fakultat fur Mathematik*.
- Huang, J., Y. Jiao, B. Jin, J. Liu, X. Lu, and C. Yang (2021). A unified primal dual active set algorithm for nonconvex sparse recovery. *Statistical Science* 36(2), 215–238.
- Huang, L., J. Jia, B. Yu, B.-G. Chun, P. Maniatis, and M. Naik (2010). Predicting execution time of computer programs using sparse polynomial regression. *Advances in Neural Information Processing Systems* 23, 883–891.
- Hunter, D. R. and R. Li (2005). Variable selection using MM algorithms. *The Annals of Statistics* 33(4), 1617–1642.
- Le Thi, H. A., T. P. Dinh, H. M. Le, and X. T. Vo (2015). DC approximation approaches for sparse optimization. *European Journal of Operational Research* 244(1), 26–46.
- Lee, J. D., Y. Sun, and M. A. Saunders (2014). Proximal Newton-type methods for minimizing composite functions. *SIAM Journal on Optimization* 24(3), 1420–1443.

- Li, G. and T. K. Pong (2018). Calculus of the exponent of Kurdyka-Lojasiewicz inequality and its applications to linear convergence of first-order methods. *Foundations of Computational Mathematics* 18(5), 1199–1232.
- Li, X., L. Yang, J. Ge, J. Haupt, T. Zhang, and T. Zhao (2017). On quadratic convergence of DC proximal Newton algorithm in nonconvex sparse learning. *Advances in Neural Information Processing Systems* 30, 2742–2752.
- Li, X. D., D. F. Sun, and K. C. Toh (2018). A Highly Efficient Semismooth Newton Augmented Lagrangian Method for Solving LASSO Problems. *SIAM Journal on Optimization* 28(1), 433–458.
- Luo, S. and Z. Chen (2014). Sequential Lasso cum EBIC for feature selection with ultra-high dimensional feature space. *Journal of the American Statistical Association* 109(507), 1229–1240.
- Mazumder, R., J. H. Friedman, and T. Hastie (2011). Sparsenet: Coordinate descent with nonconvex penalties. *Journal of the American Statistical Association* 106(495), 1125–1138.
- Meinshausen, N. and P. Bühlmann (2006). High-dimensional graphs and variable selection with the lasso. *The Annals of Statistics* 34(3), 1436–1462.
- Micchelli, C. A., L. Shen, and Y. Xu (2011). Proximity algorithms for image models: denoising. *Inverse Problems* 27(4), 045009.
- Natarajan, B. K. (1995). Sparse approximate solutions to linear systems. *SIAM Journal on Computing* 24(2), 227–234.
- Pang, J.-S., M. Razaviyayn, and A. Alvarado (2017). Computing B-stationary points of nonsmooth DC programs. *Mathematics of Operations Research* 42(1), 95–118.

- Rockafellar, R. T. (2015). Convex analysis.
- Shi, Y., J. Huang, Y. Jiao, and Q. Yang (2018). Semi-smooth Newton algorithm for non-convex penalized linear regression. *arXiv preprint arXiv:1802.08895*.
- Tang, P., C. Wang, D. Sun, and K.-C. Toh (2020). A sparse semismooth Newton based proximal majorization-minimization algorithm for nonconvex square-root-loss regression problems. *Journal of Machine Learning Research* 21(226), 1–38.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)* 58(1), 267–288.
- Wang, L., Y. Kim, and R. Li (2013). Calibrating non-convex penalized regression in ultra-high dimension. *The Annals of Statistics* 41(5), 2505–2536.
- Wu, T. T. and K. Lange (2008). Coordinate descent algorithms for lasso penalized regression. *The Annals of Applied Statistics* 2(1), 224–244.
- Zhang, C.-H. (2010a). Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics* 38(2), 894–942.
- Zhang, T. (2010b). Analysis of multi-stage convex relaxation for sparse regularization. *Journal of Machine Learning Research* 11(3), 1081–1107.
- Zhao, P. and B. Yu (2006). On model selection consistency of Lasso. *Journal of Machine Learning Research* 7, 2541–2563.
- Zou, H. and R. Li (2008). One-step sparse estimates in nonconcave penalized likelihood models. *The Annals of Statistics* 36(4), 1509–1533.

Table 2: Simulation results of CD and PMM algorithms

p	κ	σ_1	Penalty	Method	Time	MS	CM	RE
2000	0.3	0.1	MCP	CD	0.69	26.00	100%	1.48e-4
				PMM	0.08	26.00	100%	1.48e-4
			SCAD	CD	0.70	26.00	100%	1.48e-4
				PMM	0.08	26.00	100%	1.48e-4
		1	MCP	CD	0.52	26.00	100%	1.50e-3
				PMM	0.07	26.00	100%	1.50e-3
	SCAD		CD	0.52	26.01	99%	1.50e-3	
			PMM	0.08	26.00	100%	1.50e-3	
	0.5	0.1	MCP	CD	0.70	26.00	100%	1.41e-4
				PMM	0.08	26.00	100%	1.41e-4
			SCAD	CD	0.71	26.00	100%	1.41e-4
				PMM	0.08	26.00	100%	1.41e-4
		1	MCP	CD	0.53	26.00	100%	1.40e-3
				PMM	0.07	26.00	100%	1.40e-3
	SCAD		CD	0.54	26.00	100%	1.40e-3	
			PMM	0.08	26.00	100%	1.40e-3	
	0.7	0.1	MCP	CD	0.70	26.00	100%	1.46e-4
				PMM	0.10	26.00	100%	1.46e-4
			SCAD	CD	0.71	26.00	100%	1.46e-4
				PMM	0.10	26.00	100%	1.46e-4
		1	MCP	CD	0.53	26.05	97%	1.60e-3
				PMM	0.08	26.00	100%	1.50e-3
	SCAD		CD	0.53	26.10	94%	1.60e-3	
			PMM	0.07	26.00	100%	1.50e-3	
5000	0.3	0.1	MCP	CD	2.84	58.00	100%	9.29e-5
				PMM	1.07	58.00	100%	9.29e-5
			SCAD	CD	2.85	58.00	100%	9.29e-5
				PMM	1.05	58.00	100%	9.29e-5
		1	MCP	CD	2.25	58.00	100%	9.29e-4
				PMM	1.04	58.00	100%	9.29e-4
	SCAD		CD	2.48	58.00	100%	9.00e-4	
			PMM	1.03	58.00	100%	9.29e-4	
	0.5	0.1	MCP	CD	2.85	58.00	100%	9.34e-5
				PMM	1.06	58.00	100%	9.34e-5
			SCAD	CD	2.85	58.00	100%	9.34e-5
				PMM	1.03	58.00	100%	9.34e-5
		1	MCP	CD	2.45	58.00	100%	9.34e-4
				PMM	1.10	58.00	100%	9.34e-4
	SCAD		CD	2.49	58.00	100%	9.00e-4	
			PMM	1.14	58.00	100%	9.34e-4	
	0.7	0.1	MCP	CD	2.87	58.00	100%	9.77e-5
				PMM	1.31	58.00	100%	9.77e-5
			SCAD	CD	2.87	58.00	100%	9.77e-5
				PMM	1.29	58.00	100%	9.77e-5
		1	MCP	CD	2.41	58.03	99%	9.96e-4
				PMM	1.37	58.00	100%	9.77e-4
	SCAD		CD	2.28	58.03	99%	1.00e-3	
			PMM	1.40	58.00	100%	9.77e-4	

Table 3: Real results of CD and PMM algorithms

Data name	n,p	$\lambda_{max}(XX^T)$	Penalty	NNZ	Method	R^2_{kkt}	Time
log1p.E2006.train	16087,4265669	5.86e+7	MCP	6	CD	3.91e-2	4.08e+3
					PMM	9.98e-6	9.39e+1
			SCAD	6	CD	5.36e-2	4.15e+3
					PMM	9.97e-6	9.52e+1
E2006.train	16087,150348	1.91e+5	MCP	6	CD	2.50e-3	7.30e+1
					PMM	4.90e-3	4.33e+1
			SCAD	6	CD	2.50e-3	7.33e+1
					PMM	4.00e-3	3.68e+1
abalone7	4177,6435	5.21e+5	MCP	7	CD	9.66e-1	1.58e+1
					PMM	2.86e-4	1.51e+0
			SCAD	7	CD	9.64e-1	1.65e+1
					PMM	2.77e-4	1.61e+0
bodyfat7	252,116280	5.29e+4	MCP	9	CD	5.57e-2	1.18e+1
					PMM	6.60e-3	2.30e+0
			SCAD	1	CD	7.63e-2	1.18e+1
					PMM	7.70e-3	2.23e+0
cpusmall7	8192,50388	8.01e+7	MCP	1103	CD	1.00e+0	2.97e+2
					PMM	4.10e-3	2.05e+0
			SCAD	1105	CD	1.00e+0	3.78e+2
					PMM	7.30e-4	4.57e+0
housing7	506,77520	3.28e+5	MCP	44	CD	2.64e-2	1.48e+1
					PMM	1.70e-3	2.38e-1
			SCAD	46	CD	3.62e-2	1.41e+1
					PMM	8.74e-4	1.68e-1
mg9	1385,5005	4.78e+3	MCP	9	CD	6.31e-1	2.29e+0
					PMM	4.60e-3	1.56e+0
			SCAD	9	CD	6.55e-1	2.21e+0
					PMM	4.50e-3	1.96e+0
mpg7	392,3432	1.28e+4	MCP	26	CD	4.92e-2	8.24e-1
					PMM	2.50e-3	3.65e-2
			SCAD	27	CD	6.75e-2	7.69e-1
					PMM	2.30e-3	3.75e-2
pyrim5	74,169911	1.22e+6	MCP	326	CD	5.81e-2	1.07e+1
					PMM	4.92e-4	5.60e-2
			SCAD	327	CD	7.96e-2	9.94e+0
					PMM	2.48e-4	5.37e-2
space_ga9	3107,5005	4.01e+3	MCP	8	CD	1.89e-1	6.23e+0
					PMM	4.60e-2	1.99e+0
			SCAD	9	CD	2.59e-1	6.26e+0
					PMM	4.50e-2	2.24e+0
triazines4	186,557845	2.08e+7	MCP	983	CD	4.29e-2	9.61e+1
					PMM	7.39e-4	6.31e+1
			SCAD	983	CD	5.87e-2	9.48e+1
					PMM	1.86e-4	6.67e+1