# Direct Sparsity Optimization Based Feature Selection for Multi-Class Classification

Hanyang Peng[1]  and Yong Fan[2]

[1]National Laboratory of Pattern Recognition,
Institute of Automation, Chinese Academy of Sciences, China
[2]Department of Radiology, Perelman School of Medicine,
University of Pennsylvania,

July 13, 2016

## Outline

# Introduction

## Why Feature Selection?

- To remove redundant or noisy features

- To improve the generalized performance

- To reduce the computational burden

- To enhance the interpretability of intrinsic characteristics of data

# Introduction

## Fundamental Model for Feature Selection

Solving the following $l_0$-Minimization problem, subject to data fitting constraints, $Xw = y$, and then utilize the non-elements of solution to select useful features, i.e.,

$$\min_{w} \|w\|_0 \ , \ s.t., Xw = y \qquad (1)$$

## Basis Pursuit

By satisfying some assumptions(Restricted Isometry Property, RIP), the solution of Problem (1) can be obtained by solving (2), i.e.,

$$\min_{w} \|w\|_1 \ , \ s.t., \ Xw = y \qquad (2)$$

It is not so robust when the data $X$ or class label $y$ is corrupted by noise.

# Introduction

## Sparse SVM

Many Sparse SVM methods with discriminant margin are proposed to improve the robustness and enhance performance , such as $l_1$-SVM, i.e.,

$$\min_{\boldsymbol{w}} \|\boldsymbol{w}\|_1 \;, \; s.t., \boldsymbol{y} \odot \boldsymbol{X}\boldsymbol{w} \geqslant \boldsymbol{1} \qquad (3)$$

The optimization algorithm is special design for binary–class problem, hence the multi-class problem do not have compact form.

## Sparsity Regularization Based Methods

Many Sparsity Regularization Based Methods have been proposed with different sparsity regularization terms, such as Lasso

$$\min_{\boldsymbol{W}} \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|_2^2 + \lambda\|\boldsymbol{w}\|_1 \qquad (4)$$

A trade-off between a data-fitting loss function term and a sparsity term should be took, and it is sensitive to the parameter $\lambda$

# Introduction

## Sparsity Regularization Based Methods

In recent years, To learn sparse representations shared across multiple tasks or multiple classes, $l_{2,1}$ −norm based regularized method are proposed, and the class label is rearranged as one-versus-rest model, where $\boldsymbol{Y} = \left\{\boldsymbol{f}^i\right\}_{i=1}$ and $\boldsymbol{f}^i = [-1, \dots, 1, \dots, -1]$, such as Robust Feature selection (RFS),

$$\min_{\boldsymbol{W}} \|\boldsymbol{XW} - \boldsymbol{Y}\|_{2,1} + \lambda\|\boldsymbol{w}\|_{2,1} \qquad (5)$$

# Direct $L_{2,p}$-Minimization for feature selection

## The Proposed Original Model

$$\min_{W}\|W\|_{2,p}, \quad s.t., Y \odot XW \geqslant 1 \qquad (6)$$

Advantages

- $\ell_{2,p}$-norm $(0 < p < 1)$ can give rise to more sparse solutions

- No regularization term , do not need to make a compromise between residual of data-fitting and sparsity

- Enlarging discriminant margin between classes can boost generalization performance

It is difficult to solve directly.

# Direct $L_{2,p}$-Minimization for feature selection

## Equivalent Model

The optimization problem of (6) can be reformulated by introducing a slack variable $E$ whose elements have the same sign as the corresponding elements of $Y$, *i.e.,*

$$\min_{W,E}\|W\|_{2,p}, \quad s.t., XW = Y + E, Y \odot E \succcurlyeq 0 \tag{7}$$

## Direct Optimization

- Step 1: solve the linear equation $XW = Y + E$ to obtain the solution space of $W$ with variable $E$

- Step 2: directly search the solution space to find a solution to minimize $\|W\|_{2,p}$

# An Optimization Algorithm for the Model

## Solution Space of W

Gaussian Elimination

$$[X \vdots (Y + E)] = [X_1 \, X_2 \vdots (Y + E)] \xrightarrow{\text{left}-\text{multiply } L} \begin{bmatrix} I & M & \vdots & N + LE \\ 0 & 0 & \vdots & 0 \end{bmatrix} \quad (8)$$

The solution space of $W$ is

$$W = PU + Q + F = \begin{bmatrix} M \\ I \end{bmatrix} U + \begin{bmatrix} N \\ 0 \end{bmatrix} + \begin{bmatrix} LE \\ 0 \end{bmatrix} \quad (9)$$

The problem (10) can be reformulated as

$$\min_{U,E} \left\| PU + Q + \begin{bmatrix} LE \\ 0 \end{bmatrix} \right\|_{2,p}, s.t., Y \odot E \geqslant 0, \quad (10)$$

$\ell_{2,\text{p}}$-norm $(0 < p \leq 1)$ is non-smooth and non-convex when $0 < p < 1$

# An Optimization Algorithm for the Model

## Iterative Optimization Algorithm

- we alternately optimize variables $U$ and $E$ for optimization problem (14).

- Adopting Iteratively Reweighted Least Square (IRLS) straregy, $\ell_{2,\mathrm{p}}$-minimization problem can be reformulated as a least square minimization problem.

  At each iterative step, the objective function of the sub-problem can become convex and smooth.

# An Optimization Algorithm for the Model

## Optimizing Variable U

At $k$-th step

$$W^k = PU^k + Q + \begin{bmatrix} LE^k \\ 0 \end{bmatrix}$$

$$G^k = Q + \begin{bmatrix} LE^k \\ 0 \end{bmatrix}$$

$$U^{k+1} = \arg\min_U \left\| \Sigma^k (PU + G^k) \right\|_F^2$$, where the $i$-th diagonal element of $\Sigma^k$ is $1/\|w_i^k\|_2^{1-p/2}$

## Optimizing Variable E

At $k$-th step

$$V^k = -MU^{k+1} + N + LE^k$$

$$H = -MU^{k+1} + N$$

$$E^{k+1} = \arg\min_E \left\| \Lambda^k (LE + H) \right\|_F^2, \text{ s.t. } Y \odot E \geqslant 0,$$ where the $i$-th diagonal element of $\Lambda^k$ is $1/\|v_i^k\|_2^{1-p/2}$

# Proof of Convergence

## Lemma 1.

*Given any two vectors $\boldsymbol{a}$ and $\boldsymbol{b}$, we have*

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

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

## Lemma 2.

*Given an optimization problem:*

$$\min_{\boldsymbol{Z}}\;\|\,\boldsymbol{S}\,\boldsymbol{\Phi}(\boldsymbol{Z})\,\|_F^2\,,\, s.t.\;\boldsymbol{Z}\in\boldsymbol{\mathcal{F}}$$

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

$$\|\boldsymbol{\Phi}(\boldsymbol{Z}^*)\|_{2,p} \leq \|\boldsymbol{\Phi}(\boldsymbol{Z_0})\|_{2,p}$$

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

# Proof of Convergence

## Theorem 1.

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

## Theorem 2.

If sequences $\{W^k\}$ and $\{E^k\}$ produced in The Algorithm 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.
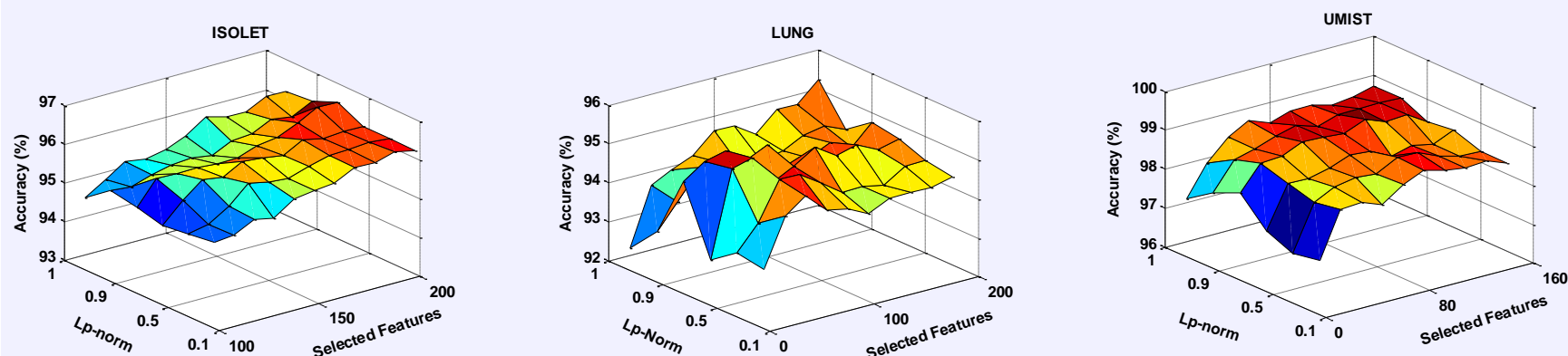
## Effect of parameter $p$



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

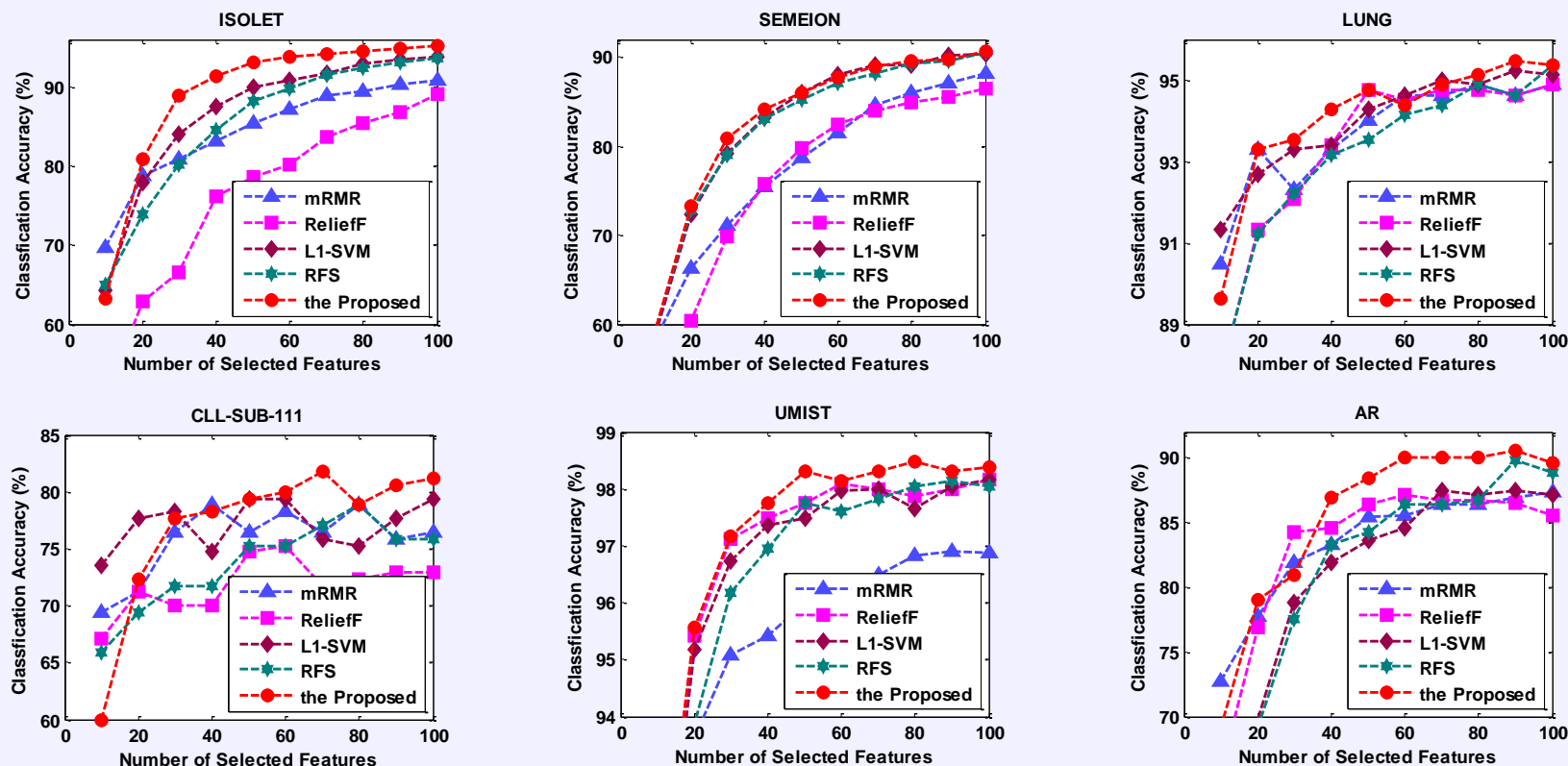# Experiments

## Effect of parameter $p$



Figure 2: Average classification accuracy of 10 trials for linear—SVM 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

# Conclusions and Discussions

## Summary

- Proposed Model: $L_{2,p}$-Minimization subject to data-fitting inequality constraints

- Outstanding Features
  - $L_{2,p}$-norm boosts more sparsity
  - No regularization term free tuning the parameter
  - Enlarging margin between classes improve the robustness to noise and generalization performance

- Optimization Approach
  - Adopting Gaussian Elimination, obtaining the solution space of $W$ with variable $E$
  - Utilizing IRLS strategy, at each iterative step, reformulating the non-convex and non-smooth problem to a least square minimization problem

Thanks for your attention