# A Fast Hybrid Algorithm for Large-Scale $\ell_1$-Regularized Logistic Regression

**Jianing Shi**                                     JS2615@COLUMBIA.EDU
*Department of Biomedical Engineering*
*Columbia University*
*New York, NY 10025, USA*

**Wotao Yin**                                    WOTAO.YIN@RICE.DU
*Department of Computational and Applied Mathematics*
*Rice University*
*Houston, TX 77005, USA*

**Stanley Osher**                                    SJO@MATH.UCLA.EDU
*Department of Mathematics*
*University of California Los Angeles*
*Los Angeles, CA 90095, USA*

**Paul Sajda**                                    PS629@COLUMBIA.EDU
*Department of Biomedical Engineering*
*Columbia University*
*New York, NY 10025, USA*

## Abstract

$\ell_1$-regularized logistic regression, also known as sparse logistic regression, is widely used in machine learning, computer vision, data mining, bioinformatics and neural signal processing. The use of $\ell_1$ regularization attributes attractive properties to the classifier, such as feature selection, robustness to noise, and as a result, classifier generality in the context of supervised learning. When a sparse logistic regression problem has large-scale data in high dimensions, it is computationally expensive to minimize the non-differentiable $\ell_1$-norm in the objective function. Motivated by recent work (Koh et al., 2007; Hale et al., 2008), we propose a novel hybrid algorithm based on combining two types of optimization iterations: one being very fast and memory friendly while the other being slower but more accurate. Called hybrid iterative shrinkage (HIS), the resulting algorithm is comprised of a fixed point continuation phase and an interior point phase. The first phase is based completely on memory efficient operations such as matrix-vector multiplications, while the second phase is based on a truncated Newton's method. Furthermore, we show that various optimization techniques, including line search and continuation, can significantly accelerate convergence. The algorithm has global convergence at a geometric rate (a Q-linear rate in optimization terminology). We present a numerical comparison with several existing algorithms, including an analysis using benchmark data from the UCI machine learning repository, and show our algorithm is the most computationally efficient without loss of accuracy.

**Keywords:** logistic regression, $\ell_1$ regularization, fixed point continuation, supervised learning, large scale

## 1. Introduction

Logistic regression is an important linear classifier in machine learning and has been widely used in computer vision (Bishop, 2007), bioinformatics (Tsuruoka et al., 2007), gene classification (Liao and Chin, 2007), and neural signal processing (Parra et al., 2005; Gerson et al., 2005; Philiastides and Sajda, 2006). $\ell_1$-regularized logistic regression or so-called sparse logistic regression (Tibshirani, 1996), where the weight vector of the classifier has a small number of nonzero values, has been shown to have attractive properties such as feature selection and robustness to noise. For supervised learning with many features but limited training samples, overfitting to the training data can be a problem in the absence of proper regularization (Vapnik, 1982, 1988). To reduce overfitting and obtain a robust classifier, one must find a sparse solution.

Minimizing or limiting the $\ell_1$-norm of an unknown variable (the weight vector in logistic regression) has long been recognized as a practical avenue for obtaining a sparse solution. The use of $\ell_1$ minimization is based on the assumption that the classifier parameters have, *a priori*, a Laplace distribution, and can be implemented using maximum-a-posteriori (MAP). The $\ell_2$-norm is a result of penalizing the mean of a Gaussian prior, while a $\ell_1$-norm models a Laplace prior, a distribution with heavier tails, and penalizes on its median. Such an assumption attributes important properties to $\ell_1$-regularized logistic regression in that it tolerates outliers and, therefore, is robust to irrelevant features and noise in the data. Since the solution is sparse, the nonzero components in the solution correspond to useful features for classification; therefore, $\ell_1$ minimization also performs feature selection (Littlestone, 1988; Ng, 1998), an important task for data mining and biomedical data analysis.

### 1.1 Logistic Regression

The basic form of logistic regression seeks a hyperplane that separates data belonging to two classes. The inputs are a set of training data $X = [x_1, \cdots, x_m]^\top \in \mathbb{R}^{m \times n}$, where each row of $X$ is a sample and samples of either class are assumed to be independently identically distributed, and class labels $b \in \mathbb{R}^m$ are of $-1/+1$ elements. A linear classifier is a hyperplane $\{x : w^\top x + v = 0\}$, where $w \in \mathbb{R}^n$ is a set of weights and $v \in \mathbb{R}$ the intercept. The conditional probability for the classifier label $b$ based on the data, according to the logistic model, takes the following form,

$$p(b_i|x_i) = \frac{\exp\left((w^T x_i + v) b_i\right)}{1 + \exp\left((w^T x_i + v) b_i\right)}, \quad i = 1, ..., m.$$

The average logistic loss function can be derived from the empirical logistic loss, computed from the negative log-likelihood of the logistic model associated with all the samples, divided by number of samples $m$,

$$l_{\mathrm{avg}}(w, v) = \frac{1}{m} \sum_{i=1}^{m} \theta\left((w^T x_i + v) b_i\right),$$

where $\theta$ is the logistic transfer function: $\theta(z) := \log(1 + \exp(-z))$. The classifier parameters $w$ and $v$ can be determined by minimizing the average logistic loss function,

$$\arg\min_{w,v} \ l_{\mathrm{avg}}(w, v).$$

Such an optimization can also be interpreted as a MAP estimate for classifier weights $w$ and intercept $v$.

### 1.2 $\ell_1$-Regularized Logistic Regression

The so-called *sparse logistic regression* has emerged as a popular linear decoder in the field of machine learning, adding the $\ell_1$-penalty on the weights $w$:

$$\arg\min_{w,v} \; l_{\text{avg}}(w,v) + \lambda \|w\|_1, \tag{1}$$

where $\lambda$ is a regularization parameter. It is well-known that $\ell_1$ minimization tends to give sparse solutions. The $\ell_1$ regularization results in logarithmic sample complexity bounds (number of training samples required to learn a function), making it an effective learner even under an exponential number of irrelevant features (Ng, 1998, 2004). Furthermore, $\ell_1$ regularization also has appealing asymptotic sample-consistency for feature selection (Zhao and Yu, 2007).

Signals arising in the natural world tend to be sparse (Parra et al., 2001). Sparsity also arises in signals represented in a certain basis, such as the wavelet transform, the Krylov subspace, etc. Exploiting sparsity in a signal is therefore a natural constraint to employ in algorithm development. An exact form of sparsity can be sought using the $\ell_0$ regularization, which explicitly penalizes the number of nonzero components,

$$\arg\min_{w,v} \; l_{\text{avg}}(w,v) + \lambda \|w\|_0. \tag{2}$$

Although theoretically attractive, problem (2) is in general NP-hard (Natarajan, 1995), requiring an exhaustive search. Due to this computational complexity, $\ell_1$ regularization has become a popular alternative, and is subtly different than $\ell_0$ regularization, in that the $\ell_1$-norm penalizes large coefficients/parameters more than small ones.

The idea of adopting the $\ell_1$ regularization for seeking sparse solutions to optimization problems has a long history. As early as the 1970's, Claerbout and Muir first proposed to use $\ell_1$ to deconvolve seismic traces (Claerbout and Muir, 1973), where a sparse reflection function was sought from bandlimited data (Taylor et al., 1979). In the 1980's, Donoho et al. quantified the ability of $\ell_1$ to recover sparse reflectivity functions (Donoho and Stark, 1989; Donoho and Logan, 1992). After the 1990s', there was a dramatic rise of applications using the sparsity-promoting property of the $\ell_1$-norm. Sparse model selection was proposed in statistics using LASSO (Tibshirani, 1996), wherein the proposed soft thresholding is related to wavelet thresholding (Donoho et al., 1995). Basis pursuit, which aims to extract sparse signal representation from overcomplete dictionaries, also underwent great development during this time (Donoho and Stark, 1989; Donoho and Logan, 1992; Chen et al., 1998; Donoho and Huo, 2001; Donoho and Elad, 2003; Donoho, 2006). In recent years, minimization of the $\ell_1$-norm has appeared as a key element in the emerging field of compressive sensing (Candés et al., 2006; Candés and Tao, 2006; Figueiredo et al., 2007; Hale et al., 2008). $\ell_1$ minimization also has far reaching impact on various applications such as portfolio optimization (Lobo et al., 2007), sparse principle component analysis (d'Aspremont et al., 2005; Zou et al., 2006), sparse interconnect wiring design (Vandenberghe et al., 1997, 1998), sparse control system design (Hassibi et al., 1999), and optimization of well-connected sparse graphs (Ghosh and Boyd, 2006). Research on total variation based image processing also shows that minimizing the $\ell_1$-norm of the intensity gradient can effectively remove random noise (Rudin et al., 1992). In the realm of machine learning, $\ell_1$ regularization exists in various forms of classifiers, including $\ell_1$-regularized logistic regression (Tibshirani, 1996), $\ell_1$-regularized probit regression (Figueiredo and Jain, 2001; Figueiredo, 2003), $\ell_1$-regularized support vector machines (Zhu et al., 2004), and $\ell_1$-regularized multinomial logistic regression (Krishnapuram et al., 2005).

### 1.3 Existing Algorithms for $\ell_1$-Regularized Logistic Regression

The $\ell_1$-regularized logistic regression problem (1) is a convex and non-differentiable problem. A solution always exists but can be non-unique. These characteristics postulate some difficulties in solving the problem. Generic methods for nondifferentiable convex optimization, such as the ellipsoid method and various sub-gradient methods (Shor, 1985; Polyak, 1987), are not designed to handle instances of (1) with data of very large scale. There has been very active development on numerical algorithms for solving the $\ell_1$-regularized logistic regression, including LASSO (Tibshirani, 1996), Gl1ce (Lokhorst, 1999), Grafting (Perkins and Theiler, 2003), GenLASSO (Roth, 2004), and SCGIS (Goodman, 2004). The IRLS-LARS (iteratively reweighted least squares least angle regression) algorithm uses a quadratic approximation for the average logistic loss function, which is consequently solved by the LARS (least angle regression) method (Efron et al., 2004; Lee et al., 2006). The BBR (Bayesian logistic regression) algorithm, described in Eyheramendy et al. (2003), Madigan et al. (2005), and Genkin et al. (2007), uses a cyclic coordinate descent method for the Bayesian logistic regression. Glmpath, a solver for $\ell_1$-regularized generalized linear models using path following methods, can also handle the logistic regression problem (Park and Hastie, 2007). MOSEK is a general purpose primal-dual interior point solver, which can solve the $\ell_1$-regularized logistic regression by formulating the dual problem, or treating it as a geometric program (Boyd et al., 2007). SMLR, algorithms for various sparse linear classifiers, can also solve sparse logistic regression (Krishnapuram et al., 2005). Recently, Koh, Kim, and Boyd proposed an interior-point method (Koh et al., 2007) for solving (1). Their algorithm takes truncated Newton steps and uses preconditioned conjugated gradient iterations. This interior-point solver is efficient and provides a highly accurate solution. The truncated Newton method has fast convergence, but forming and solving the underlying Newton systems require excessive amounts of memory for large-scale problems, making solving such large-scale problems prohibitive. A comparison of several of these different algorithms can be found in Schmidt et al. (2007).

### 1.4 Our Hybrid Algorithm

In this paper, we propose a hybrid algorithm that is comprised of two phases: the first phase is based on a new algorithm called iterative shrinkage, inspired by a fixed point continuation (FPC) (Hale et al., 2008), which is computationally fast and memory friendly; the second phase is a customized interior point method, devised by Koh et al. (2007).

Figure 1 shows a diagram of our hybrid algorithm, termed Hybrid Iterative Shrinkage (HIS) algorithm. Our algorithm requires less memory and, on mid/large-scale problems, runs faster than the interior point method. The iterative shrinkage phase only performs matrix-vector multiplications in size of $X$, as well as a very simple *shrinkage* operation (see (6) below), and therefore requires minimal memory consumption. By extending the results in Hale et al. (2008), we prove Q-linear convergence and show that the signs of $w_{\mathrm{opt}}$ (hence, the indices of nonzero elements) are obtained in a finite number of steps, typically much earlier than convergence. Based on the latter result, we propose a hybrid algorithm that is even faster and results in highly accurate solutions. Specifically, our algorithm predicts the sign changes in future shrinkage iterations, and when the signs of $w^k$ are likely to be stable, switches to the interior point method and operates on a reduced problem that is much smaller than the original. The interior point method achieves high accuracy in the solution, making our hybrid algorithm equally accurate, as will be shown in the Section 4.
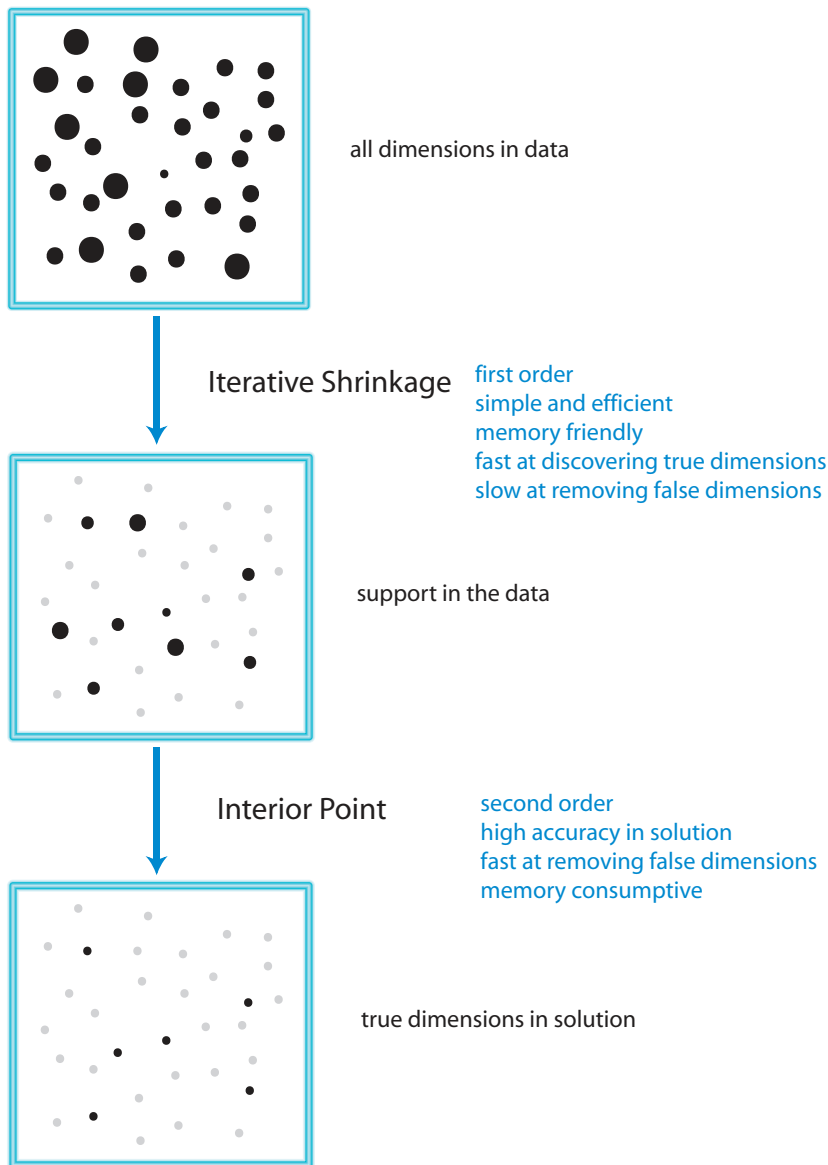
Hybrid Iterative Shrinkage (HIS)



Figure 1: A diagram of our proposed hybrid iterative shrinkage (HIS) algorithm. The HIS algorithm is comprised of two phases: the iterative shrinkage phase and the interior point phase. The iterative shrinkage is inspired by a fixed point continuation method (Hale et al., 2008), which is computationally fast and memory friendly. The interior point method is based on a second-order truncated Newton method, devised by Koh et al. (2007). Our hybrid approach takes advantage of different computational strengths of the two methods and uses them for optimal algorithm acceleration while attaining high accuracy. Black dots indicate the nonzero dimensions, gray dots indicate dimensions that are eliminated, and the size of the dots show the error that each dimension contributes to the final solution. Note that the final solution is sparse with an overall small error.

There are several novel aspects of our hybrid approach. The rationale of the hybrid approach is based on the observation that the iterative shrinkage phase reduces the algorithm to gradient projection after a finite number of iterations, which will be described in Section 3.1. We build on this observation a hybrid approach to take advantage of the two phases of the computation using two types of numerical methods. In the first phase, inspired by the FPC by Hale et al. (2008), we customize the iterative shrinkage algorithm for the sparse logistic regression, whose objective function is not quadratic. In particular, the step length in the iterative shrinkage algorithm is not constant, unlike the compressive sensing problem. Therefore, we resort to a line search strategy to avoid computing the Hessian matrix (required for finding the step length for stability). In addition, the $\ell_1$ regularization is only applied to the $w$ component and not $v$ in sparse logistic regression. This change in the model requires a different shrinkage step, as well as a careful treatment in the line search strategy.

The remainder of the paper is organized as follows. In Section 2, we present the iterative shrinkage algorithm for sparse logistic regression, and prove its global convergence and Q-linear convergence. In Section 3, we provide the rationale for the hybrid approach, together with a description of the hybrid algorithm. Numerical results are presented in Section 4. We conclude the paper in Section 5.

## 2. Sparse Logistic Regression using Iterative Shrinkage

The iterative shrinkage algorithm used in the first phase is inspired by a fixed point continuation algorithm by Hale et al. (2008).

### 2.1 Notation

For simplicity, we define $\|\cdot\| := \|\cdot\|_2$, as the Euclidean norm. The *support* of $x \in \mathbb{R}^n$ is denoted by $\mathrm{supp}(x) := \{i : x_i \neq 0\}$. We use $g$ to denote the gradient of $f$, that is, $g(x) = \nabla f(x)$, $\forall x$. For any index set $I \subseteq \{1, \ldots, n\}$ (later, we will use index sets $E$ and $L$), $|I|$ is the cardinality of $I$ and $x_I$ is defined as the sub-vector of $x$ of length $|I|$, consisting only of components $x_i$, $i \in I$. Similarly, for any vector-value mapping $h$, $h_I(x)$ denotes the sub-vector of $h(x)$ consisting of $h_i(x)$, $i \in I$.

To express the subdifferential of $\|\cdot\|_1$ we will use the signum function and multi-function (i.e., set-valued mapping). The signum function of $t \in \mathbb{R}$ is

$$\mathrm{sgn}(t) := \begin{cases} +1 & t > 0, \\ 0 & t = 0, \\ -1 & t < 0; \end{cases}$$

while the signum multi-function of $t \in \mathbb{R}$ is

$$\mathrm{SGN}(t) := \partial|t| = \begin{cases} \{+1\} & t > 0, \\ [-1, 1] & t = 0, \\ \{-1\} & t < 0, \end{cases}$$

which is also the subdifferential of $|t|$.

For $x \in \mathbb{R}^n$, we define $\mathrm{sgn}(x) \in \mathbb{R}^n$ and $\mathrm{SGN}(x) \subset \mathbb{R}^n$ component-wise as $(\mathrm{sgn}(x))_i := \mathrm{sgn}(x_i)$ and $(\mathrm{SGN}(x))_i := \mathrm{SGN}(x_i)$, $i = 1, 2, \cdots, n$, respectively. Furthermore, vector operators such as $|x|$

and $\max\{x,y\}$ are defined to operate component-wise, analogous with the definitions of sgn and SGN above. For $x,y \in \mathbb{R}^n$, let $x \odot y \in \mathbb{R}^n$ denote the component-wise product of $x$ and $y$, that is, $(x \odot y)_i = x_i y_i$. Finally, we let $X^*$ denote the set of all optimal solutions of problem (3).

## 2.2 Review of Fixed Point Continuation for $\ell_1$-minimization

A fixed-point continuation algorithm was proposed in Hale et al. (2008) as a fast algorithm for large-scale $\ell_1$-regularized convex optimization problems. The authors considered the following large-scale $\ell_1$-regularized minimization problem,

$$\min_{x \in \mathbb{R}^n} f(x) + \lambda \|x\|_1, \tag{3}$$

where $f : \mathbb{R}^n \to \mathbb{R}$ is differentiable and convex, but not necessarily strictly convex, and $\lambda > 0$. They devised a fixed-point iterative algorithm and proved its global convergence, including finite convergence for some quantities, and a Q-linear (Quotient-linear) convergence rate without assuming strict convexity of $f$ or solution uniqueness. Numerically they demonstrated Q-linear convergence in the quadratic case $f(x) = \|Ax - b\|_2^2$, where $A$ is completely dense, and applied their algorithm to $\ell_1$-regularized compressed sensing problems. As we will adopt this algorithm for solving our problem (1), we review some important and useful results here and develop some new insights in the context of $\ell_1$-regularized logistic regression.

The rationale for FPC is based on the idea of operator splitting. It is well-known in convex analysis that minimizing a function in the form of $\phi(x) = \phi_1(x) + \phi_2(x)$, where both $\phi_1$ and $\phi_2$ are convex, is equivalent to finding a zero of the subdifferential $\partial\phi(x)$, that is, seeking $x$ satisfying $0 \in T_1(x) + T_2(x)$ for $T_1 := \partial\phi_1$ and $T_2 := \partial\phi_2$. We say $(I + \tau T_1)$ is invertible if $y = x + \tau T_1(x)$ has a unique solution $x$ for any given $y$. For $\tau > 0$, if $(I + \tau T_1)$ is invertible and $T_2$ is single-valued, then

$$\begin{aligned} 0 \in T_1(x) + T_2(x) &\iff 0 \in (x + \tau T_1(x)) - (x - \tau T_2(x)) \\ &\iff (I - \tau T_2)x \in (I + \tau T_1)x \\ &\iff x = (I + \tau T_1)^{-1}(I - \tau T_2)x. \end{aligned} \tag{4}$$

This gives rise to the forward-backward splitting algorithm in the form of a fixed-point iteration,

$$x^{k+1} := (I + \tau T_1)^{-1}(I - \tau T_2)x^k. \tag{5}$$

Applying (4) to problem (1), where $\phi_1(x) := \lambda\|x\|_1$ and $\phi_2(x) := f(x)$, the authors of Hale et al. (2008) obtained the following optimality condition of $x^*$:

$$x^* \in X^* \iff 0 \in g(x^*) + \lambda\,\mathrm{SGN}(x^*) \iff x^* = (I + \tau T_1)^{-1}(I - \tau T_2)x^*,$$

where $T_2(\cdot) = g(\cdot)$, the gradient of $f(\cdot)$, and $(I + \tau T_1)^{-1}(\cdot)$ is the shrinkage operator. Therefore, the fixed-point iteration (5) for solving (3) becomes

$$x^{k+1} = s \circ h(x^k),$$

which is a composition of two mappings $s$ and $h$ from $\mathbb{R}^n$ to $\mathbb{R}^n$.

The gradient descent operator is defined as

$$h(\cdot) := I(\cdot) - \tau \nabla f(\cdot).$$

The shrinkage operator, on the other hand, can be written as

$$s(\cdot) = \text{sgn}(\cdot) \odot \max\{|\cdot| - \nu, 0\}, \tag{6}$$

where $\nu = \lambda\tau$. Shrinkage is also referred to soft-thresholding in the language of wavelet analysis:

$$(s(y))_i = \begin{cases} y_i - \nu, & y_i > \nu, \\ 0, & y_i \in [-\nu, \nu], \\ y_i + \nu, & y_i < -\nu. \end{cases}$$

In each iteration, the gradient descent step $h$ reduces $f(x)$ by moving along the negative gradient direction of $f(x^k)$ and the shrinkage step $s$ reduces the $\ell_1$-norm by "shrinking" the magnitude of each nonzero component in the input vector.

## 2.3 Iterative Shrinkage for Sparse Logistic Regression

Recall that in the sparse logistic regression problem (1), the $\ell_1$ regularization is only applied to $w$, not to $v$. Therefore, we propose a slightly different fixed point iteration. For simplicity of notation, we define column vectors $u = (w; v) \in \mathbb{R}^{n+1}$ and $c_i = (a_i; b_i) \in \mathbb{R}^{n+1}$, where $a_i = bx_i$, for $i = 1, 2, ..., m$. This reduces (1) to

$$\min_u \; l_{\text{avg}}(u) + \lambda \|u_{1:n}\|_1,$$

where $l_{\text{avg}} = \frac{1}{m}\sum_{i=1}^{m}\theta(c_i^T u)$, and $\theta$ denotes the logistic transfer function $\theta(z) = \log(1 + \exp(-z))$.

The gradient and Hessian of $l_{\text{avg}}$ with respect to $u$ is given by

$$g(u) \equiv \nabla l_{\text{avg}}(u) = \frac{1}{m}\sum_{i=1}^{m}\theta'(c_i^\top u)c_i,$$

$$H(u) \equiv \nabla^2 l_{\text{avg}}(u) = \frac{1}{m}\sum_{i=1}^{m}\theta''(c_i^\top u)c_i c_i^\top,$$

where $\theta'(z) = -(1 + e^z)^{-1}$ and $\theta''(z) = (2 + e^{-z} + e^z)^{-1}$. To guarantee convergence, we require the step length be bounded by $2(\max_u \lambda_{\max} H(u))^{-1}$.

The iterative shrinkage algorithm for sparse logistic regression is

$$\begin{aligned} u^{k+1} &= s \circ h_{1:n}(u^k), \text{ for } w \text{ component,} \\ u^{k+1} &= h_{n+1}(u^k), \text{ for } v \text{ component,} \end{aligned} \tag{7}$$

which is a composition of two mappings $h$ and $s$ from $\mathbb{R}^n$ to $\mathbb{R}^n$, where the gradient operator is

$$h(\cdot) = \cdot - \tau g(\cdot) = \cdot - \tau \nabla l_{\text{avg}}(\cdot).$$

While the authors in Hale et al. (2008) use a constant step length satisfying

$$0 < \tau < 2/\lambda_{\max}\{H_{EE}(u) : u \in \Omega\},$$

we employ line search to avoid the expensive calculation of maximum eigenvalues. We will present the convergence of the iterative shrinkage algorithm in Section 2.4. The details of the line search algorithm will be discussed in Section 2.5.

---

**Algorithm 1** Fixed-Point Continuation Algorithm

---

**Require:** $A = [c_1^\top; c_2^\top; \cdots; c_m^\top] \in \mathbb{R}^{m \times (n+1)}$, $u = (w; v) \in \mathbb{R}^{n+1}$, $f(u) = m^{-1}\phi(Au)$, task: $\min_u l_{\text{avg}}(u) + \lambda \|u\|_1$

    Initialize $u^0$

    **while** "not converge" **do**

        Armijo-like line search algorithm (Algorithm 2)

        $k = k + 1$

    **end while**

---

### 2.4 Convergence

Global convergence and finite convergence on certain quantities were proven in Hale et al. (2008) when the following conditions are met: (i) the optimal solution set $X^*$ is non-empty, (ii) $f \in C^2$ and its Hessian $H = \nabla^2 f$ is positive semi-definite in $\Omega = \{x : \|x - x^*\| \le \rho\} \subset \mathbb{R}^n$ for $\rho > 0$, and (iii) the maximum eigenvalue of $H$ is bounded on $\Omega$ by a constant $\hat{\lambda}_{\max}$ and the step length $\tau$ is uniformly less than $2/\hat{\lambda}_{\max}$. These conditions are sufficient for the forward operator $h(\cdot)$ to be non-expansive.

**Assumption 1**  *Assume problem* (1) *has an optimal solution set* $X^* \ne \emptyset$, *and there exists a set*

$$\Omega = \{x : \|x - x^*\| \le \rho\} \subset \mathbb{R}^n$$

*for some* $x^* \in X^*$ *and* $\rho > 0$ *such that* $f \in C^2(\Omega)$, $H(x) := \nabla^2 f(x) \succeq 0$ *for* $x \in \Omega$ *and*

$$\hat{\lambda}_{\max} := \sup_{x \in \Omega} \lambda_{\max}(H(x)) < \infty.$$

*For simplicity for the analysis, we choose a constant step length* $\tau$ *in the fixed-point iterations* (7): $x^{k+1} = s(x^k - \tau g(x^k))$, *where* $\nu = \tau\lambda$, *and*

$$\tau \in \left(0, 2/\hat{\lambda}_{\max}\right),$$

*which guarantees that* $h(\cdot) = I(\cdot) - \tau g(\cdot)$ *is non-expansive in* $\Omega$.

**Theorem 1**  *Under Assumption, the sequence* $\{u^k\}$ *generated by the fixed-point iterations* (7) *applied to problem* (1) *from any starting point* $x^0 \in \Omega$ *converges to some* $u^* \in U^* \cap \Omega$. *In addition, for all but finitely many iterations, we have*

$$u_i^k = u_i^* = 0, \qquad \forall i \in L = \{i : |g_i^*| < \lambda, 1 \le i \le n\}, \tag{8}$$

$$\text{sgn}(h_i(u^k)) = \text{sgn}(h_i(u^*)) = -\frac{1}{\lambda} g_i^*, \qquad \forall i \in E = \{i : |g_i^*| = \lambda, 1 \le i \le n\}, \tag{9}$$

*where as long as*

$$\omega := \min\{\nu(1 - \frac{|g_i^*|}{\lambda}) : i \in L\} > 0.$$

*The numbers of iterations not satisfying* (8) *and* (9) *do not exceed* $\|u^0 - u^*\|^2/\omega^2$ *and* $\|u^0 - u^*\|^2/\nu^2$, *respectively.*

**Proof** We sketch the proof here. First, the iteration (7) is shown to be non-expansive in $\ell_2$, that is, $\|u^k - u^*\|$ does not increase in $k$ with the assumption on the step length $\tau$. Specifically, in Assumption, the step length $\tau$ is chosen small enough to guarantee that $\|h(u^k) - h(u^*)\| \leq \|u^k - u^*\|$ (in practice, $\tau$ is determined, for example, by line search.) On the other hand, through a component-wise analysis, one can show that no matter what $\tau$ is, the shrinkage operator $s(\cdot)$ is always non-expansive, that is, $\|s(h_{1:n}(u^k)) - s(h_{1:n}(u^*))\| \leq \|h_{1:n}(u^k) - h_{1:n}(u^*)\|$. Therefore, from the definition of $u^{k+1}$ in (7), we have

$$\|u^{k+1} - u^*\| \leq \|u^k - u^*\|, \tag{10}$$

using the fact that $u^*$ is optimal if and only if $u^*$ is a fixed point with respect to (7). However, this non-expansiveness of (7) does not directly give convergence.

Next, $\{u^k\}$ is shown to have a limit point $\bar{x}$, that is, a subsequence converging to $\bar{u}$, due to the compactness of $\Omega$ and (10). (7) can be proven to converge globally to $\bar{u}$. To show this, we first get

$$\|[s \circ h_{1:n}(\bar{u}); h_{n+1}(\bar{u})] - [s \circ h_{1:n}(u^*); h_{n+1}(u^*)]\| = \|\bar{u} - u^*\|,$$

from the fact that $\bar{u}$ is a limit point, and then use this equation to show that $\bar{u} = [s \circ h_{1:n}(\bar{u}); h_{n+1}(\bar{u})]$, that is, $\bar{u}$ is a fixed point with respect to (7), and thus an optimal solution. Repeating the first step above we have $\|u^{k+1} - \bar{u}\| \leq \|u^k - \bar{u}\|$, which extends $\bar{u}$ from being the limit of a subsequence to one of the entire sequence.

Finally, to obtain the finite convergence result, we need to take a closer look at the shrinkage operator $s(\cdot)$. When (8) does not hold for some iteration $k$ at component $i$, we have $|u_i^{k+1} - u_i^*|^2 \leq |u_i^k - u_i^*|^2 - \omega^2$, and for (9), we have $|u_i^{k+1} - u_i^*|^2 \leq |u_i^k - u_i^*|^2 - \nu^2$. Obviously, there can be only a finite number of iterations $k$ in which either (8) or (9) does not hold, and such numbers do not exceed $\|u^0 - u^*\|^2/\omega^2$ and $\|u^0 - u^*\|^2/\nu^2$, respectively. ∎

A linear convergence result with a certain convergence rate can also be obtained. As long as $H_{EE}(x^*) := [H_{i,j}(x^*)]_{i,j \in E}$ has full rank *or* $f(x)$ is convex quadratic in $x$, the sequence $\{x^k\}$ converges to $x^*$ R-linearly, and $\{\|x^k\|_1 + \mu f(x^k)\}$ converges to $\|x^*\|_1 + \mu f(x^*)$ Q-linearly. Furthermore, if $H_{EE}(x^*)$ has the full rank, then R-linear convergence can be strengthened to Q-linear convergence by using the fact that the minimal eigenvalue of $H_{EE}$ at $x^*$ is strictly greater than 0.

## 2.5 Line Search

An important element of the iterative shrinkage algorithm is the step length $\tau$ at each iteration. To ensure the stability of the algorithm, we require that the step length satisfy

$$0 < \tau < 2/\lambda_{\max}\{H_{EE}(u) : u \in \Omega\}.$$

In compressive sensing, where the smooth part of the objective function is quadratic, the step length is constant. In sparse logistic regression, however, the Hessian matrix changes at each iteration. If one has to dynamically compute the step length at each iteration, this requires an expensive computation for the Hessian matrix. Therefore, we resort to an "Armijo-like" line search algorithm to avoid such a computational burden. For large-scale problems, a line search method, if used appropriately, can save tremendous CPU time and memory. Convergence of the Armijo-like line search is not proven in our paper, however heuristic results are obtained through numerical experiments.

---

**Algorithm 2** Armijo-like Line Search Algorithm

---

Compute heuristic step length $\alpha_0$
Gradient step: $u^{k-} = u^k - \alpha_0 \nabla l_{avg}(u^k)$
Shrinkage step: $u^{k+} = s_{1:n}(u^{k-}, \lambda\alpha_0)$
Obtain search direction: $p^k = u^{k+} - u^k$
**while** "j < max line search attempts" **do**
   **if** Armijo-like condition is met **then**
      Accept line search step, update $u^{k+1} = u^k + \alpha_j p^k$
   **else**
      Keep backtracking $\alpha_j = \mu\alpha_{j-1}$
   **end if**
   $j = j + 1$
**end while**

---

Let's denote the objective function for the $\ell_1$-regularized logistic regression as $\phi(u)$ for convenience:

$$\phi(u) = l_{avg}(u) + \lambda\|u_{1:n}\|_1,$$

where $l_{avg}(u) = \frac{1}{m}\sum_{i=1}^{m}\theta(c_i^T u)$ and $\theta$ is the logistic transfer function. A line search method, at each iteration, computes the step length $\alpha_k$ and the search direction $p^k$:

$$u^{k+1} = u^k + \alpha_k p^k.$$

The search direction will be described in Eqn. (12). For our sparse logistic regression, a sequence of step length candidates are identified, and a decision is made to accept one when certain conditions are satisfied. We compute a heuristic step length and gradually decrease it until a sufficient decrease condition is met.

Let's define the heuristic step length as $\alpha_0$. Ideally the choice of step length $\alpha_0$, would be a global minimizer of the smooth part of the objective function,

$$\varphi(\alpha) = l_{avg}(u^k + \alpha p^k), \quad \alpha > 0,$$

which is too expensive to evaluate, unlike the quadratic case in compressive sensing. Therefore, an inexact line search strategy is usually performed in practice to identify a step length that achieves sufficient decrease in $\varphi(\alpha)$ at minimal cost. Motivated by a similar approach in GPSR (Figueiredo et al., 2007), we compute the heuristic step length through a minimizer of the quadratic approximation for $\varphi(\alpha)$,

$$l_{avg}(u^k - \alpha\nabla l_{avg}(u^k)) \approx l_{avg}(u^k) - \alpha\nabla l_{avg}(u^k)^T \nabla l_{avg}(u^k) + 0.5\alpha^2 \nabla l_{avg}(u^k)^T H(u^k)\nabla l_{avg}(u^k).$$

Differentiating the right-hand side with respect to $\alpha$ and setting the derivative to zero, we obtain

$$\alpha_0 = \frac{\nabla l_{avg}(\bar{u}^k)^T \nabla l_{avg}(\bar{u}^k)}{\nabla l_{avg}(\bar{u}^k)^T H(\bar{u}^k)\nabla l_{avg}(\bar{u}^k)}, \tag{11}$$

where $\bar{u}_i^k = 0$, if $u_i = 0$ or $|g_i| < \lambda$ and $\bar{u}^k = u^k$, otherwise. From (11) and the strict positiveness of $\theta''$, we can see that the denominator is strictly positive as long as the gradient is nonzero. Computationally a very useful trick is not to compute the Hessian matrix directly, since we only use the vector-matrix product between the gradient vector $l_{avg}(\bar{u}^k)$ and the Hessian matrix $H(\bar{u}^k)$.

Based on the heuristic step length $\alpha_0$, we can obtain the search direction $p^k$, which is a combination of the gradient descent step and the shrinkage step:

$$u^{k-} = u^k - \alpha_0 \nabla l_{\text{avg}}(u^k),$$
$$u^{k+} = s_{1:n}(u^{k-}, \lambda\alpha_0),$$
$$p^k = u^{k+} - u^k. \tag{12}$$

It is easy to verify that $s_\nu(y)$ is the solution to the non-smooth unconstrained minimization problem $\min \frac{1}{2}\|x-y\|_2^2 + \lambda\|x\|_1$. This minimization problem is equivalent to the following smooth constrained optimization problem,

$$\min \frac{1}{2}\|x-y\|_2^2 + \nu z, \text{ subject to } (x,z) \in \Omega := \{(x,z) \mid \|x\|_1 \leq z\},$$

whose optimality condition is

$$(s(x,\nu) - x)^T (y - s(x,\nu) + \nu(z - \|s(x,\nu)\|_1) \geq 0,$$

for all $x \in \mathbb{R}^n$, $(y,z) \in \Omega$ and $\nu > 0$. Once we substitute $u - \tau g$ for $x$, $u$ for $y$, $\|u_{1:n}\|_1$ for $z$ and set $\nu = \lambda\tau$, the optimality condition becomes

$$(s_{1:n}(u-\tau g, \lambda\tau) - (u-\tau g))^T (u - s_{1:n}(u-\tau g, \lambda\tau)) + \lambda\tau(\|u_{1:n}\|_1 - \|s_{1:n}(u-\tau g, \lambda\tau)\|_1) \geq 0.$$

Using the fact $u^+ = s_{1:n}(u-\tau g, \lambda\tau)$, $p = u^+ - u$, we get

$$g^T p + \lambda(\|u_{1:n}^+\|_1 - \|u_{1:n}\|_1) \leq -p^T p/\tau,$$

which means

$$\nabla l_{\text{avg}}(u^k)^T p^k + \lambda\|u_{1:n}^{k+}\|_1 - \lambda\|u_{1:n}^k\|_1 \leq 0.$$

We then geometrically backtrack the step lengths, letting $\alpha_j = \alpha_0$, $\mu\alpha_0$, $\mu^2\alpha_0, \ldots$, until the following Armijo-like condition is satisfied:

$$\phi(u^k + \alpha_j p^k) \leq C_k + \alpha_j \Delta_k.$$

Notice that the Armijo-like condition for line search stipulates that the step length $\alpha_j$ in the search direction $p^k$ should produce a sufficient decrease of the objective function $\phi(u)$. $C_k$ is a reference value with respect to the previous objective values, while the decrease in the objective function is described as

$$\Delta_k := \nabla l_{\text{avg}}(u^k)^T p^k + \lambda\|u_{1:n}^{k+}\|_1 - \lambda\|u_{1:n}^k\|_1 \leq 0.$$

There are two types of Armijo-like conditions depending on the choice of $C_k$. One can choose $C_k = \phi(u^k)$, which makes the line search monotone. One can also derive a non-monotone line search, where $C_k$ is a convex combination of the previous value $C_{k-1}$ and the function value $\phi(u^k)$. We refer interested readers to Wen et al. (2009) for more details.

Figure 2 illustrates the computational speedup using the line search. The top panel shows the evolution of the objective function as a function of iterations. Tested on the benchmark data from the UCI repository, we see that our algorithm results in a speedup of 40 (6000 iterations without line search vs. 150 iterations with line search). The bottom panel shows the step length used in the algorithm. In the absence of the line search, we require that the step length satisfy $\tau < 2/\hat{\lambda}_{\text{max}}$. For the Armijo-like line search, we illustrate both the heuristic step length $\alpha_0$ (solid black curve) and the actual step length after backtracking (dashed red curve). Red asterisk labels the transition points on the continuation path, a concept we will discuss in the next section. Note that the step lengths can be on the order of 100 times larger for line search vs. no line search.
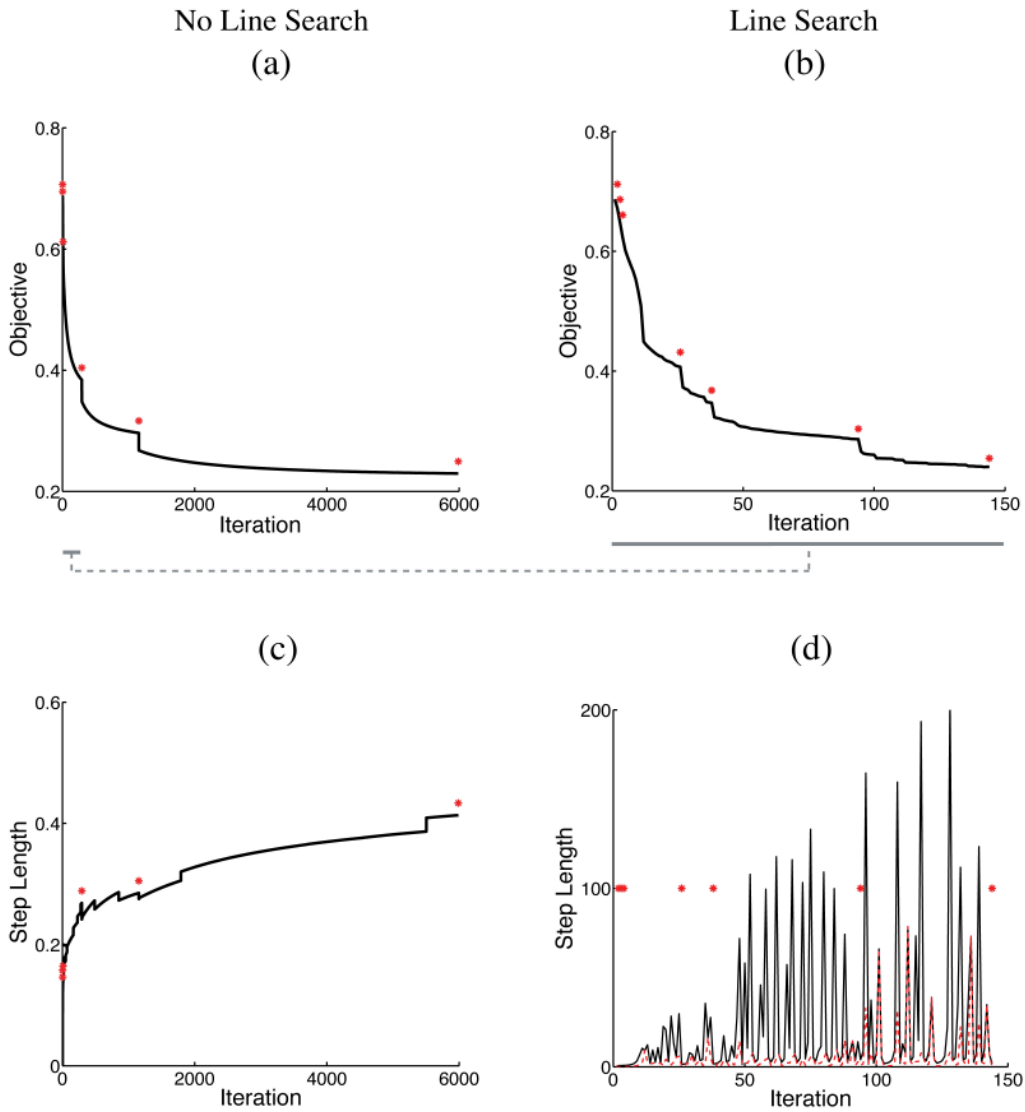
Figure 2: Illustration of the Armijo-like line search, comparing the iterative shrinkage algorithm with (right column) and without (left column) line search. (a) The objective function of the iterative shrinkage algorithm without line search, attaining convergence after 6000 iterations. (b) The objective of the iterative shrinkage algorithm with line search, converging at around 150 iterations. The gray bars under the "iteration" axes highlight the difference between the number of iterations—the gray bar in (a) represents the same number of iterations as the gray bar in (b). (c) The step length without line search is bounded by $2/\hat{\lambda}_{max}$ to ensure convergence. (d) The step length used in the Armijo-like line search, (solid black curve) heuristic step length $\alpha_0$ (Eqn. 11), (dashed red curve) actual time step after backtracking. The transition point on the continuation path is indicated in (red asterisk). Data used in this numerical experiment are the ionosphere data from the UCI machine learning data repository (http://archive.ics.uci.edu/ml/datasets/Ionosphere). Parameters used are $utol = 0.001$, $gtol = 0.01$, $\lambda_0 = 0.1$, $\lambda = 0.001$.

## 2.6 Continuation Path

A continuation strategy is adopted in our algorithm, by designing a regularization path similar to that is used in Hale et al. (2008),

$$\lambda_0 > \lambda_1 > ... > \lambda_{L-1} = \bar{\lambda}.$$

This idea is closely related to the homotopy algorithm in statistics, and has been successfully applied to the $\ell_1$-regularized quadratic case, where the fidelity term is $f(x) = \|Ax - b\|_2^2$. The rationale of using such a continuation strategy is due to a fast rate of convergence for large $\lambda$. Therefore, by taking advantage of different convergence rate for a family of regularization parameter $\lambda$, if stopped appropriately, we can speed up the convergence rate of the full path. An intriguing discussion regarding the convergence rate of fixed-point algorithm with $\lambda$ and $\omega$, the spectral properties of Hessian, was presented in Hale et al. (2008). In the case of the logistic regression, we have decided to use the geometric progression for the continuation path. We define

$$\lambda_i = \lambda_0 \beta^{i-1}, \text{for } i = 0, ..., L-1,$$

where $\lambda_0$ can be calculated based on the ultimate $\bar{\lambda}$ we are interested in and the continuation path length $L$, that is, $\lambda_0 = \bar{\lambda}/\beta^{L-1}$.

As mentioned earlier, the goal of a continuation strategy is to construct a path with different rate of convergence, with which we can speed up the whole algorithm. The solution obtained from a previous subpath associated with $\lambda_{i-1}$ is used as the initial condition for the next subpath for $\lambda_i$. Note that we design the path length $L$ and the geometric progression rate $\beta$ in such a way that the initial regularization $\lambda_0$ is fairly large, leading to a sparse solution for the initial path. Therefore, the initial condition for the whole path, considering the sparsity in solution, is a zero vector.

Another design issue regarding such a continuation strategy is we stop each subpath according to some criteria, in an endeavor to approximate the solution in the next $\lambda$ as fast as possible. This means that a strong convergence is not required in subpath's except for the final one, and we can vary the stopping criteria to "tighten" such a convergence as we proceed. The following two stopping criteria are used:

$$\frac{\|u^{k+1} - u^k\|}{\max(\|u^k\|, 1)} < utol_i,$$

$$\frac{\|\nabla l_{\text{avg}}(u^k)\|_\infty}{\lambda_i} - 1 < gtol.$$

The first stopping criterion requires that relative change in $u$ be small, while the second one is related to the optimality condition, defined in Eqn. (13). Theoretically, we would like to vary $utol_i$ to attain a seamless Q-linear convergence path. Such a choice seems to be problem dependent, and probably even data dependent in practice. It remains an important, yet difficult research topic to study the properties of different continuation strategies. We have chosen to use a geometric progression for the tolerance value, $utol_i = utol_0 * \gamma^{i-1}$, with $utol_0 = utol/\gamma^{L-1}$. In our numerical simulation, we use $utol = 10^{-4}$ and $gtol = 0.2$.

Figure 3 shows the continuation path using fixed $utol$ and a varying $utol$ following geometric progression. When we use a fixed $utol$ to ensure strong convergence each for $\lambda$ along the path, the solver spends a lot of time evolving slowly. One can see in (a) that the objective function shows a

fairly flat reduction at earlier stages of the path. Clearly by relaxing the convergence at earlier stages of the path, we can accelerate the computation, shown in (b). The choice of *utol* and *gtol* seems to be data dependent in our experience, and the result we show in (b) might be suboptimal. Further optimization of the continuation path can potentially accelerate the computation even more, which remains an open question for future research.
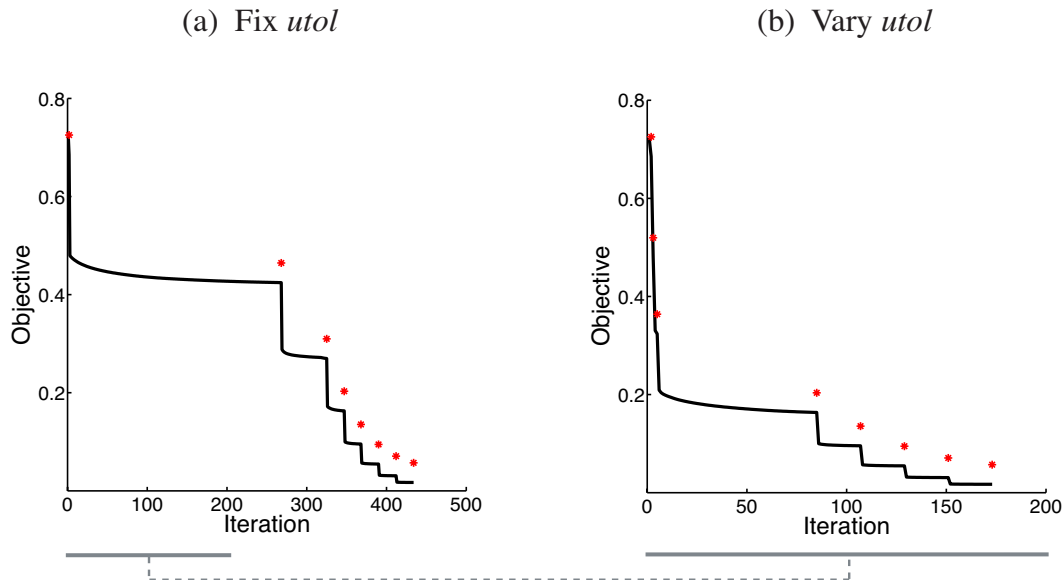
(a) Fix *utol*                    (b) Vary *utol*



Figure 3: Illustration of the continuation strategy (a) using a fixed *utol* = 0.0001 is used for the stopping criterion, (b) using a varying *utol* according to geometric progression. Note that a stronger convergence is not necessary in earlier stages on the continuation path. By using a varying *utol*, especially tightening *utol* as we move along the path, we can accelerate the fixed point continuation algorithm. Shown is the objective value (black curve) as a function of iteration, where the transition point on the regularization path is labeled in (red asterisk). Data used in this experiment has 10000 dimension and 100 samples. A continuation path of length 8, starting from 0.128 and ending at 0.001.

## 3. Hybrid Iterative Shrinkage (HIS) Algorithm

In this section we describe a hybrid approach called HIS, which uses the iterative shrinkage algorithm described previously to enforce sparsity and identify the support in the data, followed by a subspace optimization via an interior point method.

### 3.1 Why A Hybrid Approach?

The hybrid approach is based on an interesting observation for the iterative shrinkage algorithm, regarding some finite convergence properties. The optimality condition for $\min f(x) + \lambda \|x\|_1$ is the following

$$g(x) + \lambda \text{SGN}(x) \in \mathbf{0}, \tag{13}$$

which requires that $|g_i| \le \lambda$, for $i = 1, ..., n$. We define two index sets

$$L := \{i : |g_i^*| < \lambda\} \text{ and } E := \{i : |g_i^*| = \lambda\},$$

where $g^* = g(u^*)$ is constant for all $u^* \in X^*$ and $|g_i^*| \le \lambda$ for all $i$. Hence, $L \cap E = \emptyset$ and $L \cup E = \{1, \ldots, n\}$. The following holds true for all but a finite number of $k$:

$$u_i^k = u_i^* = 0, \qquad \forall i \in L,$$

$$\text{sgn}(h_i(u^k)) = \text{sgn}(h_i(u^*)) = -\frac{1}{\lambda} g_i^*, \qquad \forall i \in E.$$

Assume that the underlying problem is nondegenerate, then $L$ and $E$ equal the sets of zero and nonzero components in $x^*$. According to the above finite convergence result, the iterative shrinkage algorithm obtains $L$ and $E$, and thus the optimal support and signs of the optimal nonzero components, in a finite number of steps.

**Corollary 2** *Under Assumption 1, after a finite number of iterations, the fixed-point iteration* (7) *reduces to gradient projection iterations for minimizing* $\phi(u_E)$ *over a constraint set* $O_E$, *where*

$$\phi(u_E) := -(g_E^*)^\top u_E + f((u_E; \mathbf{0})), \text{ and}$$

$$O_E = \{u_E \in \mathbb{R}^{|E|} : -\text{sgn}(g_E^*) \odot u_E \ge 0\}.$$

*Specifically, we have* $u^{k+1} = (u_E^{k+1}; \mathbf{0})$ *in which*

$$u_E^{k+1} := P_{O_E} \left( u_E^k - \tau \nabla \phi(u_E^k) \right),$$

*where* $P_{O_E}$ *is the orthogonal projection onto* $O_E$, *and* $\nabla \phi(u_E) = -g_E^* + g_E((u_E; \mathbf{0}))$.

This corollary, see Corollary 4.6 in Hale et al. (2008), can be directly applied to sparse logistic regression. The fixed point continuation reduces to the gradient projection after a finite number of iterations. The proof of this corollary is in general true for the $u_{1:n}$, that is, the $w$ component in our problem.

Corollary 2 implies an important fact: there are two phases in the fixed point continuation algorithm. In the first phase, the number of nonzero elements in the $x$ evolve rapidly, until after a finite number of iterations, when the support (non-zero elements in a vector) is found. Precisely, it means that for all $k > K$, the nonzero entries in $u^k$ include all true nonzero entries in $u^*$ with the matched signs. However, unless $k$ is large, $u^k$ typically also has extra nonzeros. At this point, the fixed point continuation reduces to the gradient projection, starting the second phase of the algorithm. In the second phase, the zero elements in the vector stay unaltered, while the magnitude of the nonzero elements (support) keeps evolving.

The above observation is a general statement for any $f$ that is convex. Recall the quadratic case, where $f = \|y - Ax\|_2^2$, the second phase is very fast in terms of convergence rate. This is due to the quadratic function, and in an application to compressive sensing, the fixed point continuation algorithm alone has resulted in super-fast performance for large-scale problems (Hale et al., 2008). In the case of sparse logistic regression, we have a non-strictly convex $f$, the average logistic regression. This results in a fairly slow convergence rate when the algorithm reaches the second phase. In view of the continuation strategy we have, this greatly affects the speed of the last subpath, with the

regularization parameter $\bar{\lambda}$ of interest. In some sense, we have designed a continuation path that is super-fast until it reaches the second phase of the final subpath. This is not surprising given that the fixed point continuation algorithm is based on gradient descent and shrinkage operator. We envision that by switching to a Newton's method, we can accelerate the second phase.

Based on this intuition, we are now in a position to describe a hybrid algorithm: a fixed point continuation plus an interior point truncated Newton method. For the latter part we resort to the customized interior point in Koh et al. (2007). We modified the source code of the *l1logreg* software (written in C), and built an interface to our MATLAB code. This hybrid approach, based on our observation of the two phases, enables us to attain a good balance of speed and accuracy.

## 3.2 Interior Point Phase

The second phase of our HIS algorithm used an interior point method developed by Koh et al. (2007). We directly used a well-developed software package *l1logreg*[1] and modified the source code to build an interface to MATLAB. We review some key points for the interior point method here.

In Koh et al. (2007), the authors overcome the difficulty of non-differentiability of the objective function by transforming the original problem into an equivalent one with linear inequality constraints,

$$
\begin{aligned}
\min \quad & \frac{1}{m}\sum_{i=1}^{m} l_{\mathrm{avg}}(w^T a_i + v b_i) + \lambda \mathbf{1}^T u \\
\text{s.t.} \quad & -u_i \le w_i \le u_i, \ \ i = 1,...,n.
\end{aligned}
$$

A logarithmic barrier function, smooth and convex, is further constructed for the bound constraints,

$$
\rho(w,u) = -\sum_{i=1}^{n} \log(u_i + w_i) - \sum_{i=1}^{n} \log(u_i - w_i),
$$

defined on the domain $\{(w,u) \in \mathbb{R}^n \times \mathbb{R}^n \,|\, |w_i| < u_i, i = 1,...,n\}$. The following optimization problem can be obtained by augmenting the logarithmic barrier,

$$
\psi_t(v,w,u) = t l_{\mathrm{avg}}(v,w) + t\lambda \mathbf{1}^T u + \rho(w,u),
$$

where $t > 0$. The resulting objective function is smooth, strictly convex and bounded below, and therefore has a unique minimizer $(v^*(t), w^*(t), u^*(t))$. This defines a curve in $\mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^n$, parameterized by $t$, called the central path. The optimal solution is also shown to be dual feasible. In addition, $(v^*(t), w^*(t))$ is $2n/t$-suboptimal.

As a primal interior-point method, the authors computed a sequence of points on the central path, for an increasing sequence of values of $t$, and minimized $\psi_t(v,w,u)$ for each $t$ using a truncated Newton's method. The interior point method was customized by the authors in several ways: 1) the dual feasible point and the associated duality gap was computed in a cheap fashion, 2) the central path parameter $t$ was updated to achieve a robust convergence when combined with the preconditioned conjugate gradient (PCG) algorithm, 3) an option for solving the Newton's system was given for problems of different scales, where small and medium dense problems were solved by direct methods (Cholesky factorization), while large problems were solved using iterative methods (conjugate gradients). Interested readers are referred to Koh et al. (2007) for more details.

---

1. Software can be downloaded at `http://www.stanford.edu/~boyd/papers/l1_logistic_reg.html`.

## 3.3 The Hybrid Algorithm

The hybrid algorithm leverages the computational strengths of both the iterative shrinkage solver and the interior point solver.

---
**Algorithm 3** Hybrid Iterative Shrinkage (**HIS**) Algorithm

---
**Require:** $A = [c_1^\top; c_2^\top; \cdots; c_m^\top] \in \mathbb{R}^{m \times n+1}$, $u = (w; v) \in \mathbb{R}^{n+1}$, $f(u) = m^{-1}\phi(Au)$
  task: $\min_u l_{\mathrm{avg}}(u) + \lambda\|u\|_1$
  Initialize $u^0$
  **PHASE 1 : ITERATIVE SHRINKAGE**
  Select $\lambda_0$ and $utol_0$
  **while** "not converge" **do**
    **if** "the last continuation path", $i == (L-1)$ and "transition condition" **then**
      "transit into PHASE 2"
    **else**
      Update $\lambda_i = \lambda_{i-1}\beta$, $utol_i = utol_{i-1}\gamma$
      Compute heuristic step length $\alpha_0$
      Gradient descent step: $u^{k-} = u^k - \alpha_0 \nabla l_{\mathrm{avg}}(u^k)$
      Shrinkage step: $u^{k+} = s_{1:n}(u^{k-}, \lambda\alpha_0)$
      Obtain line search direction: $p^k = u^{k+} - u^k$
      **while** "j < max line search attempts" **do**
        **if** Armijo-like condition is met **then**
          Accept line search step, update $u^{k+1} = u^k + \alpha_j p^k$
        **else**
          Keep backtracking $\alpha_j = \mu\alpha_{j-1}$
        **end if**
        $j = j+1$
      **end while**
    **end if**
  **end while**
  **PHASE 2 : INTERIOR POINT**
  Initialize $\tilde{w} = w_{nonzero}$ get subproblem $\min \psi_t(v, \tilde{w}, u)$
  **while** "not converged" $\eta > \epsilon$ **do**
    Solve the Newton system : $\nabla^2 \psi_t^k(v, \tilde{w}, u)[\Delta v, \Delta\tilde{w}, \Delta u] = -\nabla\psi_t^k(v, \tilde{w}, u)$
    Backtracking line search : find the smallest integer $j \geq 0$ that satisfies
    $\psi_t^k(v + \alpha_j\Delta v, \tilde{w} + \alpha_j\Delta\tilde{w}, u + \alpha_j\Delta u) \leq \psi_t^k(v, \tilde{w}, u) + c\alpha_j\nabla\psi_t^k(v, \tilde{w}, u)^T[\Delta v, \Delta\tilde{w}, \Delta u]$
    Update $\psi_t^{k+1}(v, \tilde{w}, u) = \psi_t^k(v, \tilde{w}, u) + \alpha_j(\Delta v, \Delta\tilde{w}, \Delta u)$
    Check dual feasibility
    Evaluate duality gap $\eta$
    $k = k+1$
  **end while**

---

In the first phase, we use the iterative shrinkage solver, due to its computational efficiency and memory friendliness. It is especially beneficial to have a memory friendly solver for the initial phase when one is dealing with large-scale data sets. Recall that we use a continuation strategy for the iterative shrinkage phase, where a sequence of $\lambda$'s is used along a regularization path. In the

last subpath where $\lambda$ is the desired one, we transit to the interior point when the true support of the vector is found. The corollary in Section 3.1 states that iterative shrinkage recovers the true support in a finite number of steps. In addition, iterative shrinkage obtains all true nonzero components long before the true support is obtained. Therefore, as long as the iterative shrinkage seems to stagnate, which can be observed when the objective function evolves very slowly, it is highly likely that all true nonzero components are obtained. This indicates that the algorithm is ready for switching to the interior point.

In practice, we require the following transition condition,

$$\frac{\|u^{k+1} - u^k\|}{\max(\|u^k\|, 1)} < utol_t,$$

and extract the nonzero components in $w$ as the input to the interior point solver. By doing so, we reduce the problem to a subproblem where the dimension is much smaller, and solve the subproblem using the interior point method.

The resulting hybrid algorithm achieves high computational speed while attaining the same numerical accuracy as the interior point method, as demonstrated with empirical results in the next section.

## 4. Numerical Results

In this section we present numerical results, on a variety of data sets, to demonstrate the benefits of our hybrid framework in terms of computational efficiency and accuracy.

### 4.1 Benchmark

We carried out a numerical comparison of the HIS algorithm with several existing algorithms in literature for $\ell_1$-regularized logistic regression. Inspired by a comparison study on this topic by Schmidt et al. (2007),[2] we compared our algorithm with 10 algorithms, including a generalized version of Gauss-Seidel, Shooting, Grafting, Sub-Gradient, epsL1, Log-Barrier, Log-Norm, SmoothL1, EM, ProjectionL1 and Interior-Point method. In the numerical study, we replaced the interior point solver by the one written by Koh et al. (2007). Benchmark data were taken from the publicly available UCI machine learning repository.[3] We used 10 data sets of small to median size (internetad1, arrhythmia, glass, horsecolic, indiandiabetes, internetad2, ionosphere, madelon, pageblock, spambase, spectheart, wine).

All of the methods were run until the same convergence criteria was met, where appropriate, for instance the step length, change in function value, negative directional derivative, optimality condition, convergence tolerance is less than $10^{-6}$. We treated each algorithm solver as a black box and evaluated both the computation time and the sparsity (measured by cardinality of solution). We set an upper limit of 250 iterations, meaning we stop the solver when the number of iteration exceeds 250. Since different algorithm has different speed for each iterate (usually a Newton step is more expensive than a gradient descent step), we think the computation time is a more appropriate evaluation criterion than number of iterations. The ability of the algorithm to find a sparse solution, measured by the cardinality, was also evaluated in this process.

---

2. Source code is available at `http://www.cs.wisc.edu/~gfung/GeneralL1`.
3. UCI machine learning repository is at `http://www.ics.uci.edu/~mlearn/MLRepository.html`.

Figure 4 shows the benchmark result using data from the UCI machine learning repository. All numerical results shown are averaged over a regularization path. The parameters for the regularization path are calculated according to each data set, where the maximal regularization parameter is calculated as follows:

$$\lambda_{max} = \frac{1}{m} \left\| \frac{m_-}{m} \sum_{b_i=+1} a_i + \frac{m_+}{m} \sum_{b_i=-1} a_i \right\|_\infty, \tag{14}$$

where $m_-$ is the number of training samples with label $-1$ and $m_+$ is the number of training samples with label $+1$ (Koh et al., 2007). $\lambda_{max}$ is an upper bound for the useful range of regularization parameter. When $\lambda \geq \lambda_{max}$, the cardinality of the solution will be zero. In this case, we test a regularization path of length 10, that is, $\lambda_{max}, 0.9\lambda_{max}, 0.8\lambda_{max} \dots 0.1\lambda_{max}$. Among all the numerical solvers, our HIS algorithm is the most efficient. HIS achieves comparable cardinality in the solution, compared to the interior point solver.

We also evaluated the accuracy of the solution by looking at the classification performance using Kfold cross-validation. Table 1 summaries the accuracy of the solution using the HIS algorithm, compared to the interior point (IP) algorithm. Clearly, HIS algorithm achieves comparable accuracy compared to IP, an algorithm that is recognized for its high accuracy.
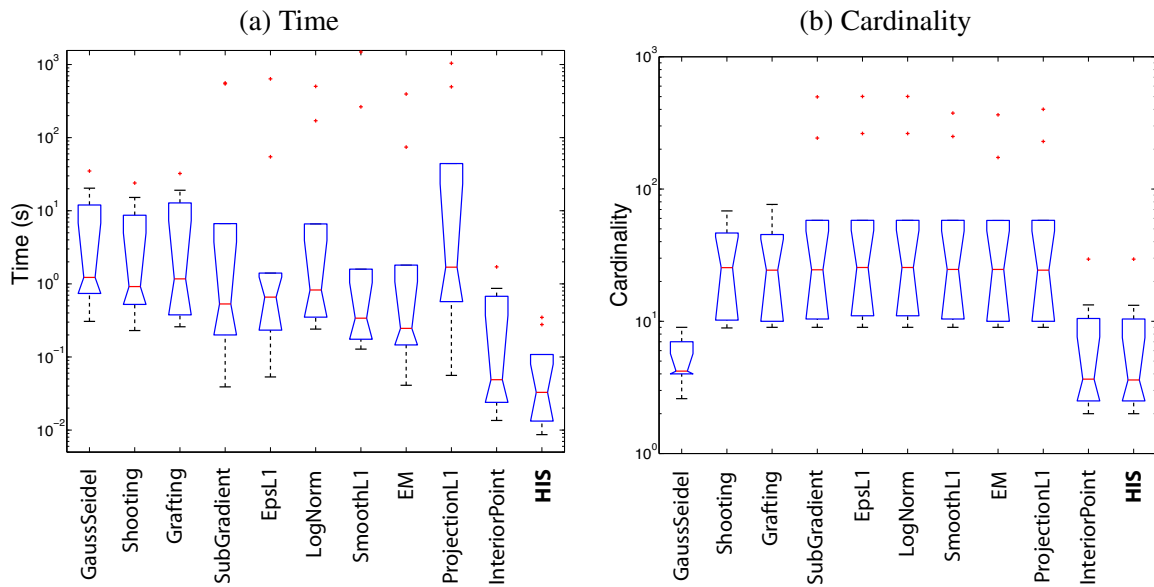


Figure 4: Comparison of our hybrid iterative shrinkage (HIS) method with several other existing methods in literature. Benchmark data were taken from the UCI machine learning repository, including 10 publicly available data sets. (a) Distribution of computation time across 10 data sets, (b) Distribution of cardinality for the solution across 10 data sets, averaged over a regularization path.

## 4.2 Scaling Result

Numerical experiments were carried out to study how our algorithm scales with the problem size. For the sake of generality, we used simulated data whose dimension ranges from 64 to 131072. The

Accuracy Comparison
($Az \in [0.5, 1.0]$)

| dataname | accuracy(HIS) | accuarcy(IP) |
|---|---|---|
| arrhythmia | 0.7363 | 0.7363 |
| glass | 0.6102 | 0.6102 |
| horsecolic | 0.5252 | 0.5252 |
| ionosphere | 0.5756 | 0.5756 |
| madelon | 0.6254 | 0.6254 |
| spectheart | 0.5350 | 0.5350 |
| wine | 0.6102 | 0.6102 |
| internetad | 0.8486 | 0.8486 |

Table 1: Comparison of solution accuracy for our hybrid iterative shrinkage (HIS) algorithm and the interior point (IP) algorithm. Accuracy of the solution was measured by Az value, resulted from Kfold cross-validation, where Kfold is 10. A regularization path of varying $\lambda$ were computed to determine the maximum generalized Az value. The data sets were taken from the UCI machine learning repository.

data is drawn from a Normal distribution, where the mean of the distribution is shifted by a small amount for each class (0.1 for samples with label 1, and $-0.1$ for samples with label $-1$). The number of samples is the same for both classes and chosen to be smaller than the dimension of the data. Experiments for each dimension were carried out on 100 different sets of random data. We compared the mean and variances of the computation time, and compared our HIS algorithm to the IP algorithm.

Table 2 summarizes the computational speed for the HIS algorithm and the IP algorithm. It is noteworthy that the HIS algorithm improves the efficiency of computation, while maintaining comparable accuracy to the IP algorithm. Figure 5 plots the computation result as a function of dimension for better illustration. In (a) one can clearly see the speedup we gain from the HIS algorithm (red), compared to the IP algorithm (blue). We also show the solution quality in (b), where the weights we get from both solvers, is comparable.

### 4.3 Regularization Parameter

In general, the regularization parameter $\lambda$ affects the number of iterations to converge for any solver. As $\lambda$ becomes smaller, the cardinality of the solution increases, and the computation time needed for convergence also increases. Therefore when one seeks a solution with less sparsity (small $\lambda$), it is more computationally expensive.

In practice, when one carries out classification on a set of data, the optimal regularization parameter is often unknown. Speaking of optimality, we refer to a regularization parameter that results in the best classification result evaluated using Kfold cross-validation. One would run the algorithm along a regularization path, $\lambda_{max}, ..., \lambda_{min}$, where $\lambda_{max}$ is computed by Eqn. (14) and where $\lambda_{min}$ is supplied by the user.

Figure 6 shows the evolution of solution along the regularization path, using a small data set (ionosphere) from the UCI machine learning repository. This explores sparsity of different degrees

Speed Comparison
(in *second*)

| dimension | mean(HIS) | std(HIS) | mean(IP) | std(IP) |
|---|---|---|---|---|
| 64 | 0.0026 | 0.00069 | 0.0043 | 0.00057 |
| 128 | 0.0025 | 0.00058 | 0.0049 | 0.00037 |
| 256 | 0.0026 | 0.00075 | 0.0078 | 0.00052 |
| 512 | 0.0024 | 0.00059 | 0.018 | 0.0017 |
| 1024 | 0.0023 | 0.00056 | 0.029 | 0.0023 |
| 2048 | 0.0026 | 0.00064 | 0.054 | 0.0026 |
| 4096 | 0.0028 | 0.00057 | 0.098 | 0.0050 |
| 8192 | 0.0030 | 0.00059 | 0.19 | 0.0076 |
| 16384 | 0.0033 | 0.00055 | 0.40 | 0.018 |
| 32768 | 0.0038 | 0.00055 | 0.89 | 0.037 |
| 65536 | 0.0049 | 0.00054 | 2.01 | 0.096 |
| 131072 | 0.0077 | 0.00056 | 4.49 | 0.24 |

Table 2: Speed comparison of the HIS algorithm with the IP algorithm, based on simulated random benchmark data. Shown here is the computation speed as a function of dimension. Data used here are generated by sampling from two Gaussian distributions. Note that in the simulation, the continuation path used in the iterative shrinkage may or may not be optimal, which means that the speed profile for the HIS algorithm can be essentially accelerated even more.
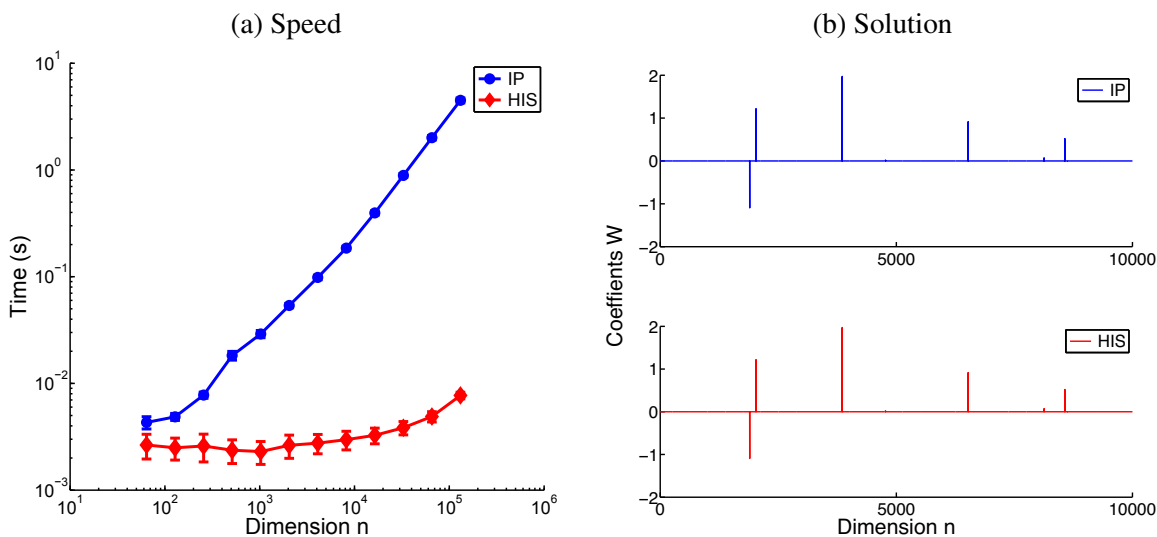


Figure 5: Comparison for the random benchmark data, between the HIS algorithm and the IP algorithm. (a) Speed profile for these two approaches: (blue curve) shows the speed profile for the IP algorithm, and (red curve) shows the speed profile for the HIS algorithm as a function of the data dimension. (b) An example of the solutions using the IP algorithm (blue) and the HIS algorithm (red).
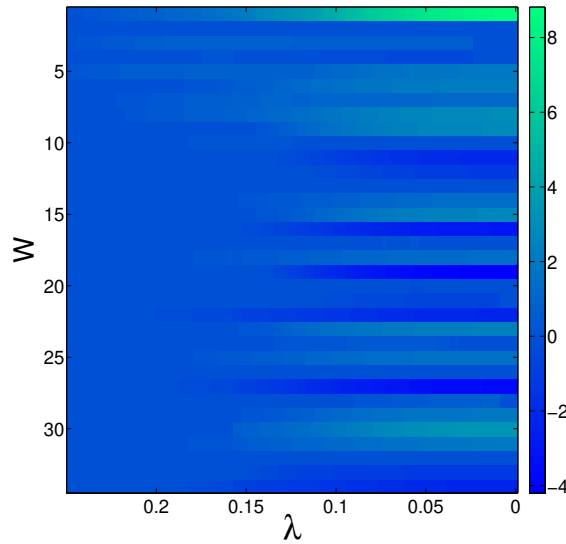
Figure 6: Solution $w$ evolves along a regularization path, following a geometric progression from $10^{-1}$ to $10^{-4}$. Data is ionosphere from UCI machine learning repository. As the $\lambda$ becomes smaller, the cardinality of the solution goes up.

in the solution, and one can determine the optimal sparsity for the data. This is an attractive property of this model, where one can search in the feature space the most informative features about discrimination.

We illustrate the effect of the regularization parameter using real data of large scale. The data concerns a two alternative force choice task for face versus car discrimination. We used a spiking neuron model of primary visual cortex to map the input into cortical space, and decoded the resulting spike trains using sparse logistic regression (Shi et al., 2009). The data has 40960 dimensions and 360 samples for each of the two classes. Kfold cross-validation was used to evaluate the classification performance, where the number of Kfolds is 10 in our simulation.

The speedup of the HIS algorithm compared to the IP algorithm is shown in Figure 7(a), where blue indicates the computation time of the IP algorithm, and red shows the HIS algorithm. The HIS algorithm results in a significant speedup over the IP algorithm, without loss of accuracy. Note that there is an issue of model selection when we apply sparse logistic regression model to the data, in a sense there exists an optimal level of sparsity that achieves the best classification result. We ran the model with a sequence of regularization parameters, which resulted in classification result (evaluated by Az value from Kfold cross-validation). Figure 7(b) illustrates the classification result as a function of the cardinality of the solution. One can see the bell shape in the curve, which provides a route to select the optimal sparsity for the solution.
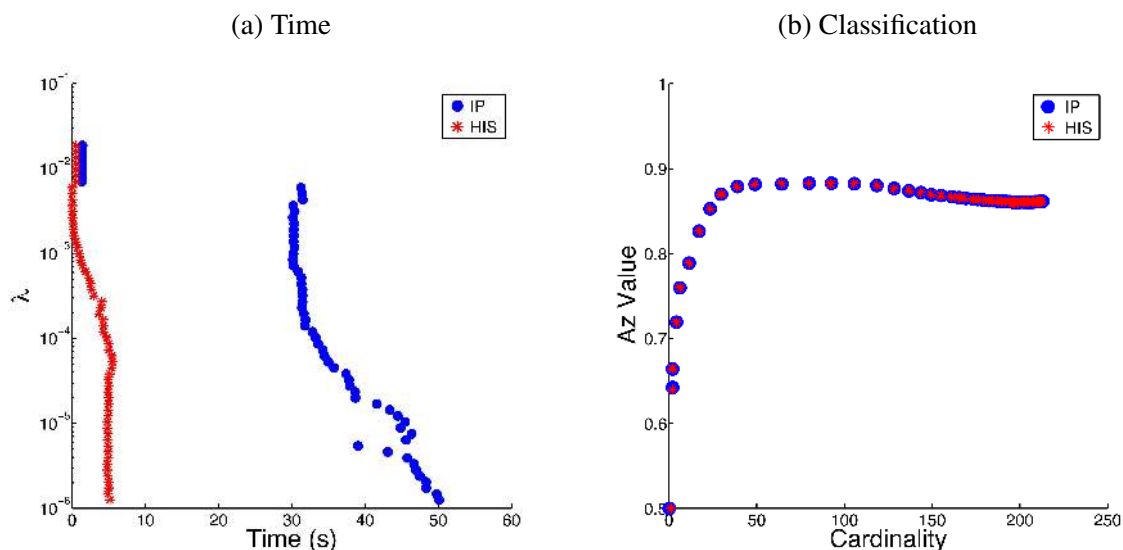
(a) Time  (b) Classification



Figure 7: An example using real data of large scale, $n = 40960$, $m = 360$. (a) Computation time along such a regularization path, where the smaller $\lambda$ requires more computation time. Note that the simulation is carried out for each $\lambda$ separately. (b) Classification performance derived from ROC analysis based on Kfold cross-validation. Data used in this simulation are neural data for a visual discrimination task (Shi et al., 2009).

### 4.4 Data Sets with Large Dimensions and Samples

We applied the HIS algorithm to some examples of real-world data that have both large dimensions $n$ and samples $m$. In this case, we considered text classification using the binary rcv1 data[4] (Lewis et al., 2004), and real-sim data.[5]

We ran the simulation on an Apple Mac Pro with two 3 GHz Quad-Core Intel processors, and 8 GB of memory. The timing of the simulation was calculated within the Matlab interface. All the operations were optimized for sparse matrix computation. Table 3 summarizes the numerical results. For both examples of text classification, we observed a speedup using the HIS algorithm while attaining the same numerical accuracy, compared with the IP algorithm. The regularization parameter does affect the computational efficiency, as we have observed in the previous section.

## 5. Conclusion

We have presented in this paper a computationally efficient algorithm for the $\ell_1$-regularized logistic regression, also called the sparse logistic regression. The sparse logistic regression is a widely used model for binary classification in supervised learning. The $\ell_1$ regularization leads to sparsity in the solution, making it a robust classifier for data whose dimensions are larger than the number of

---

4. Binary rcv1 data is available at `http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html#rcv1.binary`.

5. Real-sim data is available at `http://www.cs.umass.edu/~mccallum/code-data.html`.

Text Classification Application
(in *second*)

| $\lambda$ | rcv1 $n = 20242$ $m = 47236$ nonzero = 1498952 | | real-sim $n = 72309$ $m = 20958$ nonzero = 3709083 | |
|---|---|---|---|---|
| | Time(HIS) | Time(IP) | Time(HIS) | Time(IP) |
| $10^{-1}$ | 0.11 | 1.93 | 0.62 | 6.61 |
| $10^{-2}$ | 0.27 | 1.93 | 0.62 | 6.61 |
| $10^{-3}$ | 2.08 | 8.20 | 5.50 | 18.45 |
| $10^{-4}$ | 5.80 | 8.66 | 13.12 | 19.36 |

Table 3: Illustration of performance on text classification, where both the dimensions *n* and samples *m* are large-scale. We compare the computational efficiency of the HIS and IP algorithms. In both cases, the solution accuracy is the same.

samples. Sparsity also provides an attractive avenue for feature selection, useful for various data mining tasks.

Solving the large-scale sparse logistic regression usually requires expensive computational resources, depending on the specific solver, memory and/or CPU time. The interior point method is so far the most efficient solver in the literature, but requires expensive memory consumption. We have presented the HIS algorithm, which couples a fast shrinkage method and a slower but more accurate interior point method. The iterative shrinkage algorithm has global convergence with a Q-linear rate. Various techniques such as line search and continuation strategy are used to accelerate the computation. The shrinkage solver only involves the gradient descent and the shrinkage operator, both of which are first-order. Based solely on efficient memory operations such as matrix-vector multiplication, the shrinkage solver serves as the first phase for the algorithm. This reduces the problem to a subspace whose dimension is smaller than the original problem. The HIS algorithm then transits into the second phase, using a more accurate interior point solver. We numerically compare the HIS algorithm with other popular algorithms in the literature, using benchmark data from the UCI machine learning repository. We show that the HIS algorithm is the most computationally efficient, while maintaining high accuracy. The HIS algorithm also scales very well with dimension of the problem, making it attractive for solving large-scale problems.

There are several ways to extend the HIS algorithm. One is to extend it beyond binary classification, allowing for multiple classes (Krishnapuram and Hartemink, 2005). The other is to further improve the regularization path. When applying the HIS algorithm, one will usually explore a range of sparsity by constructing a regularization path ($\lambda_{max}$, $\lambda_1$, ..., $\lambda_{min}$). Usually the smaller the $\lambda$, the more expensive it is to employ the shrinkage algorithm. One can accelerate the computation using the Bregman regularization, inspired by Yin et al. (2008). The Bregman iterative algorithm essentially boosts the solution by solving a sequence of optimizations, resulting in a different regularization path. Bregman has also been shown to improve solution quality in the presence of noise (Burger et al., 2006; Shi and Osher, 2008; Osher et al., 2010). We will discuss such a regularization path in a future paper.

## Acknowledgments

## References

C.M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2007.

S. Boyd, S.-J. Kim, L. Vandenberghe, and A. Hassibi. A tutorial on geometric programming. *Optimization and Engineering*, 8(1):67–127, 2007.

M. Burger, G. Gilboa, S. Osher, and J. Xu. Nonlinear inverse scale space methods. *Communications in Mathematical Sciences*, 4:179–212, 2006.

E.J. Candés and T. Tao. Near optimal signal recovery from random projections: Universal encoding strategies? *IEEE Trans. Inform. Theory*, 52(2):5406–5425, 2006.

E.J. Candés, J. Romberg, and T. Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inform. Theory*, 52(2):489–509, 2006.

S. Chen, D.L. Donoho, and M.A. Saunders. Atomic decomposition by basis pursuit. *SIAM J. Scientific Computing*, 20:33–61, 1998.

J.F. Claerbout and F. Muir. Robust modeling with erratic data. *Geophysics*, 38(5):826–844, 1973.

A. d'Aspremont, L. El Ghaoui, M. Jordan, and G. Lanckriet. A direct formulation for sparse pca using semidefinite programming. In *Advances in Neural Information Processing Systems*, pages 41–48. MIT Press, 2005.

D.L. Donoho. Compressed sensing. *IEEE Trans. Inform. Theory*, 52:1289–1306, 2006.

D.L. Donoho and M. Elad. Optimally sparse representations in general nonorthogonal dictionaries by $\ell_1$ minimization. *Proc. Nat'l Academy of Science*, 100(5):2197–2202, 2003.

D.L. Donoho and X. Huo. Uncertainty principles and ideal atomic decomposition. *IEEE. Trans. Inform. Theory*, 48(9):2845–2862, 2001.

D.L. Donoho and B.F. Logan. Signal recovery and the large sieve. *SIAM J. Appl. Math.*, 52(2): 577–591, 1992.

D.L. Donoho and P.B. Stark. Uncertainty principle and signal recovery. *SIAM J. Appl. Math.*, 49(3): 906–931, 1989.

D.L. Donoho, I. Johnstone, G. Kerkyacharian, and D. Picard. Wavelet shrinkage: Asymptopia? *J. Roy. Stat. Soc. B*, 57(2):301–337, 1995.

B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Annals of Statistics*, 32(2):407–499, 2004.

S. Eyheramendy, A. Genkin, W. Ju, D. Lewis, and D. Madigan. Sparse bayesian classifiers for text categorization. Technical report, J. Intelligence Community Research and Development, 2003.

M. Figueiredo. Adaptive sparseness for supervised learning. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 25:1150–1159, 2003.

M. Figueiredo and A. Jain. Bayesian learning of sparse classifiers. In *IEEE Conf. Computer Vision and Pattern Recognition*, pages 35–41, 2001.

M. Figueiredo, R. Nowak, and S. Wright. Gradient projection for sparse reconstruction: application to compressed sensing and other inverse problems. *IEEE J. Selected Topics in Signal Processing: Special Issue on Convex Optimization Methods for Signal Processing*, 1(4):586–598, 2007.

A. Genkin, D.D. Lewis, and D. Madigan. Large-scale bayesian logistic regression for text categorization. *Technometrics*, 49(3):291–304, 2007.

A.D. Gerson, L.C. Parra, and P. Sajda. Cortical origins of response time variability during rapid discrimination of visual objects. *Neuroimage*, 28(2):342–353, 2005.

A. Ghosh and S. Boyd. Growing well-connected graphs. In *45th IEEE Conference on Decision and Control*, pages 6605–6611, 2006.

J. Goodman. Exponential priors for maximum entropy models. In *Proceedings of the Annual Meetings of the Association for Computational Linguistics*, 2004.

E. Hale, W. Yin, and Y. Zhang. Fixed-point continuation for $\ell_1$-minimization: methodology and convergence. *SIAM J. Optimization*, 19(3):1107–1130, 2008.

A. Hassibi, J. How, and S. Boyd. Low-authority controller design via convex optimization. *AIAA Journal of Guidance, Control and Dynamics*, 22(6):862–872, 1999.

K. Koh, S.-J. Kim, and S. Boyd. An interior-point method for large-scale $\ell_1$-regularized logistic regression. *J. Machine Learning Research*, 8:1519–1555, 2007.

B. Krishnapuram and A. Hartemink. Sparse multinomial logistic regression: fast algorithms and generalization bounds. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 27(6):957–968, 2005.

B. Krishnapuram, L. Carin, and M. Figueiredo. Sparse multinomial logistic regression: fast algorithms and generalization bounds. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 27(6): 957–968, 2005.

S. Lee, H. Lee, P. Abbeel, and A. Ng. Efficient $\ell_1$-regularized logistic regression. In *21th National Conference on Artificial Intelligence (AAAI)*, 2006.

D.D. Lewis, Y. Yang, T.G. Rose, and F. Li. Rcv1: A new benchmark collection for text categorization research. *J. Machine Learning Research*, 5:361–397, 2004.

J.G. Liao and K.V. Chin. Logistic regression for disease classification using microarray data: model selection in a large p and small n case. *Bioinformatics*, 23(15):1945–51, 2007.

N. Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, 2:285–318, 1988.

M. Lobo, M. Fazel, and S. Boyd. Portfolio optimization with linear and fixed transaction costs. *Annals of Operations Research*, 152(1):376–394, 2007.

J. Lokhorst. The lasso and generalised linear models. Technical report, Honors Project, Department of Statistics, University of Adelaide, South Australia, Australia, 1999.

D. Madigan, A. Genkin, D. Lewis, and D Fradkin. Bayesian multinomial logistic regression for author identification. In *Maxent Conference*, pages 509–516, 2005.

B.K. Natarajan. Sparse approximate solutions to linear system. *SIAM J. Computing*, 24(2):227–234, 1995.

A. Ng. Feature selection, $\ell_1$ vs $\ell_2$ regularization, and rotational invariance. In *International Conference on Machine Learning (ICML)*, pages 78–85. ACM Press, New York, 2004.

A. Ng. On feature selection: Learning with exponentially many irrelevant features as training examples. In *International Conference on Machine Learning (ICML)*, pages 404–412, 1998.

S. Osher, Y. Mao, B. Dong, and W. Yin. Fast linearized bregman iteration for compressive sensing and sparse denoising. *Communications in Mathematical Sciences*, 8(1):93–111, 2010.

M.Y. Park and T. Hastie. $\ell_1$ regularized path algorithm for generalized linear models. *J. R. Statist. Soc. B*, 69:659–677, 2007.

L.C. Parra, C. Spence, and P. Sajda. Higher-order statistical properties arising from the non-stationarity of natural signals. In *Advances in Neural Information Processing Systems*, volume 13, pages 786–792, 2001.

L.C. Parra, C.D. Spence, A.D. Gerson, and P. Sajda. Recipes for the linear analysis of EEG. *Neuroimage*, 28(2):326–341, 2005.

S. Perkins and J. Theiler. Online feature selection using grafting. In *Proceedings of the Twenty-First International Conference on Machine Learning (ICML)*, pages 592–599. ACM Press, 2003.

M.G. Philiastides and P. Sajda. Temporal characterization of the neural correlates of perceptual decision making in the human brain. *Cereb Cortex*, 16(4):509–518, 2006.

B. Polyak. *Introduction to Optimization*. Optimization Software, 1987.

V. Roth. The generalized lasso. *IEEE Tran. Neural Networks*, 15(1):16–28, 2004.

L. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D*, 60(1-4):259–268, 1992.

M. Schmidt, G. Fung, and R. Rosales. Fast optimization methods for l1 regularization: a comparative study and two new approaches. In *European Conference on Machine Learning (ECML)*, pages 286–297, 2007.

J. Shi and S. Osher. A nonlinear inverse scale space method for a convex multiplicative noise model. *SIAM J. Imaging Sciences*, 1(3):294–321, 2008.

J. Shi, J. Wielaard, R.T. Smith, and P. Sajda. Perceptual decision making investigated via sparse decoding of a spiking neuron model of V1. In *4th International IEEE/EMBS Conference on Neural Engineering*, pages 558–561, 2009.

N.Z. Shor. *Minimization Methods for Non-differentiable functions*. Springer Series in Computational Mathematics. Springer, 1985.

H.L. Taylor, S.C. Banks, and J.F. McCoy. Deconvolution with the $\ell_1$ norm. *Geophysics*, 44(1): 39–52, 1979.

R. Tibshirani. Regression shrinkage and selection via the lasso. *J. Roy. Stat. Soc. B*, 58(1):267–288, 1996.

Y. Tsuruoka, J. McNaught, J. Tsujii, and S. Ananiadou. Learning string similarity measures for gene/protein name dictionary look-up using logistic regression. *Bioinformatics*, 23(20):2768–74, 2007.

L. Vandenberghe, S. Boyd, and A. El Gamal. Optimal wire and transistor sizing for circuits with non-tree topology. In *IEEE/ACM International Conference on Computer Aided Design*, pages 252–259, 1997.

L. Vandenberghe, S. Boyd, and A. El Gamal. Optimizing dominant time constant in RC circuits. *IEEE Trans. Computer-Aided Design*, 17(2):110–125, 1998.

V. Vapnik. *Estimation of Dependences Based on Empirical Data*. Springer-Verlag, 1982.

V. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, 1988.

Z. Wen, W. Yin, D. Goldfarb, and Y. Zhang. A fast algorithm for sparse reconstruction based on shrinkage, subspace optimization and continuation. Technical report, Rice University CAAM TR09-01, 2009.

W. Yin, S. Osher, J. Darbon, and D. Goldfarb. Bregman iterative algorithm for $\ell_1$-minimization with applications to compressed sensing. *SIAM J. Imaging Science*, 1(1):143–168, 2008.

P. Zhao and B. Yu. On model selection consistency of lasso. *J. Machine Learning Research*, 7: 2541–2567, 2007.

J. Zhu, S. Rosset, T. Hastie, and R. Tibshirani. 1-norm support vector machines. In *Advances in Neural Information Processing Systems*, volume 16, pages 49–56. MIT Press, 2004.

H. Zou, T. Hastie, and R. Tibshirani. Sparse principle component analysis. *J. Computational and Graphical Statistics*, 15(2):262–286, 2006.