

Direct Sparsity Optimization Based Feature Selection for Multi-Class Classification

Hanyang Peng¹, Yong Fan²

¹National Laboratory of Pattern Recognition, Institute of Automaiton, Chinese Academy of Sciences, 100190, Beijing, P.R. China

²Department of Radiology, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, 19104, USA

hanyang.peng@nlpr.ia.ac.cn, yong.fan@ieee.org

Abstract

A novel sparsity optimization method is proposed to select features for multi-class classification problems by directly optimizing a $\ell_{2,p}$ -norm ($0 < p \leq 1$) based sparsity function subject to data-fitting inequality constraints to obtain large between-class margins. The direct sparse optimization method circumvents the empirical tuning of regularization parameters in existing feature selection methods that adopt the sparsity model as a regularization term. To solve the direct sparsity optimization problem that is non-smooth and non-convex when $0 < p < 1$, we propose an efficient iterative algorithm with proved convergence by converting it to a convex and smooth optimization problem at every iteration step. The proposed algorithm has been evaluated based on publicly available datasets. The experiments have demonstrated that our algorithm could achieve feature selection performance competitive to state-of-the-art algorithms.

1 Introduction

Feature selection has been an important component in machine learning models (Guyon & Elisseeff, 2003). Feature selection approaches in general can be divided into three groups: filter methods (Kira & Rendell, 1992; Lewis, 1992; Peng et al., 2005), wrapper methods (Guyon et al., 2002), and embedded methods (Cawley et al., 2006; Wang et al., 2008; Xiang et al., 2012). The filter methods use proxy measures that are independent on the machine learning models to rank features according to their relevance to the learning problem. The wrapper methods search subsets of features to optimize a given learning model's performance, and typically have higher computational cost than the filter methods. The embedded methods integrate the feature selection task into the model learning, and are able to achieve good performance with a moderate computational cost. Particularly, sparse linear model based feature selection methods have achieved promising performance. The sparse linear model based

methods typically adopt ℓ_1 -norm regularization, and many variants have been proposed with different sparsity regularization terms, such as Lasso (Tibshirani, 1996) and sparse support vector machine (SVM) (Bradley & Mangasarian, 1998; Mangasarian, 2006).

In multi-task learning, various $\ell_{2,1}$ -norm (Liu et al., 2009; Nie et al., 2010; Obozinski et al., 2006) and $\ell_{\infty,1}$ -norm (Liu et al., 2009) based regularization models have been investigated for selecting features with joint sparsity across different tasks. Moreover, group Lasso based methods (Kong & Ding, 2013; Kong et al., 2014) have also been proposed in recent years. In fact, the multi-task feature selection algorithms share similarities with group lasso based methods. Since non-convex ℓ_p -norm or $\ell_{r,p}$ -norm ($0 < p < 1$) based regularization models can yield sparser solutions than ℓ_1 -norm or $\ell_{r,1}$ -norm based models, they have gained increasing attention in recent studies (Chartrand & Staneva, 2008; Liu et al., 2007; Zhang et al., 2014).

The sparse linear model based feature selection algorithms typically take a trade-off between a data-fitting loss function term and a sparsity term, therefore there inevitably exists residual in the loss function. However, little is known about such a residual's impact on the feature selection. Most of these algorithms are developed for classification problems. However, their data-fitting loss functions adopted are typically based on the regression of class labels, rather than surrogates of 0-1 loss that are more appropriate for classification problems (Bartlett et al., 2004). To overcome the limitations of existing sparse model based feature selection methods, we present a novel feature selection method via directly optimizing a linear model's sparsity with $\ell_{2,p}$ -norm ($0 < p \leq 1$), subject to data-fitting inequality constraints, instead of adopting the sparsity as a regularization term. Our constrained optimization formulation circumvents the difficulty of regularization parameter setting, and the inequality constraint based data-fitting loss function enables large between-class margins, similar to SVMs, but capable of handling multi-class problems. We propose an efficient algorithm to solve the optimization problem associated with the direct sparsity optimization, which is non-convex and non-smooth

when $0 < p < 1$, by transforming it to a Frobenius-norm induced problem at each iteration step, which has been proved to converge to Karush–Kuhn–Tucker (KKT) points.

The proposed algorithm has been evaluated based on 6 publicly available datasets, and extensive comparison experiments have demonstrated that our algorithm could achieve feature selection performance competitive to state-of-the-art algorithms, including ReliefF (Kira & Rendell, 1992), Min-Redundancy Max-Relevance (mRMR) (Peng, et al., 2005), ℓ_1 -norm based ℓ_1 -SVM (Bradley & Mangasarian, 1998; Mangasarian, 2006), and $\ell_{2,1}$ -norm based Robust Feature Selection (RFS) (Nie, et al., 2010; Xiang, et al., 2012).

2 Direct Sparsity Optimization Based Feature Selection (DSO-FS)

Throughout this paper, matrices are written in bold uppercase, vectors are written in bold lowercase, and all the scalars are denoted by normal letters. \mathbf{I} denotes an identity matrix and $\mathbf{1}$ denotes a vector or matrix with all the elements equal to 1. Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, the $\ell_{r,p}$ -norm ($r > 0, p > 0$)¹ of \mathbf{A} is defined as:

$$\|\mathbf{A}\|_{r,p} = \left(\sum_{i=1}^m \left(\sum_{j=1}^n |A_{i,j}|^r \right)^{\frac{p}{r}} \right)^{\frac{1}{p}} = \left(\sum_{i=1}^m (\|\mathbf{a}_i\|_r)^p \right)^{\frac{1}{p}}, \quad (1)$$

where $\|\mathbf{a}_i\|_r$ denotes ℓ_r -norm of the i -th row vector of \mathbf{A} .

Given m training samples $\{\mathbf{x}^i, y^i\}_{i=1}^m$ where $\mathbf{x}^i \in \mathbb{R}^n$ is a data point and y^i is its associated class label in c ($c \geq 2$) classes, the multiclass classification problem can be modeled as a linear learning problem. For simplicity, the bias of the standard linear regression is absorbed into \mathbf{W} as an additional dimension with all elements equal to 1.

$$\mathbf{XW} = \mathbf{Y}, \quad (2)$$

where $\mathbf{X} = [\mathbf{x}^1; \mathbf{x}^2; \dots; \mathbf{x}^m; \mathbf{1}]$, $\mathbf{W} \in \mathbb{R}^{(n+1) \times c}$ is the weight matrix to be learned, and $\mathbf{Y} = [\mathbf{f}^1; \mathbf{f}^2; \dots; \mathbf{f}^i; \dots; \mathbf{f}^m] \in \mathbb{R}^{m \times c}$ is a class label matrix with labels rearranged using a one-versus-rest model, i.e., $\mathbf{f}^i = [-1, \dots, 1, \dots, -1] \in \mathbb{R}^c$ (the j -th element is 1 and others

For selecting a subset of features, \mathbf{W} should have sparse columns and share joint sparsity along its row direction since each row of \mathbf{W} corresponds to the same feature. Therefore, we model the feature selection problem as

$$\min_{\mathbf{W}} \|\mathbf{W}\|_{2,0}, \quad s. t., \mathbf{XW} = \mathbf{Y}. \quad (3)$$

where $\|\mathbf{W}\|_{2,0}$ is the number of non-zero rows in \mathbf{W} , of which not all the elements are zero.

It is NP-hard to solve the optimization problem of Eqn. (3). Therefore $\ell_{2,p}$ -norm ($0 < p \leq 1$) can be adopted instead, resulting in a relaxed sparsity optimization problem:

$$\min_{\mathbf{W}} \|\mathbf{W}\|_{2,p}, \quad s. t., \mathbf{XW} = \mathbf{Y}, \quad (4)$$

where $0 < p \leq 1$.

Many studies (Liu, et al., 2009; Nie, et al., 2010; Obozinski, et al., 2006) assumed that the $\ell_{2,1}$ -norm based problems is equivalent to $\ell_{2,0}$ -norm based problems under certain conditions (Candes & Tao, 2005). In practice, $\ell_{2,p}$ -norm ($0 < p < 1$) can lead to sparser solutions in most cases although it is non-convex (Chartrand & Staneva, 2008; Liu, et al., 2007; Zhang, et al., 2014).

It is desired that the classification model's margins between classes are as large as possible for obtaining improved generalization performance. Accordingly, the equality constraint in Eqn. (4) is relaxed to be inequality constraints, i.e.,

$$\min_{\mathbf{W}} \|\mathbf{W}\|_{2,p}, \quad s. t., \mathbf{Y} \odot \mathbf{XW} \geq \mathbf{1}, \quad (5)$$

where $\mathbf{1} \in \mathbb{R}^{m \times c}$, \odot is a Hadamard product operator for element-wise multiplication between two matrices of the same dimensions, and \geq denotes that the elements of the left matrix are greater than or equal to their corresponding ones of the right matrix.

The optimization problem of Eqn. (5) can be reformulated by introducing a slack variable $\mathbf{E} \in \mathbb{R}^{m \times c}$ whose elements have the same positive or negative sign as the corresponding elements of \mathbf{Y} , as

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{E}} \quad & \|\mathbf{W}\|_{2,p} \\ s. t. \quad & \mathbf{XW} = \mathbf{Y} + \mathbf{E} \\ & \mathbf{Y} \odot \mathbf{E} \geq \mathbf{0}, \end{aligned} \quad (6)$$

where $\mathbf{E} \in \mathbb{R}^{m \times c}$.

To solve the optimization problem of Eqn. (6), we first solve the linear equation $\mathbf{XW} = \mathbf{Y} + \mathbf{E}$ to obtain the solution space of \mathbf{W} , and then directly search the solution space to find a solution with the minimum of $\|\mathbf{W}\|_{2,p}$. Note that if $\mathbf{XW} = \mathbf{Y} + \mathbf{E}$ is inconsistent, especially when the number of data samples m is greater than the number of features n , least-square solution space of the equation can be used as a substitute. Actually, we only need to solve $\mathbf{XW} = \mathbf{X}\mathbf{X}^+(\mathbf{Y} + \mathbf{E})$, where \mathbf{X}^+ is pseudo-inverse. This equation is compatible when \mathbf{X} is row full rank, since $\mathbf{X}\mathbf{X}^+ = \mathbf{I}$ on this occasion. Gaussian Elimination is a simple and efficient way to obtain the solution space of \mathbf{W} . Without loss the generality, we assume that the first m_0 column vectors of \mathbf{X} are linearly independent, i.e.

$$\begin{aligned} [\mathbf{X} : \mathbf{X}\mathbf{X}^+(\mathbf{Y} + \mathbf{E})] &= [\mathbf{X}_1 \ \mathbf{X}_2 : \mathbf{X}\mathbf{X}^+(\mathbf{Y} + \mathbf{E})] \\ \xrightarrow{\text{left-multiply } \mathbf{D}} & \begin{bmatrix} \mathbf{I} & \mathbf{M} & \mathbf{N} + \mathbf{LE} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}, \end{aligned} \quad (7)$$

where the rank of \mathbf{X} is m_0 , $\mathbf{X}_1 \in \mathbb{R}^{m \times m_0}$, $\mathbf{X}_2 \in \mathbb{R}^{m \times (n - m_0)}$, $\mathbf{I} \in \mathbb{R}^{m_0 \times m_0}$, $\mathbf{M} \in \mathbb{R}^{m_0 \times (n - m_0)}$, $\mathbf{N} \in \mathbb{R}^{m_0 \times c}$, and $\mathbf{D} \in \mathbb{R}^{m_0 \times m_0}$ is product matrix of a series of elemental matrices, $\mathbf{L} = \mathbf{D}\mathbf{X}\mathbf{X}^+ \in \mathbb{R}^{m_0 \times m_0}$, and $n_0 = n + 1 - m_0$. Thus, the solution space of \mathbf{W} is

¹ If $0 < r < 1$ or $0 < p < 1$, $\ell_{r,p}$ -norm does not satisfy triangle inequality. However, this does not affect the proposed algorithm.

$$\mathbf{W} = \mathbf{P}\mathbf{U} + \mathbf{Q} + \mathbf{F} = \begin{bmatrix} \mathbf{M} \\ \mathbf{I} \end{bmatrix} \mathbf{U} + \begin{bmatrix} \mathbf{N} \\ \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{L}\mathbf{E} \\ \mathbf{0} \end{bmatrix} \quad (8)$$

where $\mathbf{P} \in \mathbb{R}^{n \times n_0}$, $\mathbf{Q} \in \mathbb{R}^{n \times c}$, $\mathbf{F} \in \mathbb{R}^{n \times c}$, $\mathbf{I} \in \mathbb{R}^{n_0 \times n_0}$, and $\mathbf{U} \in \mathbb{R}^{n_0 \times c}$ is an arbitrary matrix.

Finally, the optimization problem of Eqn. (5) can be reformulated as

$$\min_{\mathbf{U}, \mathbf{E}} \left\| \mathbf{P}\mathbf{U} + \mathbf{Q} + \begin{bmatrix} \mathbf{L}\mathbf{E} \\ \mathbf{0} \end{bmatrix} \right\|_{2,p}, \text{ s.t.}, \mathbf{Y} \odot \mathbf{E} \succcurlyeq \mathbf{0} \quad (9)$$

where $0 < p \leq 1$.

3 An Iterative Algorithm for DSO-FS

We propose an iterative algorithm to solve the optimization problem of Eqn. (9) due to that no analytical solution is available. At each iteration step, we alternately optimize variables \mathbf{U} and \mathbf{E} . An objective function with $\ell_{2,p}$ -norm ($0 < p \leq 1$) is non-smooth and non-convex when $0 < p < 1$. To efficiently solve this problem, $\ell_{2,p}$ -norm is reformulated by Frobenius-norm (\mathcal{F} -norm) that is smooth and convex, as

$$\|\mathbf{A}\|_{2,p}^p = \|\boldsymbol{\Sigma}\mathbf{A}\|_F^2, \quad (10)$$

where $\boldsymbol{\Sigma}$ is a diagonal matrix with the i -th diagonal element $\Sigma_{ii} = 1/\|\mathbf{a}_i\|_2^{1-p/2}$, and $\|\mathbf{a}_i\|_2$ is defined in Eqn.(1)². Note that $\mathbf{Z}^* = \operatorname{argmin}_{\mathbf{Z}} \|\mathbf{A}(\mathbf{Z})\|_{2,p}^p = \operatorname{argmin}_{\mathbf{Z}} \|\mathbf{A}(\mathbf{Z})\|_{2,p}$, where $\mathbf{A}(\mathbf{Z})$ is a function of \mathbf{Z} . The similar strategies has also been adopted in FOCal Underdetermined System Solver (FOCUSS) (Cotter et al., 2005; Gorodnitsky & Rao, 1997), and various optimization techniques for $\ell_{2,1}$ -norm regularization based feature selection algorithms (Hou et al., 2011; Krishnapuram et al., 2005; Merchante et al., 2012; Nie, et al., 2010; Xiang, et al., 2012; Yi Yang et al., 2011) can be considered as its special case when $p = 1$. Our solution is summarized in Algorithm 1.

Algorithm 1. Feature Selection via Direct Sparsity Optimization (DSO-FS)

Input: data points $\{\mathbf{x}^i\}_{i=1}^m$ ($\mathbf{x}^i \in \mathbb{R}^n$) and their corresponding label $\{y^i\}_{i=1}^m$; norm power p ; number of features d to be selected.

Construct \mathbf{X} and \mathbf{Y} according to Eqn. (2), and \mathbf{L} , \mathbf{M} , \mathbf{N} , \mathbf{P} and \mathbf{Q} according to Eqn. (7) and Eqn. (8).

Set $k = 0$ and initialize \mathbf{E}^0 and $\boldsymbol{\Sigma}^0 \in \mathbb{R}^{n \times n}$ with a zero matrix and an identity matrix respectively.

repeat

$$\mathbf{G} = \mathbf{Q} + \begin{bmatrix} \mathbf{L}\mathbf{E}^k \\ \mathbf{0} \end{bmatrix}$$

$$\mathbf{U}^{k+1} = \operatorname{argmin}_{\mathbf{U}} \|\boldsymbol{\Sigma}^k(\mathbf{P}\mathbf{U} + \mathbf{G})\|_F^2$$

$$\mathbf{V}^k = -\mathbf{M}\mathbf{U}^{k+1} + \mathbf{N} + \mathbf{L}\mathbf{E}^k.$$

Update diagonal matrix \mathbf{A}^k , where the i -th di-

² $\boldsymbol{\Sigma} = \boldsymbol{\Pi}^{-1}$, where $\boldsymbol{\Pi}$ is a diagonal matrix with the i -th diagonal element $\Pi_{ii} = \|\mathbf{a}_i\|_2^{1-p/2}$. If any row vector $\mathbf{a}_i = \mathbf{0}$, we define $\boldsymbol{\Sigma} = \boldsymbol{\Pi}^+$.

agonal element is $\frac{1}{\|\mathbf{v}_i^k\|_2^{1-p/2}}$

$$\mathbf{H} = -\mathbf{M}\mathbf{U}^{k+1} + \mathbf{N}$$

$$\mathbf{E}^{k+1} = \operatorname{argmin}_{\mathbf{E}} \|\mathbf{A}^k(\mathbf{L}\mathbf{E} + \mathbf{H})\|_F^2, \text{ s.t. } \mathbf{Y} \odot \mathbf{E} \succcurlyeq \mathbf{0}$$

$$\text{Update } k = k + 1, \mathbf{W}^k = \mathbf{P}\mathbf{U}^k + \mathbf{Q} + \begin{bmatrix} \mathbf{L}\mathbf{E}^k \\ \mathbf{0} \end{bmatrix}$$

Update $\boldsymbol{\Sigma}^k$, where the i -th diagonal element is $\frac{1}{\|\mathbf{w}_i^k\|_2^{1-p/2}}$

until convergence

Sort all features according to $\|\mathbf{w}_i\|_2$ and select the largest d features.

Note that at each iteration step, for solving $\mathbf{U}^{k+1} = \operatorname{argmin}_{\mathbf{U}} \|\boldsymbol{\Sigma}^k(\mathbf{P}\mathbf{U} + \mathbf{G})\|_F^2$, an analytical solution is available, i.e.,

$$\mathbf{U}^{k+1} = -(\mathbf{P}^t \boldsymbol{\Sigma} \mathbf{P})^{-1} \mathbf{P}^t \mathbf{S} \mathbf{G}. \quad (11)$$

where $\mathbf{S} = (\boldsymbol{\Sigma}^k)^2$.

According to Eqn. (8), Eqn. (11) can be reformulated as

$$\mathbf{U}^{k+1} = (\mathbf{M}^t \mathbf{S}_1 \mathbf{M} + \mathbf{S}_2)^{-1} \mathbf{M}^t \mathbf{S}_1 \mathbf{K}, \quad (12)$$

where $\mathbf{S}_1 \in \mathbb{R}^{m_0 \times m_0}$ is a diagonal matrix, of which the diagonal elements are the first m_0 elements of \mathbf{S} , $\mathbf{S}_2 \in \mathbb{R}^{n_0 \times n_0}$ is a diagonal matrix, of which the diagonal elements are the last n_0 elements of \mathbf{S} , and $\mathbf{K} = \mathbf{N} + \mathbf{L}\mathbf{E}^k$ ($\in \mathbb{R}^{m_0 \times c}$).

If $m_0 < n_0$, the optimization problem's computational cost can be further reduced by a simple matrix operation

$$\mathbf{U}^{k+1} = \mathbf{T}(\mathbf{M}\mathbf{T} + \mathbf{I})^{-1} \mathbf{K}, \quad (13)$$

where $\mathbf{T} = \mathbf{S}_2^{-1} \mathbf{M}^t \mathbf{S}_1$ and $\mathbf{I} \in \mathbb{R}^{m_0 \times m_0}$.

For feature selection problems with the number of samples less than the number of features, i.e., $m_0 \leq m \ll n_0$, updating Eqn. (13) has reduced computational cost than Eqn. (11). It is worth noting that it is not necessary to directly calculate $\mathbf{C} = (\mathbf{M}\mathbf{T} + \mathbf{I})^{-1} \mathbf{K}$, which is computationally more expensive than solving the following linear equation

$$(\mathbf{M}\mathbf{T} + \mathbf{I})\mathbf{C} = \mathbf{K}. \quad (14)$$

Although no analytical solution is available for $\mathbf{E}^{k+1} = \operatorname{argmin}_{\mathbf{E}} \|\mathbf{A}^k(\mathbf{L}\mathbf{E} + \mathbf{H})\|_F^2$ (s.t. $\mathbf{Y} \odot \mathbf{E} \succcurlyeq \mathbf{0}$), it is a smooth and convex optimization problem. Therefore, it can be efficiently solved by existing tools, such as CVX (CVX Research, 2011).

4 Convergence proof

Algorithm 1 makes $\|\mathbf{W}\|_{2,p}$ to monotonically decrease at every iteration step and finally converges.

Lemma 1. Given any two vectors \mathbf{a} and \mathbf{b} , we have

$$(1 - \theta)\|\mathbf{a}\|_2^2 + \theta\|\mathbf{b}\|_2^2 \geq \|\mathbf{a}\|_2^{2-2\theta}\|\mathbf{b}\|_2^{2\theta}, \quad (15)$$

where $0 < \theta < 1$ and the equality holds if and only if $\mathbf{a} = \mathbf{b}$.

Proof. Since $\ln(x^2)$ is concave, we have

$$\ln((1-\theta)x_1^2 + \theta x_2^2) \geq (1-\theta)\ln(x_1^2) + \theta\ln(x_2^2), \quad (16)$$

where $0 < \theta < 1$. The equality holds if and only if $x_1 = x_2$, indicating that

$$(1-\theta)x_1^2 + \theta x_2^2 \geq x_1^{2-2\theta} x_2^{2\theta}. \quad (17)$$

Then, we have

$$(1-\theta)\|\mathbf{a}\|_2^2 + \theta\|\mathbf{b}\|_2^2 \geq \|\mathbf{a}\|_2^{2-2\theta}\|\mathbf{b}\|_2^{2\theta}, \quad (18)$$

where the equality holds if and only if $\mathbf{a} = \mathbf{b}$. \square

Lemma 2. Given an optimization problem:

$$\min_{\mathbf{Z}} \|\mathbf{S}\Phi(\mathbf{Z})\|_{\mathcal{F}}^2, \quad \text{s. t. } \mathbf{Z} \in \mathcal{F},$$

where $\Phi(\mathbf{Z})$ is a function of \mathbf{Z} , \mathcal{F} is the feasible region, and \mathbf{S} is a diagonal matrix whose i -th diagonal element is $1/\|\Phi(\mathbf{Z}_0)_i\|_2^{1-p/2}$ (\mathbf{Z}_0 could be any object in \mathcal{F} , $\Phi(\mathbf{Z}_0)_i$ is the i -th row vector of $\Phi(\mathbf{Z}_0)$ and $0 < p \leq 2$), we have

$$\|\Phi(\mathbf{Z}^*)\|_{2,p} \leq \|\Phi(\mathbf{Z}_0)\|_{2,p}, \quad (19)$$

where \mathbf{Z}^* is the optimal solution of Eqn. (19) and the equality holds if and only if $\Phi(\mathbf{Z}^*) = \Phi(\mathbf{Z}_0)$

Proof. Since \mathbf{Z}^* is the optimal solution, we have

$$\|\mathbf{S}\Phi(\mathbf{Z}^*)\|_{\mathcal{F}}^2 \leq \|\mathbf{S}\Phi(\mathbf{Z}_0)\|_{\mathcal{F}}^2. \quad (20)$$

Then

$$\sum_i \frac{\|\Phi(\mathbf{Z}^*)_i\|_2^2}{\|\Phi(\mathbf{Z}_0)_i\|_2^{2-p}} \leq \sum_i \|\Phi(\mathbf{Z}_0)_i\|_2^p. \quad (21)$$

where $\Phi(\mathbf{Z}_0)_i$ and $\Phi(\mathbf{Z}^*)_i$ are the i -th row vector of $\Phi(\mathbf{Z}_0)$ and $\Phi(\mathbf{Z}^*)$, respectively.

According to Lemma 1, we have

$$\begin{aligned} & \left(1 - \frac{p}{2}\right) \|\Phi(\mathbf{Z}_0)_i\|_2^2 + \frac{p}{2} \|\Phi(\mathbf{Z}^*)_i\|_2^2 \\ & \geq \|\Phi(\mathbf{Z}_0)_i\|_2^{2-p} \|\Phi(\mathbf{Z}^*)_i\|_2^p. \end{aligned} \quad (22)$$

Then dividing the both sides by $\|\Phi(\mathbf{Z}_0)_i\|_2^{2-p}$, we have

$$\begin{aligned} \|\Phi(\mathbf{Z}^*)_i\|_2^p & \leq \left(1 - \frac{p}{2}\right) \|\Phi(\mathbf{Z}_0)_i\|_2^p \\ & \quad + \frac{p}{2} \frac{\|\Phi(\mathbf{Z}^*)_i\|_2^2}{\|\Phi(\mathbf{Z}_0)_i\|_2^{2-p}} \end{aligned} \quad (23)$$

It indicates that

$$\begin{aligned} \sum_i \|\Phi(\mathbf{Z}^*)_i\|_2^p & \leq \left(1 - \frac{p}{2}\right) \sum_i \|\Phi(\mathbf{Z}_0)_i\|_2^p \\ & \quad + \frac{p}{2} \sum_i \frac{\|\Phi(\mathbf{Z}^*)_i\|_2^2}{\|\Phi(\mathbf{Z}_0)_i\|_2^{2-p}}. \end{aligned} \quad (24)$$

Combining Eqns. (21) and (24), we obtain

$$\sum_i \|\Phi(\mathbf{Z}^*)_i\|_2^p \leq \sum_i \|\Phi(\mathbf{Z}_0)_i\|_2^p. \quad (25)$$

Therefore,

$$\|\Phi(\mathbf{Z}^*)\|_{2,p} \leq \|\Phi(\mathbf{Z}_0)\|_{2,p}, \quad (26)$$

where the equality holds if and only if $\Phi(\mathbf{Z}^*) = \Phi(\mathbf{Z}_0)$. \square

Theorem 1. The sequence $\{\mathbf{W}^k\}$ produced via Algorithm 1 has the following properties: $\|\mathbf{W}^k\|_{2,p}$ is non-increasing at successive iteration steps and $\{\|\mathbf{W}^k\|_{2,p}\}$ converges to a limited value.

Proof. Supposing we have obtained the solution \mathbf{U}^k , \mathbf{E}^k , and the objective function \mathbf{W}^k at the $(k+1)$ -th iteration step, we solve the optimization problem $\min_{\mathbf{U}} \|\Sigma^k(\mathbf{P}\mathbf{U} + \mathbf{G})\|_{\mathcal{F}}^2$ to obtain \mathbf{U}^{k+1} by fixing \mathbf{E}^k . According to Lemma 2, we have

$$\begin{aligned} \|\mathbf{V}^k\|_{2,s} & = \left\| \mathbf{P}\mathbf{U}^{k+1} + \mathbf{Q} + \begin{bmatrix} \mathbf{L}\mathbf{E}^k \\ \mathbf{0} \end{bmatrix} \right\|_{2,p} \\ & \leq \|\mathbf{W}^k\|_{2,p}. \end{aligned} \quad (27)$$

Then we fix \mathbf{U}^{k+1} , and solve the optimization problem $\min_{\mathbf{E}} \|\Lambda^k(\mathbf{L}\mathbf{E} + \mathbf{H})\|_{\mathcal{F}}^2$ to obtain \mathbf{E}^{k+1} . According to Lemma 2, we have

$$\begin{aligned} \|\mathbf{W}^{k+1}\|_{2,p} & = \left\| \mathbf{P}\mathbf{U}^{k+1} + \mathbf{Q} + \begin{bmatrix} \mathbf{L}\mathbf{E}^{k+1} \\ \mathbf{0} \end{bmatrix} \right\|_{2,p} \\ & \leq \|\mathbf{V}^k\|_{2,p}. \end{aligned} \quad (28)$$

Combining (27) and (28), we obtain

$$\|\mathbf{W}^{k+1}\|_{2,p} \leq \|\mathbf{W}^k\|_{2,p}, \quad (29)$$

where the equality holds if and only if $\mathbf{W}^{k+1} = \mathbf{V}^k = \mathbf{W}^k$. Since the lower bound of $\|\mathbf{W}\|_{2,p}$ is limited, $\{\|\mathbf{W}^k\|_{2,p}\}$ will converge. \square

Theorem 2. If sequences $\{\mathbf{W}^k\}$ and $\{\mathbf{E}^k\}$ produced in Algorithm 1 have limit points, the limit points satisfy the Karush–Kuhn–Tucker (KKT) conditions of Eqn. (6). When $p \geq 1$, the limited points are globally optimal.

Due to space limitations, the proof is provided in supplementary document (Peng & Fan, 2016)..

It is worth noting that Lemma 2 works for sparsity regularization based feature selection algorithms, *i.e.*,

$$\min_{\mathbf{Z}} \sum_i \|\mathbf{A}_i\mathbf{Z} + \mathbf{B}_i\|_{2,p} + \lambda\|\mathbf{Z}\|_{2,p}, \quad (30)$$

where λ is the regularization coefficient.

Eqn. (30) can be reformulated as

$$\min_{\mathbf{Z}} \left[(\mathbf{A}_1\mathbf{Z} + \mathbf{B}_1)^t, \dots, (\mathbf{A}_i\mathbf{Z} + \mathbf{B}_i)^t, \dots, \lambda\mathbf{Z}^t \right]_{2,p}. \quad (31)$$

The Robust Feature Selection (RFS) (Nie, et al., 2010; Xiang, et al., 2012) is a special case of Eqn. (30) when $i = 1$ and $p = 1$.

5 Experiments

5.1 Experimental datasets and settings

The proposed algorithm has been evaluated based on 6 publicly available datasets. In particular, 2 datasets were obtained from UCI, including ISOLET and SEMEION. ISOLET is a speech recognition data set with 7797 samples from 26 classes, and each sample has 617 features; SEMEION contains 1593 handwritten images from ~80 people, stretched in a rectangular box of 16×16 . For these 2 datasets, the number of features is less than the number of their samples. Another 2 datasets were microarray data, including LUNG and CLL-SUB-111. In particular, LUNG consists of 203 samples, each having 3312 genes with standard deviations larger than 50 expression units (Bhattacharjee et al., 2001); CLL-SUB-111 consists of 111 samples with 11340 features (featureselection.asu.edu). Our algorithm has also been validated based on 2 image datasets, including UMIST and AR. In particular, UMIST includes face images with a resolution of 56×46 from 20 different people, and AR has 130 samples with 2400 features.

We compared our method with sparsity regularization based feature selection methods, including ℓ_1 -SVM (Bradley & Mangasarian, 1998; Mangasarian, 2006) and $\ell_{2,1}$ -norm based Robust Feature Selection (RFS) (Nie, et al., 2010; Xiang, et al., 2012), with an objective function $\mathcal{J}(\mathbf{W}) = \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|_{2,1} + \lambda\|\mathbf{W}\|_{2,1}$. We also compared our algorithm with well-known filter methods, including ReliefF (Kira & Rendell, 1992) and mRMR (Peng, et al., 2005).

The feature selection methods were evaluated based on their classification accuracy. Particularly, linear SVM (Chang & Lin, 2011) was adopted to build classifiers based on the selected features. The parameter C of linear SVM classifiers were tuned using a cross-validation strategy by searching a candidate set $[10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10^1, 10^2]$. The regularization parameter of ℓ_1 -SVM and RFS were tuned using the same cross-validation strategy by searching a candidate set $[10^{-3}, 10^{-2}, 10^{-1}, 1, 10^1, 10^2, 10^3]$.

In our experiments, we first normalized all the data to have

0 mean and unit standard deviation for each feature. 10 trials were performed on each dataset. In each trial, the samples of each dataset were randomly spitted into training and testing subsets with a ratio of 6:4. For tuning parameters, a 3-fold was used for datasets with less than 200 training samples, and an 8-fold cross-validation was used for other datasets.

5.2 Effect of parameter p

To investigate how the classification performance is affected by the parameter p of DSO-FS, we performed experiments based on datasets ISOLET, LUNG, and UMIST. We obtained solutions of \mathbf{W} with different settings of $p \in [0.1, 0.3, 0.5, 0.7, 0.9, 1]$ using Algorithm 1 based on the above 3 datasets, and then selected top ranked features according to ℓ_2 -norm $\|\mathbf{w}_i\|_2$ to build classifiers. Figure 1 shows the classification accuracy as a function of the number of the selected features and the parameter p . The results shown in Figure 1 indicated that p played an important role in the classification.

5.3 Comparisons with state-of-the-art methods

Since the classification performance of linear SVM classifiers built on the features selected by our method is hinged on the parameter p , we used a cross-validation strategy to select an optimal value from $[0.1, 0.3, 0.5, 0.7, 0.9, 1]$ for our method. Figure 2 shows the average classification performance of classifiers built on the features selected by different methods in 10 trials. In particular, the average classification accuracy is shown as a function of the number of features used in the classification model. Compared with other methods, the proposed method achieved higher classification accuracy on most datasets, indicating that our method had overall better performance than other algorithms.

Table 1 summarizes mean and standard deviation of the classification rates in 10 trails for classifiers built on the top 100 features selected by the algorithms under comparison. These results demonstrated that our algorithm had the overall best classification accuracy on all the 6 datasets.

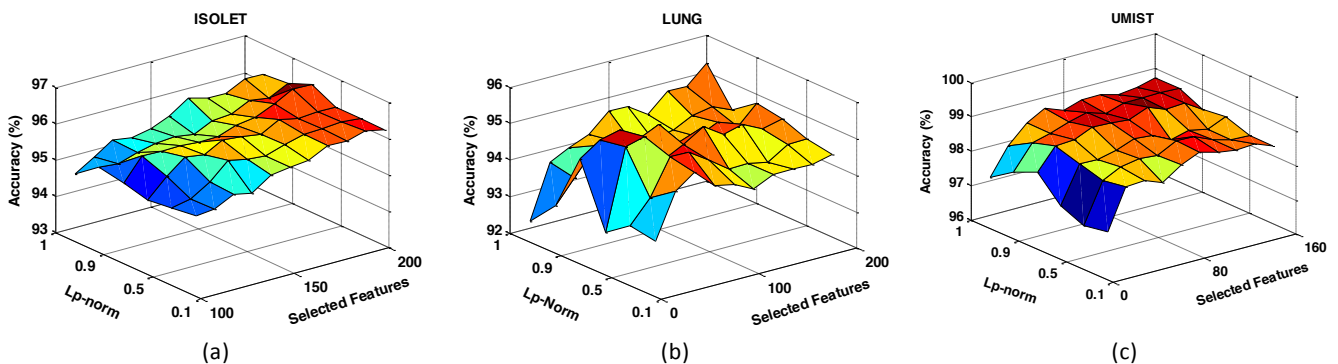


Figure 1: Classification accuracy with different numbers of features selected with different values of p . The results shown were obtained based on datasets: (a) ISOLET, (b) LUNG, and (c) UMIST.

Table1: Mean and standard deviation of the classification Accuracy (%) of Linear-SVM built on the selected top 100 features by different algorithms for datasets: ISOLET, SEMEION, LUNG, CLL-SUB-111, UMIST, and AR.

	mRMR	ReliefF	L1-SVM	RFS	DSO-FS
ISOLET	90.82±0.72	89.13±0.53	93.90±0.53	93.72±0.55	95.23±0.44
SEMEION	88.18±1.01	86.54±0.89	90.36±0.33	90.45±1.00	90.64±0.93
LUNG	94.88±1.20	94.88±1.42	95.12±1.22	95.37±1.42	95.37±1.42
CLL-SUB-111	76.47±11.0	72.94±9.18	79.41±4.15	75.88±6.80	81.18±6.60
UMIST	96.87±0.84	98.17±0.92	98.17±0.80	98.04±0.73	98.39±1.16
AR	87.31±5.02	85.58±6.15	87.12±9.26	88.85±6.76	89.62±3.74

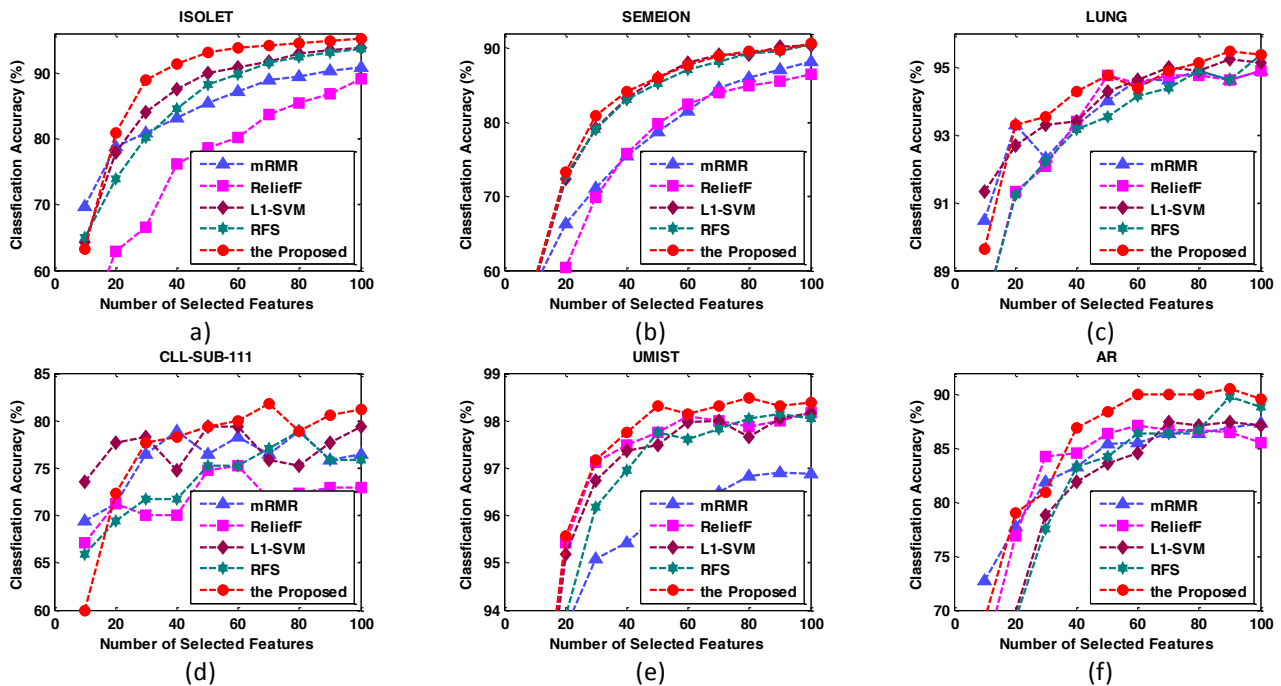


Figure 2: Average classification accuracy of 10 trials for classifiers built on the selected top 100 features by different algorithms. The results shown were obtained based on datasets: (a) ISOLET, (b) SEMEION, (c) LUNG, (d) CLL-SUB-111, (e) UMIST, and (f) AR.

6 Conclusions and Discussions

In this paper, a novel feature selection algorithm via direct sparsity optimization was proposed. Our method directly optimizes a large margin linear classification model's sparsity with $\ell_{2,p}$ -norm ($0 < p \leq 1$) subject to a data-fitting constraint, different from the sparse regularization based algorithms that typically take a trade-off between a data-fitting loss function term and a sparsity term. Since little is known about such a trade-off's impact on the feature selection, empirical parameter tuning is typically adopted for choosing the trade-off parameter. Our feature selection method adopts data-fitting inequality constraints to obtain large between-class margins, which could improve the generalization ability of the selected features, rather than data-fitting loss functions built upon the regression of class labels (Bartlett, et al., 2004). We also proposed an efficient algorithm to solve the non-convex

($0 < p < 1$) and non-smooth optimization problem associated with the feature selection problem. Extensive experiments based on 6 datasets have demonstrated that the proposed method could achieve better performance than state-of-the-art feature selection algorithms. Furthermore, our algorithm can be easily extended for solving other sparsity regularization algorithms. In particular, our algorithm could be used to solve ℓ_0 and $\ell_{2,0}$ based optimization problems subject to linear constraints by setting p close to 0. Furthermore, similar to SVM, we could also add soft-margin to our algorithm for better performance in the future work.

Acknowledgments

This work was supported in part by National Key Basic Research and Development Program of China (No. 2015CB856404), NSFC (No. 81271514, 61473296), and NIH grants MH070365 and AG014971.

References

- [Bartlett, P.L., et al., 2004] Large margin classifiers: convex loss, low noise, and convergence rates. *Advances in Neural Information Processing Systems* 16. 16: 1173-1180, 2004.
- [Bhattacharjee, A., et al., 2001] Classification of human lung carcinomas by mRNA expression profiling reveals distinct adenocarcinoma subclasses. *Proceedings of the National Academy of Sciences of the United States of America*. 98: 13790-13795, 2001.
- [Bradley, P.S. and Mangasarian, O.L., 1998] Feature selection via concave minimization and support vector machines. *Proceedings of the International Conference on Machine Learning*: 82-90, 1998.
- [Candes, E.J. and Tao, T., 2005] Decoding by linear programming. *Ieee Transactions on Information Theory*. 51: 4203-4215, 2005.
- [Cawley, G.C., et al., 2006] Sparse multinomial logistic regression via Bayesian l1 regularisation. *Advances in Neural Information Processing Systems*: 209-216, 2006.
- [Chang, C.C. and Lin, C.J., 2011] LIBSVM: A Library for Support Vector Machines. *Acm Transactions on Intelligent Systems and Technology*. 2, 2011.
- [Chartrand, R. and Staneva, V., 2008] Restricted isometry properties and nonconvex compressive sensing. *Inverse Problems*. 24, 2008.
- [Cotter, S.F., et al., 2005] Sparse solutions to linear inverse problems with multiple measurement vectors. *Ieee Transactions on Signal Processing*. 53: 2477-2488, 2005.
- Cvx Research, I. CVX: Matlab software for disciplined convex programming, version 2.0, <http://cvxr.com/cvx>, 2011.
- [Gorodnitsky, I.F. and Rao, B.D., 1997] Sparse signal reconstruction from limited data using FOCUSS: A re-weighted minimum norm algorithm. *Ieee Transactions on Signal Processing*. 45: 600-616, 1997.
- [Guyon, I. and Elisseeff, A., 2003] An introduction to variable and feature selection. *Journal of Machine Learning Research*. 3: 1157-1182, 2003.
- [Guyon, I., et al., 2002] Gene selection for cancer classification using support vector machines. *Machine Learning*. 46: 389-422, 2002.
- [Hou, C., et al., 2011] Feature Selection via Joint Embedding Learning and Sparse Regression. *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence*, 2011.
- [Kira, K. and Rendell, L.A., 1992] A Practical Approach to Feature-Selection. *Machine Learning*: 249-256, 1992.
- [Kong, D.G. and Ding, C., 2013] Efficient Algorithms for Selecting Features with Arbitrary Group Constraints via Group Lasso. 2013 *Ieee 13th International Conference on Data Mining (Icdm)*: 379-388, 2013.
- [Kong, D.G., et al., 2014] Exclusive Feature Learning on Arbitrary Structures via L1,2-norm. *Advances in Neural Information Processing Systems*, 2014.
- [Krishnapuram, B., et al., 2005] Sparse multinomial logistic regression: Fast algorithms and generalization bounds. *Ieee Transactions on Pattern Analysis and Machine Intelligence*. 27: 957-968, 2005.
- [Lewis, D.D., 1992] Feature-Selection and Feature-Extraction for Text Categorization. *Speech and Natural Language*: 212-217, 1992.
- [Liu, H., et al., 2009] Blockwise coordinate descent procedures for the multi-task lasso, with applications to neural semantic basis discovery. *International Conference on Machine Learning*, 2009.
- [Liu, J., et al., 2009] Multi-Task Feature Learning Via Efficient L2,1-Norm Minimization. *Uncertainty in Artificial Intelligence*, 2009.
- [Liu, Y.F., et al., 2007] Support vector machines with adaptive L-q penalty. *Computational Statistics & Data Analysis*. 51: 6380-6394, 2007.
- [Mangasarian, O.L., 2006] Exact l1-norm support vector machines via unconstrained convex differentiable minimization. *Journal of Machine Learning Research*: 1517-1530, 2006.
- [Merchant, L.F.S., et al., 2012] An Efficient Approach to Sparse Linear Discriminant Analysis. *International Conference on Machine Learning*, 2012.
- [Nie, F., et al., 2010] Efficient and Robust Feature Selection via Joint L2,1-Norms Minimization. *Advances in Neural Information Processing Systems*, 2010.
- [Obozinski, G., et al., 2006] Multi-task feature selection. Technical report, Department of Statistics, University of California, Berkeley, 2006.
- [Peng, H. and Fan, Y., 2016] Supplementary material of "Direct Sparsity Optimization Based Feature Selection for Multi-Class Classification", 2016.
- [Peng, H.C., et al., 2005] Feature selection based on mutual information: Criteria of max-dependency, max-relevance, and min-redundancy. *Ieee Transactions on Pattern Analysis and Machine Intelligence*. 27: 1226-1238, 2005.
- [Tibshirani, R., 1996] Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society Series B-Methodological*. 58: 267-288, 1996.
- [Wang, L., et al., 2008] Hybrid huberized support vector machines for microarray classification and gene selection. *Bioinformatics*. 24: 412-419, 2008.
- [Xiang, S.M., et al., 2012] Discriminative Least Squares Regression for Multiclass Classification and Feature Selection. *Ieee Transactions on Neural Networks and Learning Systems*. 23: 1738-1754, 2012.
- [Yi Yang, et al., 2011] L2,1-Norm Regularized Discriminative Feature Selection for Unsupervised Learning. *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence*, 2011.
- [Zhang, M., et al., 2014] Feature Selection at the Discrete Limit. *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence*, 2014.