

On Regularized Sparse Logistic Regression

Mengyuan Zhang
Clemson University
Clemson, SC, USA
mengyuz@clemson.edu

Kai Liu
Clemson University
Clemson, SC, USA
kail@clemson.edu

Abstract—Sparse logistic regression is for classification and feature selection simultaneously. Although many studies have been done to solve ℓ_1 -regularized logistic regression, there is no equivalently abundant work on solving sparse logistic regression with nonconvex regularization term. In this paper, we propose a unified framework to solve ℓ_1 -regularized logistic regression, which can be naturally extended to nonconvex regularization term, as long as certain requirement is satisfied. In addition, we also utilize a different line search criteria to guarantee monotone convergence for various regularization terms. Empirical experiments on binary classification tasks with real-world datasets demonstrate our proposed algorithms are capable of performing classification and feature selection effectively at a lower computational cost.

Index Terms—logistic regression, sparsity, feature selection

I. INTRODUCTION

Logistic regression has been applied widely in many areas as a method of classification. The goal of logistic regression is to maximize the likelihood based on the observation of training samples, with its objective function formulated as follows with a natural and meaningful probabilistic interpretation:

$$\min_{\beta} \sum_{i=1}^n -\ln p(y_i | \mathbf{x}_i; \beta) = \sum_{i=1}^n \ln(1 + \exp(\beta^T \mathbf{x}_i)) - y_i \beta^T \mathbf{x}_i, \quad (1)$$

where \mathbf{x}_i and y_i denote the i -th sample and its label.

Though logistic regression is straightforward and effective, its performance can be diminished due to over-fitting [10], especially when the dimensionality of features is very high compared to the number of available training samples. Therefore, regularization term is usually introduced to alleviate over-fitting issue [23]. Also, in applications with high-dimensional data, it's desirable to obtain sparse solutions, since in this way we are conducting classification and feature selection at the same time. Therefore, ℓ_1 -regularized logistic regression has received more attention with the sparsity-inducing property and its superior empirical performance [1].

More recent studies show that ℓ_1 -norm regularization may suffer from implicit bias problem that would cause significantly biased estimates, and such bias problem can be mitigated by a nonconvex penalty [9]. Therefore, nonconvex regularization has also been studied to induce sparsity in logistic regression [29]. Solving sparse logistic regression with convex ℓ_1 -norm regularizer or with nonconvex term using a unified algorithm has been studied in [19]. However, it imposes strong regularity condition on the nonconvex term and it transfers the nonconvex problem into an ℓ_1 -norm regularized surrogate convex function,

which limits its generality. As a contribution, we extend the scope of nonconvex penalties to a much weaker assumption and compare the performance of different regularization terms with a unified optimization framework.

In this paper, we solve ℓ_1 -regularized (sparse) logistic regression by proposing a novel framework, which can be applied to non-convex regularization term as well. The idea of our proposed method stems from the well know Iterative Shrinkage Thresholding Algorithm (ISTA) and its accelerated version Fast Iterative Shrinkage Thresholding Algorithm (FISTA) [3], upon which we modify the step-size setting and line search criteria to make the algorithm applicable for both convex and nonconvex regularization terms with empirical faster convergence rate. To be clear, we call any logistic regression with regularization term that can produce sparse solutions as *sparse logistic regression*, therefore, the term is not only limited to ℓ_1 -norm regularization.

II. RELATED WORK

Due to the NP-hardness of ℓ_0 -norm constraint problem, being its tightest convex envelope, ℓ_1 -norm is widely taken as an alternative to induce sparsity [17], [18], [32]. The main drawback of ℓ_1 -norm is it's non-differentiable, which makes the computation challenging compared to squared ℓ_2 -norm. Sub-gradient is an option but can be very slow even if we disregard the fact that it is non-monotone. Besides, there has been active research on numerical algorithms to solve ℓ_1 -regularized logistic regression. Among these, an intuitive idea is as the challenge originates from the non-smoothness of ℓ_1 -norm, we can make it 'smoothable'. For example, in [24], the ℓ_1 -norm is approximated by a smooth function that is readily solvable by any applicable optimization method. The *iteratively reweighted least squares least angle regression* (IRLS-LARS) method converts the original problem into a smooth function by the equivalent ℓ_1 -norm constrained ball [14]. In addition, coordinate descent method has been utilized to solve the ℓ_1 -regularized logistic regression as well: in [7] and [20], cyclic coordinate descent is used for optimizing Bayesian logistic regression. Besides, the interior-point method is another option with truncated Newton step and conjugated gradient iterations [12]. We refer readers to the papers and references therein.

The bias of ℓ_1 -norm is implicitly introduced as it penalizes the parameters with larger coefficients more than the smaller ones. Recently, nonconvex regularization terms have drawn considerable interest in sparse logistic regression since it is able

TABLE I
NONCONVEX REGULARIZATION TERMS IN SPARSE LOGISTIC REGRESSION.

Name	Formulation
ℓ_p -norm, $0 < p < 1$	$\lambda \beta_i ^q$
Capped ℓ_1 -norm	$\lambda \min(\beta_i , \epsilon)$
SCAD	$\begin{cases} \lambda \beta_i , & \beta_i \leq \lambda \\ \frac{-\beta_i^2 + 2\theta\lambda \beta_i - \lambda^2}{2(\theta-1)}, & \lambda < \beta_i \leq \theta\lambda \\ \frac{(\theta+1)\lambda^2}{2}, & \beta_i > \theta\lambda \end{cases}$
MCP	$\begin{cases} \lambda \beta_i - \frac{\beta_i^2}{2\theta}, & \text{if } \beta_i \leq \theta\lambda \\ \frac{\theta\lambda^2}{2}, & \text{if } \beta_i > \theta\lambda \end{cases}$

to ameliorate the bias problem of ℓ_1 -norm, and it acts as one main driver of recent progress in nonconvex and nonsmooth optimization [28].

It is worth noting that there have been many studies trying to solve the ℓ_0 -regularized problems directly. The most common method is to conduct coordinate descent for nonconvex regularized logistic regression [5], [22]. The alternating direction method of multipliers (ADMM) inspires the development of incremental aggregated proximal ADMM to solve nonconvex optimization problems, and it achieves good results in sparse logistic regression [11]. In [21], the momentumized iterative shrinkage thresholding (MIST) algorithm is proposed to minimize the nonconvex criterion for linear regression problems, and similar ideas can be applied to logistic regression as well. Besides the ℓ_0 -norm, other nonconvex penalties are also explored. The minimax concave penalty (MCP) is studied in [31] together with a penalized linear unbiased selection (PLUS) algorithm, and it shows the MCP is unbiased with superior selection accuracy. The smoothly clipped absolute deviation (SCAD) penalty is proposed in [6], which corresponds to a quadratic spline function, and the study shows SCAD penalty outperforms the ℓ_1 -norm regularizer significantly, and it has the best performance in selecting significant variables without introducing excessive biases.

III. OPTIMIZATION ALGORITHMS

A. Algorithms for ℓ_1 Regularized Sparse Logistic Regression

We first consider sparse logistic regression problem as:

$$\min_{\beta} f(\beta) = \underbrace{\sum_i \ln(1 + \exp(\mathbf{x}_i^T \beta)) - y_i(\mathbf{x}_i^T \beta)}_{l(\beta)} + \underbrace{\lambda \|\beta\|_1}_{g(\beta)}. \quad (2)$$

Theorem III.1. $\nabla l(\beta)$ in Eq (2) is Lipschitz continuous, and its Lipschitz constant is $L := \frac{1}{4} \lambda_{\max}(\mathbf{X}\mathbf{X}^T)$.

Proof. Let \mathbf{X} be the matrix of training samples, where the i -th column \mathbf{x}_i represents the i -th sample. We have

$$\nabla l(\beta) = \mathbf{X}(\mathbf{p} - \mathbf{y}), \text{ where } p_i = \frac{1}{1 + \exp(-\langle \beta, \mathbf{x}_i \rangle)} \quad (3)$$

and

$$\nabla^2 l(\beta) = \mathbf{X} \text{diag}((p_i(1-p_i))) \mathbf{X}^T. \quad (4)$$

By the mean value theorem, we know there exists $\mathbf{c} \in (\beta, \beta + \Delta)$ such that $\nabla l(\beta + \Delta) - \nabla l(\beta) = \nabla^2 l(\mathbf{c})\Delta$. Since

$$\begin{aligned} \|\nabla^2 l(\beta)\| &= \lambda_{\max}(\mathbf{X} \text{diag}((p_i(1-p_i))) \mathbf{X}^T) \\ &= \max_{\|\mathbf{v}\|=1} \mathbf{v}^T \mathbf{X} \text{diag}((p_i(1-p_i))) \mathbf{X}^T \mathbf{v} \\ &= \max_{\|\mathbf{v}\|=1} \mathbf{v}^T [\sum (\mathbf{x}_i p_i (1-p_i) \mathbf{x}_i^T)] \mathbf{v} \\ &\leq \frac{1}{4} \max_{\|\mathbf{v}\|=1} \mathbf{v}^T [\sum (\mathbf{x}_i \mathbf{x}_i^T)] \mathbf{v} = \frac{1}{4} \lambda_{\max}(\mathbf{X}\mathbf{X}^T), \end{aligned} \quad (5)$$

similarly $\|\nabla^2 l(\mathbf{c})\| \leq \frac{1}{4} \lambda_{\max}(\mathbf{X}\mathbf{X}^T)$, and we have

$$\begin{aligned} \|\nabla l(\beta + \Delta) - \nabla l(\beta)\| &= \|\nabla^2 l(\mathbf{c})\Delta\| \\ &\leq \|\nabla^2 l(\mathbf{c})\| \|\Delta\| \leq \frac{1}{4} \lambda_{\max}(\mathbf{X}\mathbf{X}^T) \|\Delta\|. \end{aligned} \quad (6)$$

Thus $\nabla l(\beta)$ is Lipschitz continuous with $\frac{1}{4} \lambda_{\max}(\mathbf{X}\mathbf{X}^T)$. \square

We use ISTA and FISTA with backtracking line search to solve Eq (2), which are described in Algorithm 1 and 2 respectively, where $p_L(\beta)$ represents the proximal operator defined as $p_L(\beta) = \arg\min_{\mathbf{w}} \frac{L}{2} \|\mathbf{w} - (\beta - \frac{1}{L} \nabla l(\beta))\|^2 + g(\mathbf{w})$. The line search stopping criterion is: $f(p_L(\beta_{k-1})) \leq q_L(p_L(\beta_{k-1}), \beta_{k-1})$, where $f(\beta)$ is defined in Eq (2), and

$$\begin{aligned} &q_L(p_L(\beta), \beta) \\ &= l(\beta) + \langle p_L(\beta) - \beta, \nabla l(\beta) \rangle + \frac{L}{2} \|p_L(\beta) - \beta\|^2 + g(p_L(\beta)). \end{aligned} \quad (7)$$

One potential drawback of the vanilla ISTA and FISTA is the initial step size which is randomly set as $L_0 > 0$ and keep increasing L_k during update to satisfy the line search stopping criterion. In case L_0 is larger than the Lipschitz constant L of $\nabla l(\beta)$, the step size can be too small to obtain optimal solution rapidly [33]. Thus, different from vanilla ISTA, in Algorithm 1, we first initialize the stepsize with Lipschitz continuous constant and then utilize Barzilai-Borwein (BB) rule to serve as a starting point for backtracking line search:

$$\delta_k = \beta_{k-1} - \beta_{k-2}, \mathbf{v}_k = \nabla l(\beta_{k-1}) - \nabla l(\beta_{k-2}), L_k = \frac{\langle \delta_k, \mathbf{v}_k \rangle}{\langle \delta_k, \delta_k \rangle}. \quad (8)$$

Experiments on synthetic data show that with BB rule, ISTA with randomly initialized step size will admit faster convergence during update, which is demonstrated in Figure 1. In Algorithm 2, we also initialize the step size as $\frac{1}{L}$ to avoid too small step size. However, to guarantee $O(\frac{1}{k^2})$ convergence rate for FISTA, we still need the step size to be monotonically nonincreasing, thus BB rule cannot be utilized in FISTA.

Besides the step size setting in BB rule aforementioned, another option proposed by us is to find the largest step-size by searching *reversely*: we start by setting the step-size to $\frac{1}{L}$ in each iteration and keep enlarging it until the line search condition is not satisfied and take the last step-size satisfying the criterion. In this way, we are able to find the largest step

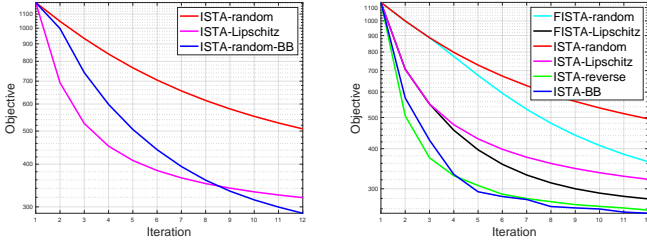


Fig. 1. ISTA with different settings. Fig. 2. ℓ_1 -norm: convergence plots.

Algorithm 1 ISTA-BB: ISTA with Lipschitz constant and BB rule to solve Eq (2).

Initialize β_0 , step size $\frac{1}{L_0}$ as $\frac{1}{L}$, where L is the Lipschitz constant of $\nabla l(\beta)$, set $\eta > 1$;

repeat

- 1) Start from $k = 2$, update the step size $\frac{1}{L_k}$ using the Barzilai-Borwein (BB) rule
- 2) Find the smallest nonnegative integer i_k such that with $\bar{L} = \eta^{i_k} L_k$ we have $f(p_{\bar{L}}(\beta_{k-1})) \leq q_{\bar{L}}(p_{\bar{L}}(\beta_{k-1}), \beta_{k-1})$.
- 3) Set $L_k = \eta^{i_k} L_k$ and update $\beta_k = p_{L_k}(\beta_{k-1})$

until convergence

Algorithm 2 FISTA-Lipschitz: FISTA with Lipschitz constant to solve Eq (2).

Initialize β_0 , step size $\frac{1}{L_0}$ as $\frac{1}{L}$, where L is the Lipschitz constant of $\nabla l(\beta)$, set $\eta > 1$, $\mathbf{w}_1 = \beta_0$, $t_1 = 1$;

repeat

- 1) Find the smallest nonnegative integer i_k such that with $\bar{L} = \eta^{i_k} L_{k-1}$ we have $f(p_{\bar{L}}(\mathbf{w}_k)) \leq q_{\bar{L}}(p_{\bar{L}}(\mathbf{w}_k), \mathbf{w}_k)$.
- 2) Set $L_k = \eta^{i_k} L_{k-1}$ and update

$$\begin{aligned} \beta_k &= p_{L_k}(\mathbf{w}_k) \\ t_{k+1} &= \frac{1 + \sqrt{1 + 4t_k^2}}{2} \\ \mathbf{w}_{k+1} &= \beta_k + \frac{t_k - 1}{t_{k+1}}(\beta_k - \beta_{k-1}) \end{aligned} \quad (9)$$

until convergence

size satisfying the line search criterion for each update iteration. The proposed method is summarized in Algorithm 3. Again, unlike the conventional ISTA, the step size in Algorithm 3 is not decreasing monotonically, and to compute the largest singular value with lower computational cost, we can utilize the power iteration method [4] to get rid of the computationally expensive singular value decomposition.

By Algorithm 3, the objective decreases much faster than vanilla ISTA. From Figure 2 we can see that for FISTA, when we initialize the step-size with $\frac{1}{L}$ (FISTA-Lipschitz), it has better convergence performance than initialized with a random number (FISTA-random), which might be smaller than $\frac{1}{L}$. For ISTA, the two algorithms proposed by us (ISTA-BB and ISTA-reverse) have similar performance and they obviously

Algorithm 3 ISTA-reverse: ISTA with Lipschitz constant and reverse step size searching to solve Eq (2).

Initialize β_0 randomly, step size $\frac{1}{L_0}$ as $\frac{1}{L}$, where L is the Lipschitz constant of $\nabla l(\beta)$, set $\eta > 1$;

repeat

- 1) Find the smallest nonnegative integer i_k such that with $\bar{L} = L_0/\eta^{i_k}$ we have $f(p_{\bar{L}}(\beta_{k-1})) > q_{\bar{L}}(p_{\bar{L}}(\beta_{k-1}), \beta_{k-1})$.
- 2) Set $L_k = L_0/\eta^{i_k-1}$ and update $\beta_k = p_{L_k}(\beta_{k-1})$

until convergence

outperform vanilla ISTA with backtracking line search, either ISTA-random or ISTA-Lipschitz.

The convergence proof of Algorithm 1, Algorithm 2, and Algorithm 3 can be easily adapted from the proof in [3], here we only present the key theorems of the convergence rate.

Theorem III.2. Let $\{\beta_k\}$ be the sequence generated by either Algorithm 1 or Algorithm 3, for any $k > 1$,

$$f(\beta_k) - f(\beta^*) = O\left(\frac{1}{k}\right), \quad (10)$$

where β^* is an optimal solution to Eq (2).

Theorem III.3. Let $\{\beta_k\}$ be the sequence generated by Algorithm 2, for any $k > 1$,

$$f(\beta_k) - f(\beta^*) = O\left(\frac{1}{k^2}\right), \quad (11)$$

where β^* is an optimal solution to Eq (2).

B. Algorithms for Nonconvex Regularized Sparse Logistic Regression

While the ℓ_1 -norm regularization is convenient since it's convex, several studies show that sometimes nonconvex regularization term can have better performance [28] though it turns the objective to nonconvex and even nonsmooth, which is challenging to obtain optimal solution. Current literature lacks a unified yet simple framework that works for both convex and a wide class of nonconvex regularization terms. With this consideration, we would like to list a bunch of nonconvex regularization terms that on one hand can ameliorate the bias problem, and on the other hand, can be solved with the same algorithms for the convex term.

Nonconvex regularization terms can be written as the difference between two convex functions as long as the Hessian is bounded [30]. For such nonconvex penalties, we are able to solve by ISTA with slight modifications.

$$\min_{\beta} f(\beta) = \underbrace{\sum_i \ln(1 + \exp(\mathbf{x}_i^T \beta)) - y_i(\mathbf{x}_i^T \beta)}_{l(\beta)} + \underbrace{g_1(\beta) - g_2(\beta)}_{g(\beta)}, \quad (12)$$

where the Hessian of $g(\beta)$ is bounded. We summarize our methods in Algorithm 4 and 5 with modified backtracking line search criteria:

$$f(\beta_{k+1}) \leq f(\beta_k) - \frac{L_k}{2} \|\beta_{k+1} - \beta_k\|^2. \quad (13)$$

TABLE II
EXAMPLES OF REGULARIZATION TERMS AND CORRESPONDING PROXIMAL OPERATORS.

Penalty	$g_1(\beta_i)$	$g_2(\beta_i)$	Proximal operator
SCAD	$\lambda \beta_i $	$\begin{cases} 0, & \beta_i \leq \lambda \\ \frac{\beta_i^2 - 2\lambda \beta_i + \lambda^2}{2(\theta-1)}, & \lambda < \beta_i \leq \theta\lambda \\ \lambda \beta_i - \frac{(\theta+1)\lambda^2}{2}, & \beta_i > \theta\lambda \end{cases}$	$\begin{cases} \text{sign}(t) \max(t - \lambda), 0, & t \leq 2\lambda \\ \frac{(\theta-1)t - \text{sign}(t)\theta\lambda}{\theta-2}, & 2\lambda < t \leq \theta\lambda \\ t, & t > \theta\lambda \end{cases}$
MCP	$\lambda \beta_i $	$\begin{cases} \frac{\beta_i^2}{2\theta}, & \text{if } \beta_i \leq \theta\lambda \\ \lambda \beta_i - \frac{\theta\lambda^2}{2}, & \text{if } \beta_i > \theta\lambda \end{cases}$	$\begin{cases} 0, & t \leq \lambda \\ \frac{\text{sign}(t)(t - \lambda)}{1 - 1/\theta}, & \lambda < t \leq \theta\lambda \\ t, & t > \theta\lambda \end{cases}$
ℓ_1 -norm			$\text{sign}(t) \max(t - \lambda, 0)$

Algorithm 4 ISTA-BB: ISTA with Lipschitz constant and BB rule to solve Eq (12).

Initialize β_0 , step size $\frac{1}{L_0}$ as $\frac{1}{L}$, where L is the Lipschitz constant of $\nabla l(\beta)$, set $\eta > 1$;

repeat

- 1) Start from $k = 2$, update the step size $\frac{1}{L_k}$ using the Barzilai-Borwein (BB) rule
- 2) Find the smallest nonnegative integer i_k such that with $\bar{L} = \eta^{i_k} L_k$ we have $f(p_{\bar{L}}(\beta_{k-1})) \leq f(\beta_{k-1}) - \frac{\bar{L}}{2} \|p_{\bar{L}}(\beta_{k-1}) - \beta_{k-1}\|^2$.
- 3) Set $L_k = \eta^{i_k} L_k$ and update $\beta_k = p_{L_k}(\beta_{k-1})$

until convergence

Algorithm 5 ISTA-reverse: ISTA with Lipschitz constant and reverse step size searching to solve Eq (12).

Initialize β_0 randomly, step size $\frac{1}{L_0}$ as $\frac{1}{L}$, where L is the Lipschitz constant of $\nabla l(\beta)$, set $\eta > 1$;

repeat

- 1) Find the smallest nonnegative integer i_k such that with $\bar{L} = L_0/\eta^{i_k}$ we have $f(p_{\bar{L}}(\beta_{k-1})) > f(\beta_{k-1}) - \frac{\bar{L}}{2} \|p_{\bar{L}}(\beta_{k-1}) - \beta_{k-1}\|^2$.
- 2) Set $L_k = L_0/\eta^{i_k-1}$ and update $\beta_k = p_{L_k}(\beta_{k-1})$

until convergence

The convergence of ISTA with different step-size searching methods is illustrated in Figure 3 with SCAD.

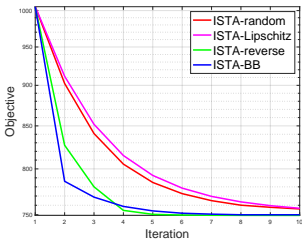


Fig. 3. Nonconvex: convergence plots.

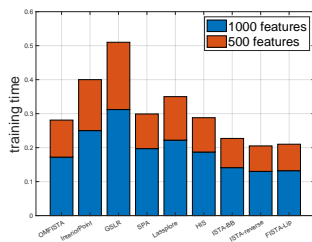


Fig. 4. \mathcal{L}_1 -norm regularized: training time comparison.

Theorem III.4. Let $\{\beta_k\}$ be the sequence generated by either Algorithm 4 or Algorithm 5, then all limit points of the sequence

$\{\beta_k\}$ are critical points of the problem in Eq (12). More specifically, after n iterations, we have

$$\min_{0 \leq k \leq n} \|\beta_{k+1} - \beta_k\|^2 \leq \frac{2(f(\beta_0) - f(\beta^*))}{nL_{\min}}, \quad (14)$$

where β^* is a limit point of the sequence $\{\beta_k\}$, L_{\min} is the minimum L among the n iterations.

Proof. With the line search criterion, we have

$$\frac{L_k}{2} \|\beta_{k+1} - \beta_k\|^2 \leq f(\beta_k) - f(\beta_{k+1}), \quad (15)$$

sum the above inequality we have

$$\frac{L_{\min}}{2} \sum_{k=0}^n \|\beta_{k+1} - \beta_k\|^2 \leq f(\beta_0) - f(\beta_{n+1}), \quad (16)$$

With $f(\beta_n) \geq f(\beta^*)$, we have

$$\frac{L_{\min}}{2} \sum_{k=0}^n \|\beta_{k+1} - \beta_k\|^2 \leq f(\beta_0) - f(\beta^*), \quad (17)$$

based on which we will obtain the desired conclusion. \square

IV. EXPERIMENTS

The empirical studies are conducted on the following 5 benchmark classification datasets, which can be found in UCI machine learning repository [2]: Wine, Specheart, Ionosphere, Madelon, and Dorothea. For the logistic regression with convex ℓ_1 -norm penalty, we compare our proposed methods with the following counterparts: OMFISTA [34], Interior-Point method [12], GSLR [15], SPA [26], Lassplore [16], HIS [25], and proximal Newton [13]. For the logistic regression with nonconvex penalties, we compare our proposed methods with GPGN [27] and HONOR [8], and the nonconvex penalty we utilize in the experiment is SCAD. In this benchmark result, the classification performance is measured by the average testing accuracy obtained with k -fold cross-validation, in our experiment, we set $k = 5$. It's known there exists a valid upper threshold for λ_{max} in logistic regression [12], when the regularization parameter is larger than that, the cardinality of the solution will be zero. Therefore λ is usually selected as a fraction proportion of λ_{max} . We choose a 10-length path for λ , where the fraction is 0.01, 0.02, 0.05, 0.07, 0.1, 0.2, 0.3, 0.5, 0.7, 0.8

TABLE III
TESTING ACCURACY OBTAINED FROM 5-FOLD CROSS-VALIDATION, WITH ℓ_1 -NORM REGULARIZATION.

	λ/λ_{max}	OMFISTA	Interior-Point	GSLR	SPA	Lassplore	HIS	Newton	ISTA-BB	ISTA-rev	FISTA-Lip
Wine	0.02	0.920	0.920	0.845	0.902	0.898	0.905	0.902	0.919	0.922	0.922
	0.1	0.907	0.911	0.841	0.882	0.886	0.899	0.901	0.911	0.913	0.909
	0.5	0.892	0.891	0.825	0.856	0.852	0.896	0.895	0.908	0.902	0.899
Specheart	0.02	0.739	0.705	0.732	0.752	0.711	0.585	0.739	0.752	0.751	0.758
	0.1	0.721	0.691	0.716	0.735	0.702	0.571	0.722	0.739	0.739	0.731
	0.5	0.685	0.672	0.701	0.698	0.681	0.559	0.659	0.688	0.697	0.701
Ionosphere	0.02	0.835	0.821	0.832	0.828	0.808	0.576	0.851	0.855	0.857	0.858
	0.1	0.802	0.809	0.825	0.815	0.789	0.573	0.809	0.818	0.825	0.822
	0.5	0.796	0.791	0.809	0.792	0.761	0.558	0.778	0.781	0.801	0.801
Madelon	0.02	0.615	0.611	0.612	0.621	0.601	0.605	0.615	0.621	0.621	0.621
	0.1	0.603	0.601	0.601	0.611	0.592	0.601	0.611	0.615	0.611	0.615
	0.5	0.592	0.582	0.583	0.591	0.579	0.589	0.592	0.601	0.601	0.602
Arrhythmia	0.02	0.611	0.591	0.601	0.608	0.599	0.589	0.608	0.615	0.615	0.615
	0.1	0.591	0.588	0.589	0.595	0.591	0.578	0.588	0.595	0.592	0.595
	0.5	0.567	0.567	0.567	0.577	0.562	0.551	0.572	0.579	0.579	0.579

TABLE IV
TESTING ACCURACY FROM CROSS-VALIDATION WITH SCAD.

	λ/λ_{max}	GPGN	HONOR	ISTA-BB	ISTA-rev
Wine	0.02	0.921	0.925	0.929	0.931
	0.1	0.912	0.912	0.915	0.917
	0.5	0.897	0.905	0.905	0.907
Specheart	0.02	0.751	0.761	0.761	0.763
	0.1	0.732	0.735	0.737	0.739
	0.5	0.691	0.711	0.711	0.711
Ionosphere	0.02	0.855	0.857	0.859	0.857
	0.1	0.825	0.827	0.829	0.831
	0.5	0.792	0.795	0.795	0.799
Madelon	0.02	0.621	0.625	0.631	0.628
	0.1	0.612	0.616	0.619	0.615
	0.5	0.585	0.597	0.601	0.603
Arrhythmia	0.02	0.595	0.609	0.611	0.618
	0.1	0.581	0.592	0.597	0.595
	0.5	0.566	0.577	0.579	0.581

respectively. All the other parameters are set as suggested in the original papers. In Table III and Table IV, we show the testing accuracy with various λ for ℓ_1 -norm and SCAD regularizer item, respectively.

We compare the efficiency of our proposed methods with other methods. We show the computation time for the training with 1000 and 500 features. The number of samples is fixed at 1000 and with $\lambda/\lambda_{max} = 0.1$. From Figure 4 we can see that in the ℓ_1 -norm regularized logistic regression problem, our proposed methods (ISTA-BB, ISTA-reverse, and FISTA-Lip) require less computation time to converge. The proximal Newton method is not included in the figure because its running time is way higher than the others, making it hard to visualize the time in the same figure. The nonconvex regularized logistic regression follows a similar path, ISTA-BB and ISTA-reverse have better performance than GPGN and HONOR in terms of less time consumption. We also conduct numerical experiments to show the scalability of our methods. We varied the number of features and the number of samples to show how our algorithms perform with the size of samples increase. The results are illustrated in Figure 5. For our proposed methods, the computation time is quite low even with a large number

of features and samples, the trends of increasing computation time with increased features/samples are quite stable, there is no explosive increase in the computation time with the rapid growth in features/samples.

We present the convergence of our proposed methods along the regularization path, using a *Ionosphere* data from the UCI repository, which is shown in Figure 6. Similarly, we also include the original ISTA and FISTA methods to show the superiority of our proposed. The figure shows that our methods have similar performance in terms of convergence, but require less iterations to converge than that of vanilla ISTA and FISTA. We also explore the sparsity of the solutions, which is shown in Figure 7. The numbers of nonzero entries obtained by our methods are pretty similar to each other, and the solutions are more sparse than the counterparts. In general, we can see that the regularization parameter λ has effects on the number of iterations to converge and the sparsity of optimal solutions. Also, we note that the proximal Newton method is not able to induce sparsity even with large λ .

V. CONCLUSIONS

In this paper, we propose new optimization frameworks to solve sparse logistic regression problem which work for both convex and nonconvex regularization terms. Experimental results on benchmark datasets with both types of regularizers demonstrate the advantages of our proposed algorithms compared to others in terms of both accuracy and efficiency.

REFERENCES

- [1] Felix Abramovich and Vadim Grinshtein. High-dimensional classification by sparse logistic regression. *IEEE Transactions on Information Theory*, 65(5):3068–3079, 2018.
- [2] Arthur Asuncion and David Newman. Uci machine learning repository, 2007.
- [3] Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- [4] Thomas E Booth. Power iteration method for the several largest eigenvalues and eigenfunctions. *Nuclear science and engineering*, 154(1):48–62, 2006.
- [5] Patrick Breheny and Jian Huang. Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection. *The annals of applied statistics*, 5(1):232, 2011.

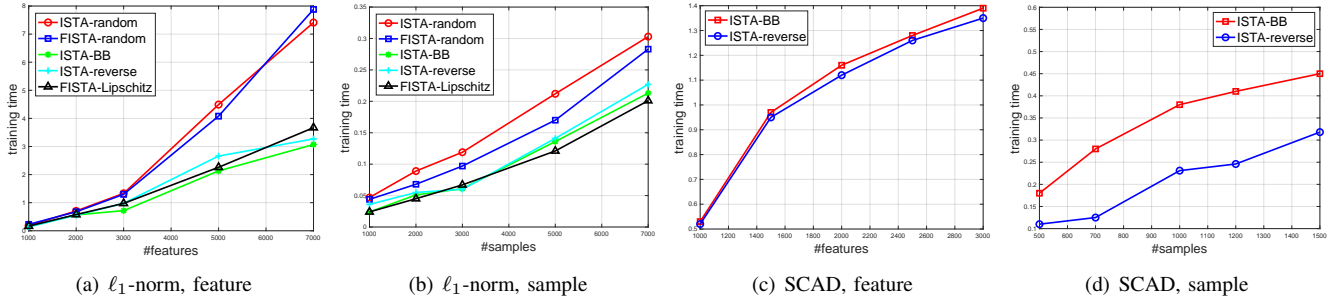


Fig. 5. Scalability study – the influence of number of features and samples on training time for each method to converge. From left to right: ℓ_1 -norm regularized with varying features, ℓ_1 -norm regularized with varying samples, SCAD regularized with varying features and SCAD regularized with varying samples.

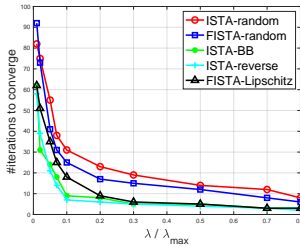


Fig. 6. \mathcal{L}_1 -norm regularized: comparison of iterations for each method to converge with varying λ .

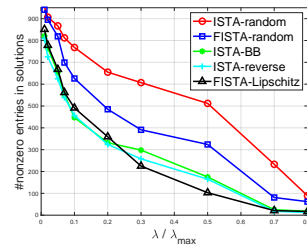


Fig. 7. \mathcal{L}_1 -norm regularized: comparison of sparsity of solutions obtained by each method with varying λ .

[6] Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360, 2001.

[7] Alexander Genkin, David D Lewis, and David Madigan. Large-scale bayesian logistic regression for text categorization. *technometrics*, 49(3):291–304, 2007.

[8] Pinghua Gong and Jieping Ye. Honor: Hybrid optimization for non-convex regularized problems. *Advances in Neural Information Processing Systems*, 28, 2015.

[9] Trevor Hastie, Robert Tibshirani, and Martin Wainwright. Statistical learning with sparsity. *Monographs on statistics and applied probability*, 143:143, 2015.

[10] Douglas M Hawkins. The problem of overfitting. *Journal of chemical information and computer sciences*, 44(1):1–12, 2004.

[11] Zehui Jia, Jieru Huang, and Zhongming Wu. An incremental aggregated proximal admm for linearly constrained nonconvex optimization with application to sparse logistic regression problems. *Journal of Computational and Applied Mathematics*, 390:113384, 2021.

[12] Kwangmoo Koh, Seung-Jean Kim, and Stephen Boyd. An interior-point method for large-scale ℓ_1 -regularized logistic regression. *Journal of Machine learning research*, 8(Jul):1519–1555, 2007.

[13] Jason D Lee, Yuekai Sun, and Michael A Saunders. Proximal newton-type methods for minimizing composite functions. *SIAM Journal on Optimization*, 24(3):1420–1443, 2014.

[14] Su-In Lee, Honglak Lee, Pieter Abbeel, and Andrew Y Ng. Efficient ℓ_1 regularized logistic regression. In *Aaai*, volume 6, pages 401–408, 2006.

[15] Alexander LeNail, Ludwig Schmidt, Johnathan Li, Tobias Ehrenberger, Karen Sachs, Stefanie Jegelka, and Ernest Fraenkel. Graph-sparse logistic regression. *arXiv preprint arXiv:1712.05510*, 2017.

[16] Jun Liu, Jianhui Chen, and Jieping Ye. Large-scale sparse logistic regression. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 547–556, 2009.

[17] Kai Liu and Hua Wang. High-order co-clustering via strictly orthogonal and symmetric ℓ_1 -norm nonnegative matrix tri-factorization. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*, 2018.

[18] Kai Liu, Hua Wang, Feiping Nie, and Hao Zhang. Learning multi-instance enriched image representations via non-greedy ratio maximization of the ℓ_1 -norm distances. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7727–7735, 2018.

[19] Po-Ling Loh and Martin J Wainwright. Regularized m-estimators with nonconvexity: Statistical and algorithmic theory for local optima. *Advances in Neural Information Processing Systems*, 26, 2013.

[20] David Madigan, Alexander Genkin, David D Lewis, and Dmitry Fradkin. Bayesian multinomial logistic regression for author identification. In *AIP conference proceedings*, volume 803, pages 509–516. American Institute of Physics, 2005.

[21] Goran Marjanovic, Magnus O Ulfarsson, and Alfred O Hero. Mist: ℓ_1 sparse linear regression with momentum. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3551–3555. IEEE, 2015.

[22] Rahul Mazumder, Jerome H Friedman, and Trevor Hastie. Sparsenet: Coordinate descent with nonconvex penalties. *Journal of the American Statistical Association*, 106(495):1125–1138, 2011.

[23] Andrew Y Ng. Feature selection, ℓ_1 vs. ℓ_2 regularization, and rotational invariance. In *Proceedings of the twenty-first international conference on Machine learning*, page 78, 2004.

[24] Mark Schmidt, Glenn Fung, and Rmer Rosales. Fast optimization methods for ℓ_1 regularization: A comparative study and two new approaches. In *European Conference on Machine Learning*, pages 286–297. Springer, 2007.

[25] Jianing Shi, Wotao Yin, Stanley Osher, and Paul Sajda. A fast hybrid algorithm for large-scale ℓ_1 -regularized logistic regression. *The Journal of Machine Learning Research*, 11:713–741, 2010.

[26] Maxime Vono, Nicolas Dobigeon, and Pierre Chainais. Sparse bayesian binary logistic regression using the split-and-augmented gibbs sampler. In *2018 IEEE 28th International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6. IEEE, 2018.

[27] Rui Wang, Naihua Xiu, and Chao Zhang. Greedy projected gradient-newton method for sparse logistic regression. *IEEE transactions on neural networks and learning systems*, 31(2):527–538, 2019.

[28] Fei Wen, Lei Chu, Peilin Liu, and Robert C Qiu. A survey on nonconvex regularization-based sparse and low-rank recovery in signal processing, statistics, and machine learning. *IEEE Access*, 6:69883–69906, 2018.

[29] Min Yuan and Yitian Xu. Feature screening strategy for non-convex sparse logistic regression with log sum penalty. *Information Sciences*, 2023.

[30] Alan L Yuille and Anand Rangarajan. The concave-convex procedure (cccp). *Advances in neural information processing systems*, 14, 2001.

[31] Cun-Hui Zhang. Nearly unbiased variable selection under minimax concave penalty. *The Annals of statistics*, 38(2):894–942, 2010.

[32] Mengyuan Zhang and Kai Liu. Enriched robust multi-view kernel subspace clustering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1993–2002, 2022.

[33] Mengyuan Zhang and Kai Liu. Multi-task learning with prior information. In *Proceedings of the 2023 SIAM International Conference on Data Mining (SDM)*, pages 586–594. SIAM, 2023.

[34] Marcelo VW Zibetti, Elias S Helou, and Daniel R Pipa. Accelerating overrelaxed and monotone fast iterative shrinkage-thresholding algorithms with line search for sparse reconstructions. *IEEE Transactions on Image Processing*, 26(7):3569–3578, 2017.